

Problem Set 2

$$R: S \times A \times S \rightarrow \mathbb{R}.$$

$$MDP: \langle S, A, R, \gamma, \phi \rangle, \quad M' = \langle S, A, R', \gamma, \phi \rangle$$

$$\text{where } R'(s, a, s') = R(s, a, s') + F(s, a, s')$$

$$F: S \times A \times S \rightarrow \mathbb{R} = \text{feedback function}.$$

$$F(s, a, s') = \gamma \phi(s') - \phi(s); \quad \phi: S \rightarrow \mathbb{R};$$

$$\forall (s, a, s') \in S \times A \times S.$$

$$\textcircled{1} \quad Q_M^*, Q_{M'}^*$$

$$Q_M^*(s, a) - \phi(s) = Q_{M'}^*(s, a)$$

$$Q_{M'}^*(s, a) = \max_a R'(s, a, s') + \gamma \sum P(s'|s, a) V(s')$$

$$= \max_a (R(s, a, s') + F(s, a, s'))$$

$$Q_M^* = R'(s, a, s') + \gamma \max_{a'} Q_{M'}(s', a')$$

$$= R(s, a, s') + F(s, a, s') + \gamma \max_{a'} Q_{M'}(s', a')$$

$$= R(s, a, s') + \gamma \phi(s') - \phi(s) + \gamma \max_{a'} Q_{M'}(s', a')$$

$$= R(s, a, s')$$

$$Q_{M'}^* + \phi(s) = R(s, a, s') + \gamma \max [Q_{M'}^*(s', a') + \phi(s')]$$

$$= R(s, a, s') + \gamma \max Q_M(s', a') = Q_M(s, a);$$

$$\pi_{M'}^*(s) = \arg \max_{a'} Q_{M'}^*(s, a) = \arg \max_{a'} Q_M^*(s, a) - \phi(s)$$

$$= \arg \max_{a'} Q_M^*(s, a)$$

$$\leadsto \pi_{M'}^*(s) = \pi_M^*(s).$$

② $Q_{\mu}^0(s, a)$ & $Q_{\mu'}^0(s, a)$; $q_{init} \in \mathbb{R}$.

$$Q_{\mu}^0(s, a) = q_{init} + \phi(s); \quad Q_{\mu'}^0(s, a) = q_{init} = Q_{\mu}^0(s, a) - \phi(s)$$

$$Q_{\mu}(s, a) = Q_{\mu}^0(s, a) + \Delta Q_{\mu}(s, a)$$

$$Q_{\mu'}(s, a) = Q_{\mu}^0(s, a) + \Delta Q_{\mu'}(s, a)$$

$$\Delta Q_{\mu}(s, a) = \Delta Q_{\mu'}(s, a)$$

$$\Rightarrow Q_{\mu}(s, a) = Q_{\mu'}(s, a) + \phi(s)$$

$$Q_{\mu}(s, a) = R(s, a, s') + \gamma \max_{a' \in A'} Q_{\mu}(s', a')$$

$$= R(s, a, s') + \gamma \max_{a' \in A'} [Q_{\mu'}(s', a') + \phi(s')]$$

$\Rightarrow \argmax Q_{\mu}(s, a) \equiv \argmax Q_{\mu'}(s, a) \Rightarrow$ update coincides.

Review policy update process:

$$\pi(s) = 0$$

while $\pi(s)$ and $\pi'(s) \leq \epsilon$:

- Policy Eval (π) \Rightarrow get V ; Q ; ...

- Policy Update \Rightarrow update using policy iteration

$$\Rightarrow \pi = \argmax_a Q^{\pi_i}(s, a) \geq Q^{\pi_i}(s, a) = \pi_i(s).$$

Policy Evaluation:

- Analytical Method (know model)
- Monte Carlo Method (doesn't know model but can sample).
- Temporal Difference.

Temporal Difference.

$$V_t^{\pi} = \alpha V_t^{\pi} + \alpha (R_t + \gamma V_{t+1}^{\pi} - V_t^{\pi})$$

$$= \alpha (R_t + \gamma V_{t+1}^{\pi}) + (1 - \alpha) V_t^{\pi}$$

⊕ Also use sampling.

⊕ Same with $Q(s, a)$: $S \times A \rightarrow \mathbb{R}$.

With exploration: $\pi(a|s) = \begin{cases} 1 - \epsilon & \argmax \\ \frac{\epsilon}{|A|} & \text{otherwise.} \end{cases}$

\Rightarrow policy update with probability.

SARSA
Q-Learning:
Other sampling methods.