

Ex 11-1

(a) Prediction = policy evaluation

Control = optimal ~~poly~~ policy approximation

(b)

(1) Model-based: model of environment is known, i.e. transition probability is known

Model-free: model is not needed.

(2) Both evaluate value function

(3) Model-based: DP.

Model-free: MC, TD

(c) optimal stat value $u_*(s) = \max_a \sum_{s'} p(s'/s, a) [R(s') + \gamma u_*(s')]$

optimal action value $q_*(s, a) = \sum_{s'} p(s'/s, a) [R(s') + \gamma \max_{a'} q_*(s', a')]$

$$u_*(s) = \max_a q_{\pi_*}(s, a)$$

(d) on-policy: evaluate/optimize ~~policy~~ policy

- SARSA (on-policy TD control)

off-policy: evaluate/optimize target policy from samples generated by behavior policy

- Q-learning (off-policy TD control)

- MC/TD - prediction/control via importance sample.

$$(e) \quad Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, A_t)]$$

Q-learning: back up best Q-value

from state reached in the observed transition

SARSA: wait until an action is actually taken

backup the Q-value for that action

$$\begin{cases} Q(s, A) \leftarrow^{\alpha} R + \gamma Q(s', A') \\ Q(s, A) \leftarrow^{\alpha} R + \gamma \max_{a' \in A} Q(s', a') \end{cases} \quad x \leftarrow^{\alpha} y \equiv x \leftarrow x + \alpha(y - x)$$

(f) Exploration: unknown environment, unknown state effect.
 \Rightarrow better future, discover new action

Large reward \Leftrightarrow ~~the~~ prefer experience from past.

Exploitation: ~~try new discover new actions~~ \Rightarrow ~~better future~~
 obtain reward

$$(g) \quad G_t^{\lambda} = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)} \quad \begin{matrix} G_t^0 \Rightarrow TD(0) \Rightarrow TD \\ G_t^1 \Rightarrow MC \end{matrix} \quad \left. \begin{matrix} \\ \end{matrix} \right\} \text{backward}$$
~~$$\lambda=0 \Rightarrow G_t^{\lambda} = \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$~~

Ex 11-2. (a) DP: $U(s) \leftarrow \mathbb{E}[R + \gamma U(s') | s] = \sum_{s'} p(s'|s) [R + \gamma U(s')]$

MC: $U(s) \leftarrow U(s) + \alpha (G_t - U(s))$

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

TD: $U(s) \leftarrow U(s) + \alpha (R + \gamma U(s') - U(s))$

(b) DP: + low variance

+ ~~high~~ efficient

- model based

- high bias

MC: - high variance

- low efficient

+ model-free

+ no-bias

TD: + low variance

+ efficient

+ model-free

+ high bias