Actor - Critic    Ch 13 - Sutton & Barto

policy $p(a|s;\theta)$

on-policy samples

$\sum_{t=0}^{T} r(s_t, a_t)$

payoff $J(\theta) = E_{\tau \sim p(\tau;\theta)}[r(\tau)]$

$(s_0, a_0, \dots, s_T, a_T, s_{T+1})$

$= E_{\tau \sim \gamma}[p(\tau;\theta) r(\tau)]$ ← importance sampling

$E_{x \sim p(x)}[f(x)] = E_{x \sim q(x)}\left[\frac{p(x)}{q(x)} f(x)\right]$

$\nabla_\theta J(\theta) = \nabla_\theta E_{\tau \sim \gamma}[p(\tau;\theta) r(\tau)]$

$= E_{\tau \sim \gamma}[\nabla_\theta p(\tau;\theta) r(\tau)]$ ← linearity of $E$, $\nabla_\theta$

$= E_{\tau \sim \gamma}[p(\tau;\theta) \nabla_\theta \log p(\tau;\theta) r(\tau)]$ ← chain rule

$\nabla \log p = \frac{\nabla p}{p}$

$= E_{\tau \sim p(\tau;\theta)}[\nabla_\theta \log p(\tau;\theta) r(\tau)]$

$$\nabla_\theta \log p(\tau;\theta) = \nabla_\theta \log\left(p(s_0) \prod_{t=0}^{T} p(a_t \mid s_t;\theta)\, p(s_{t+1} \mid s_t, a_t)\right)$$

$$= \nabla_\theta \left(\log p(s_0) + \sum_{t=0}^{T} \log p(a_t \mid s_t;\theta)\right.$$

$$\left. + \sum_{t=0}^{T} \log p(s_{t+1} \mid s_t, a_t)\right)$$

$$= \sum_{t=0}^{T} \nabla_\theta \log p(a_t \mid s_t;\theta)$$

$$\nabla_\theta J(\theta) = E_{\tau \sim p(\tau;\theta)}\left[\left(\sum_{t=0}^{T} \nabla_\theta \log p(a_t \mid s_t;\theta)\right)\left(\sum_{t=0}^{T} r(s_t, a_t)\right)\right]$$

Causality $\longrightarrow \quad = E_{\tau \sim p(\tau;\theta)}\left[\sum_{t=0}^{T}\left(\nabla_\theta \log p(a_t \mid s_t;\theta) \sum_{k=t}^{T} r(s_k, a_k)\right)\right]$

$$\nabla_\theta J(\theta) = E_{\tau \sim p(\tau;\theta)} \left[ \sum_{t=0}^{T} \left( \nabla_\theta \log p(a_t | s_t; \theta) \sum_{k=t}^{T} r(s_k, a_k) \right) \right]$$

$$= \sum_{t=0}^{T} E_{(s_0, a_0, \cdots, s_T, a_T, s_{T+1})} \left[ \nabla_\theta \log p(a_t | s_t; \theta) \sum_{k=t}^{T} r(s_k, a_k) \right]$$

<span style="color:green">$\leftarrow \sim p(\tau; \theta)$</span>

$$= \sum_{t=0}^{T} E_{(s_0, a_0, \cdots, s_t, a_t)} \left[ \nabla_\theta \log p(a_t | s_t; \theta) \right.$$

<span style="color:green">$\sim p(\tau_{0:t}; \theta)$</span>

$$\left. \times \underbrace{E_{(s_{t+1}, a_{t+1}, \cdots, s_T, a_T, s_{T+1})} \left[ \sum_{k=t}^{T} r(s_k, a_k) \right]}_{} \right]$$

<span style="color:red">$Q(s_t, a_t)$</span>

<span style="color:green">$\sim p(\tau_{t+1:T} | \tau_{0:t}; \theta)$</span>

<span style="color:green">$$E[\cdots | \tau_{0:t}] = E[\cdots | s_t, a_t]$$</span>

<span style="color:green">Markov property of MDPs</span>

<span style="color:green">$$E_{(x,y)} \left[ f(x) g(x,y) \right] = E_x \left[ E_y \left[ f(x) g(x,y) | x \right] \right]$$</span>

<span style="color:green">$$= E_x \left[ f(x) E_y \left[ g(x,y) | x \right] \right]$$</span>

$$\nabla_\theta J(\theta) = \sum_{t=0}^{T} E_{\tau_{0:t} \sim p(\tau_{0:t};\theta)} \left[ \nabla_\theta \log p(a_t | s_t; \theta) \, Q_\theta^t(s_t, a_t) \right]$$

$$Q_\theta^t(s_t, a_t) = E_{\tau_{t:T}} \left[ \sum_{k=t}^{T} r(s_t, a_t) \Big| s_t, a_t \right]$$

$$E_x[f(x)] = E_{(x,y)}[f(x)] = E_x\left[ f(x) \, E_y[1 | x] \right]$$

$$\nabla_\theta J(\theta) = E_{\tau \sim p(\tau;\theta)} \left[ \sum_{t=0}^{T} \nabla_\theta \log p(a_t | s_t; \theta) \, Q_\theta^t(s_t, a_t) \right]$$

actor          critic

learn actor $p(a|s)$ and critic $Q(s,a)$ together
using policy-gradient and Q-learning

this is good because $Q$ is lower variance than
$\sum_t r(s_t, a_t)$

Now consider infinite time, average reward case (no discounting).

$J_{avg}(\theta)$ = average reward per step

$$= \lim_{T \to \infty} \frac{1}{T} E_{\tau_{0:T} \sim p(\tau_{0:T}; \theta)} \left[ \sum_{t=0}^{T} r(s_t, a_t) \right]$$

previous payoff

$$= \sum_{s \in S} \sum_{a \in A} p(s, a; \theta) \, r(s, a)$$

ergodicity

probability of being in $(a, s)$ as $t \to \infty$

invariant distribution of the MDP

$$= \sum_{s \in S} p(s; \theta) \sum_{a \in A} p(a \mid s; \theta) \, r(s, a)$$

Gradient of payoff $J_{avg}(\theta)$ in infinite-time, average-reward:

S&B (13.5), p 326:

$$\nabla_\theta J(\theta) = \sum_s p(s;\theta) \sum_a Q_\theta(s,a) \nabla_\theta p(a|s;\theta)$$

$$= \sum_s p(s;\theta) \sum_a Q_\theta(s,a) \nabla_\theta \log p(a|s;\theta) \, p(a|s;\theta)$$

$$= \sum_s \sum_a p(s,a;\theta) \nabla_\theta \log p(a|s;\theta) Q_\theta(s,a)$$

$$= E_{(s,a) \sim p(s,a;\theta)} \left[ \nabla_\theta \log p(a|s;\theta) Q_\theta(s,a) \right]$$

either samples over S×A or samples from on-policy trajectories

$$E_\pi [ \cdots ]$$

policy gradient theorem

$$\nabla_\theta J(\theta) = E_{(s,a) \sim p(s,a;\theta)} \left[ \nabla_\theta \log p(a|s;\theta) Q_\theta(s,a) \right]$$

# Actor – critic – Q method

$$\nabla_\theta J(\theta) = E_{(s,a) \sim p(s,a;\theta)} \left[ \nabla_\theta \log p(a|s;\theta) \, Q_\theta(s,a) \right]$$

*parameters of $p$*

Approximate $Q_\theta(s,a) \approx Q(s,a;w)$

*parameters of $Q$*

loop over episodes:
    $s \sim p(s)$
    $a \sim p(a|s;\theta)$
    loop over time:

*learning rate for $p$*    *$\approx \nabla_\theta J(\theta)$*    *policy gradient*

$$\theta \leftarrow \theta + \alpha \nabla_\theta \log p(a|s;\theta) \, Q(s,a;w)$$

$$s' \sim p(s'|s,a)$$
$$a' \sim p(a'|s';\theta) \quad \leftarrow \text{on policy}$$

$$\delta \leftarrow r(s,a) + Q(s',a';w) - Q(s,a;w) \Big\} \text{ Q-learning}$$
$$w \leftarrow w + \beta \, \delta \, \nabla_w Q(s,a;w)$$

$$s \leftarrow s', \quad a \leftarrow a'$$

*learning rate for $Q$*