Markov Decision Process

$(S, A, R, P, \gamma)$

State   action   reward   discount factor

action → reward

Policy $\pi$:    $\pi^* = \sum_t \gamma^t r$     $\pi^* = \max \sum_t \gamma^t r$

$\pi^* = \underset{\pi}{\arg\max} \; E\left[ \sum_{t \geq 0} \gamma^t r_t \mid \pi \right]$

Balman Equation.

$$Q^*(s,a) = E_{s' \sim \varepsilon}\left[ r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right].$$

Value iteration:

$$Q_{i+1}(s,a) = E\left[ r + \gamma \max_{a'} Q_i(s', a') \mid s, a \right].$$

$\underbrace{Q(s,a;\theta)}_{\downarrow} \approx Q^*(s,a)$.

approximator.

forward
$$L_i(\theta_i) = E_{s, a \sim \rho(\cdot)}\left[ (y_i - Q(s,a;\theta_i))^2 \right].$$

where: $y_i = E_{s' \sim \varepsilon}\left[ r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) \mid s, a \right].$

backward.
$$\nabla_{\theta_i} L_i(\theta_i) = E_{s, a \sim \rho(\cdot); s' \sim \varepsilon}\left[ r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s,a;\theta_i) \nabla_{\theta_i} Q(s,a;\theta_i) \right]$$

parametrized policies $\pi = \{ \pi_\theta, \theta \in R^m \}$
for each policy: $J(\theta) = E_{\tau \sim p(\tau;\theta)}\left[ r(\tau) \right]$
$$= \int_\tau r(\tau) p(\tau;\theta) d\tau$$

$r(\tau)$ is the reward of a trajectory $\tau = (s_0, a_0, r_0, s_1)$

trick   $\nabla_\theta p(\tau;\theta) = p(\tau;\theta) \dfrac{\nabla_\theta p(\tau;\theta)}{p(\tau;\theta)} = p(\tau;\theta) \cdot \nabla_\theta \log p(\tau;\theta)$

$$\nabla_\theta J(\theta) = \int_\tau (r(\tau) \nabla_\theta \log p(\tau;\theta)) p(\tau;\theta) d\tau$$
$$= E_{\tau \sim p(\tau;\theta)}\left[ r(\tau) \nabla_\theta \log p(\tau;\theta) \right]. \quad \text{Monte Carlo.}$$

$$p(\tau;\theta) = \prod_{t \geq 0} p(s_{t+1} \mid s_t, a_t) \pi_\theta(a_t \mid s_t)$$

Thus: $\log p(\tau;\theta) = \sum_{t \geq 0} \log p(s_{t+1} \mid s_t, a_t) + \log \pi_\theta(a_t \mid s_t)$

$\nabla_\theta \log p(\tau;\theta) = \sum_{t \geq 0} \nabla_\theta \log \pi_\theta(a_t \mid s_t)$

$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_\theta \log \pi_\theta(a_t \mid s_t)$

baseline:
$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} \left( \sum_{t' \geq t} \gamma^{t'-t} r_{t'} - b(s_t) \right) \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

Q function as baseline.
$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} (Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t)) \nabla_\theta \log \pi_\theta(a_t \mid s_t)$$

Universal approximation theorem.