


Chapter 7 TD learning of state values.

- TD learning

1. The data/experience required:

$(s_0, r_1, s_1, \dots, s_t, r_{t+1}, s_{t+1}, \dots)$ or $\{(s_t, r_{t+1}, s_{t+1})\}_t$
generated following the given policy π .

~~1.~~ $V_t(s) = V_t(s_t) - \alpha_t(s_t) [V_t(s_t) - [r_{t+1} + \gamma V_t(s_{t+1})]]$

2. $V_{t+1}(s) = V_t(s)$

$$V_{t+1}(s_t) = \underbrace{V_t(s_t)}_{\text{New estimate}} - \underbrace{\alpha_t(s_t)}_{\text{old estimate}} \left[\underbrace{V_t(s_t)}_{\text{TD error } \delta_t} - \underbrace{[r_{t+1} + \gamma V_t(s_{t+1})]}_{\text{TD target } \bar{V}_t} \right]$$

$$\delta_t = V_t(s_t) - [r_{t+1} + \gamma V_t(s_{t+1})]$$

$$\bar{V}_t = r_{t+1} - \gamma V_t(s_{t+1})$$

3. TD target:

$$V_{t+1}(s_t) = V_t(s_t) - \alpha_t(s_t) [V_t(s_t) - \bar{V}_t]$$

$$V_{t+1}(s_t) - \bar{V}_t = V_t(s_t) - \bar{V}_t - \alpha_t(s_t) [V_t(s_t) - \bar{V}_t]$$

$$V_{t+1}(s_t) - \bar{V}_t = [1 - \alpha_t(s_t)] [V_t(s_t) - \bar{V}_t]$$

$$0 < 1 - \alpha_t(s_t) < 1$$

$$\therefore |V_{t+1}(s_t) - \bar{V}_t| \leq |V_t(s_t) - \bar{V}_t|$$

4. TD error :

$$\delta_t = \underline{V(S_t)} - \underline{[r_{t+1} + \gamma V(S_{t+1})]}$$

$$V_t = V_\pi \quad \delta_t = 0$$

$$\delta_{\pi,t} = V_\pi(S_t) - [r_{t+1} + \gamma V_\pi(S_{t+1})]$$

$$E[\delta_{\pi,t} | S_t = s_t] = 0$$

TD error : innovation . 新信息.

5. Other properties :

TD only estimates the value of a given policy

二. Convence and comparation.

TD algorithm solves the Bellman equation of a given policy π . without model.

$$V_\pi(s) = E[R + \gamma G | S=s]$$

$$E[G | S=s] = \sum_a \pi(a|s) \sum_{S'} p(S'|s,a) V_\pi(S') = E[V_\pi(S') | S=s]$$

$$V_\pi(s) = E[R + \gamma V_\pi(S') | S=s], s \in S \quad ①$$

Using RM algorithm solve equation ①

$$g(V(s)) = V(s) - E[R + \gamma V_{\pi}(s')] \geq 0$$

$$g(V(s)) = 0$$

obtain samples: $r \rightarrow R$ and $s \rightarrow s'$

$$\tilde{g}(V(s)) = V(s) - [r + \gamma V_{\pi}(s')]$$

$$= (\underbrace{V(s) - E}_{g(V(s))}) + (\underbrace{E - [r + \gamma V_{\pi}(s')]}_{\eta})$$

RM algorithm for solving $g(V(s)) \equiv 0$

$$V_{k+1}(s) = V_k(s) - \alpha_k \tilde{g}(V_k(s))$$

$$V_{k+1}(s) = V_k(s) - \alpha_k (V(s) - [r_k + \gamma V_{\pi}(s')])$$

$$\{(s, r, s')\} \rightarrow \{(s_t, r_{t+1}, s_{t+1})\}$$

访问到什么状态更新什么状态，其余保持不动。

$$V_{\pi}(s') \Rightarrow V_{\pi}(s'_k)$$

2. TD learning of action value.: Sarsa

1. a given π : $\{(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})\}$

$$\left\{ \begin{array}{l} q_{t+1}(s, a) = q_t(s, a), \quad t(s, a) \neq (s_t, a_t) \\ q_t(s_t, a_t) \rightarrow v \end{array} \right.$$

where, $t = 0, 1, 2 \dots$

$q_t(s_t, a_t)$ is a estimate of $q_\pi(s_t, a_t)$

2. Sarsa = State - action - reward - state - action.
 $s_t \quad a_t \quad r_{t+1} \quad s_{t+1} \quad a_{t+1}$

3. TD solve

$$\text{Sarsa solve: } q_\pi(s, a) = E[R + \gamma q_\pi(s', a') | s, a]$$

$$V_\pi(s) = E[R + \gamma G | S=s]$$

4. combine Sarsa with policy improvement.

① Collect experience($s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}$)
following $\pi_t(s_t)$

② Update q-value

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) + \alpha_t(s_t, a_t) [q_t(s_t, a_t) - [r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})]]$$

③ Update policy:

$$\pi_{t+1}(a|s_t) = 1 - \frac{\epsilon}{|A|} (|A|-1)$$

$$\text{if } a = \arg \max_a q_{t+1}(s_t, a)$$

$$\pi_{t+1}(a|s_t) = \frac{\epsilon}{|A|} \text{ other wise.}$$

ϵ -greedy

三. Expected Sarsa: / n-step Sarsa.

① Expected Sarsa:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) + \alpha_t(s_t, a_t) [q_t(s_t, a_t) - [r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})]]$$

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) + \alpha_t(s_t, a_t) [q_t(s_t, a_t) - [r_{t+1} + \gamma E[q_t(s_{t+1}, A)]]]$$

$$E[q_t(s_{t+1}, A)] = \sum_a \pi_t(a | s_{t+1}) q_t(s_{t+1}, a) = V_t(s_{t+1})$$

② n-step Sarsa:

$$q_\pi(s, a) = E[G_t | S_t = s, A_t = a]$$

$$\text{Sarsa} \leftarrow G_t^{(1)} = R_{t+1} + \gamma q_\pi(s_{t+1}, a_{t+1})$$

$$G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 q_\pi(s_{t+1}, a_{t+1})$$

$$\text{nStep Sarsa } G_t^{(n)} = R_{t+1} + \dots + \gamma^n q_\pi(s_{t+n}, a_{t+n})$$

$$MC \leftarrow G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots$$

四 Q-Learning:

estimate. Optimal action value

直接估计 ~~或~~ ^b action value

$$\left\{ \begin{array}{l} 1. Q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha(s_t, a_t) [q_t(s_t, a_t) - [r_{t+1} + \gamma \max_{a \in A} q_t(s_{t+1}, a)]] \\ q_{t+1}(s, a) = q_t(s, a) \end{array} \right.$$

TD target in Q-learning: $r_{t+1} + \gamma \max_{a \in A} q_t(s_{t+1}, a)$

TD target in Sarsa: $r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})$

2. Q-learning solve Bellman optimality equation

$$q(s, a) = E [R_{t+1} + \gamma \max_a q(s_{t+1}, a) | s_t = s, a_t = a]$$

3. Off-policy & on-policy (对冲的强化学习)

$$\left\{ \begin{array}{l} \text{behavior policy} = \text{generate experience} \\ \text{target policy} = \text{constantly update towards optimal policy} \end{array} \right.$$

on policy : behavior policy = target policy

off policy : behavior policy \neq target policy

Sarsa - on policy. Q-learning - off policy

4. How to judge =

① Sarsa: on policy :

solve: $q_{\pi}(s, a) = E[R + \gamma q_{\pi}(s', a') | s, a]$

where, $R \sim p(R|s, a)$, $s' \sim p(s'|s, a)$, $a' \sim \pi(a'|s')$

algorithm:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) + \alpha(s_t, a_t) [q_t(s_t, a_t) - [r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})]]$$

requires : $s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}$
- $\underbrace{s_t}_{\text{given}} \underbrace{a_t}_{p(r|s,a)} \underbrace{r_{t+1}}_{p(s'|s,a)} \underbrace{s_{t+1}}_{\pi_t(s_{t+1})} \underbrace{a_{t+1}}$

π_t both target policy and behavior policy.

$$\pi_t \rightarrow \exp \rightarrow q_{\pi(t)}$$

② Monte Carlo learning: on policy.

$$q_{\pi}(s, a) = E[R_{t+1}, \gamma R_{t+2}, \dots | S_t=s, A_t=a]$$

$$q(s, a) \approx r_{t+1} + \gamma r_{t+2} + \dots$$

$\pi \rightarrow$ trajectory $\rightarrow q(\pi) \rightarrow \pi$.

③ Q-learning: off policy:

$$q(s, a) = E[R_{t+1}, \gamma \max_a q(s_{t+1}, a) | S_t=s, A_t=a], \forall s, a$$

bellman optimal function.

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) [q_t(s_t, a_t) - [r_{t+1} + \gamma \max_{a'} q_t(s_{t+1}, a')]]$$

$$\frac{s_t, a_t}{\text{given}} \quad \frac{s_{t+1}, a_{t+1}}{p(s'|s, a)} \quad p(s' | s, a)$$

Q-learning = behavior policy $s_t, a_t \rightarrow s_{t+1}, a_{t+1}$
target policy π_{optimal}

5. Pseudocode: Q-learning

given: $\pi_b \rightarrow \text{episode: } \{s_0, a_0, r_1, s_1, a_1, r_2, \dots\}$

For each step $t = 0, 1, 2, \dots$ of the episode

Update q-value:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) [q_t(s_t, a_t) - [r_{t+1} + \gamma \max_{a \in A} q_t(s_{t+1}, a)]]$$

Update target policy:

greedy select $\pi_{T, t+1}(a|s_t) = 1 \quad \text{if } a = \arg \max_a q_{t+1}$

不需要具有一定探索性.

① behavior policy selected. (探索性比较好)

② Update.

2.

unified expression:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t)[q_t(s_t, a_t) - \bar{q}_t]$$

All the algorithms we introduced in this lecture can be expressed in a unified expression:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t)[q_t(s_t, a_t) - \bar{q}_t],$$

where \bar{q}_t is the TD target.

Different TD algorithms have different \bar{q}_t .

Algorithm	Expression of \bar{q}_t
Sarsa	$\bar{q}_t = r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})$
n -step Sarsa	$\bar{q}_t = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^n q_t(s_{t+n}, a_{t+n})$
Expected Sarsa	$\bar{q}_t = r_{t+1} + \gamma \sum_a \pi_t(a s_{t+1}) q_t(s_{t+1}, a)$
Q-learning	$\bar{q}_t = r_{t+1} + \gamma \max_a q_t(s_{t+1}, a)$
Monte Carlo	$\bar{q}_t = r_{t+1} + \gamma r_{t+2} + \dots$

The MC method can also be expressed in this unified expression by setting

$\alpha_t(s, a) = 1$ and hence $\alpha_t(s, a) = \bar{q}_t$.

所以这个是可以看成是n-step Sarsa的一个特殊情况

All the algorithms can be viewed as stochastic approximation algorithms solving the Bellman equation or Bellman optimality equation:

Algorithm	Equation aimed to solve
Sarsa	BE: $q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) S_t = s, A_t = a]$
n -step Sarsa	BE: $q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n q_{\pi}(s_{t+n}, a_{t+n}) S_t = s, A_t = a]$
Expected Sarsa	BE: $q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma \mathbb{E}_{A_{t+1}}[q_{\pi}(S_{t+1}, A_{t+1})] S_t = s, A_t = a]$
Q-learning	BOE: $q(s, a) = \mathbb{E}[R_{t+1} + \max_a q(S_{t+1}, a) S_t = s, A_t = a]$
Monte Carlo	BE: $q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots S_t = s, A_t = a]$