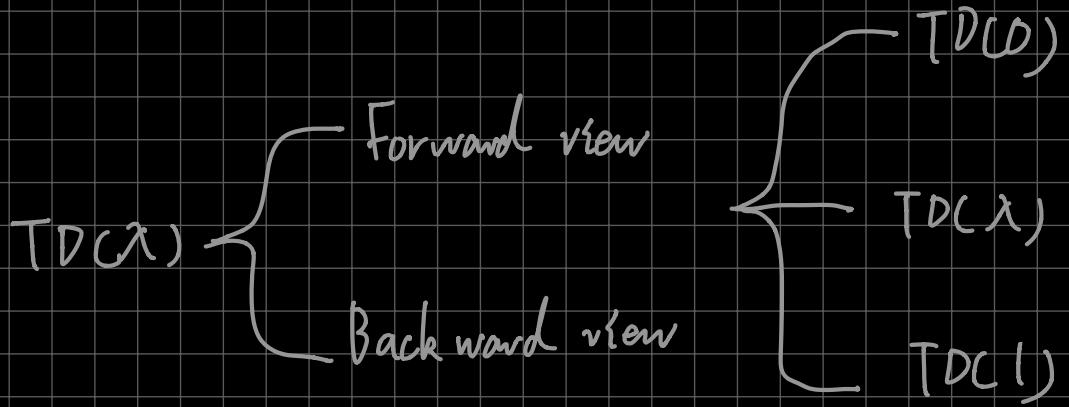


Lecture 5 TD(λ)



Forward view TD(λ):

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t^\lambda - V(s_t))$$

$$G_t^\lambda = (1-\lambda) \cdot \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

↪ λ return
 ↪ n step return
 ↪ 注意
 二者的节
 写: G_t^λ ,
 $G_t^{(n)}$

$$G_t^{(n)} = R_{t+1} + \gamma \cdot R_{t+2} \dots + \underbrace{\gamma^{n-1} \cdot R_{t+n}}_{n \text{ step}} + V(s_{t+n})$$

$\downarrow n \rightarrow \infty$

$$MC: G_t^\infty = R_{t+1} + \gamma \cdot R_{t+2} \dots + \gamma^{T-1} \cdot \underbrace{R_T}_{\text{terminal}}$$

Forward: $(s_t) \rightarrow (s_{t+1}) \rightarrow (s_{t+2}) \dots \rightarrow (s_T)$

Backward view TD(λ):

$$V(\underline{s}_t) \leftarrow \underline{V(s_t)} + \alpha \cdot \underline{s_t} \cdot E_t(s) \quad \text{小写}$$

① TD error ② Eligibility trace

$$\Delta s_t = R_{t+1} + \gamma \cdot \underline{V(s_{t+1})} - \underline{V(s_t)} \quad \text{大写}$$

$$\begin{cases} E_t(s) = \underbrace{\gamma \cdot \lambda \cdot E_{t-1}(s)}_{\text{Recency Heuristic}} + \underbrace{I(s_t=s)}_{\text{Frequency heuristic}}; \\ E_0(s) = 0 \end{cases}$$



Summary of Forward view TD(λ) and Backward view TD(λ)

Offline updates	$\lambda = 0$	$\lambda \in (0, 1)$	$\lambda = 1$
Backward view	TD(0) 	TD(λ) 	TD(1)
Forward view	TD(0)	Forward TD(λ)	MC
Online updates	$\lambda = 0$	$\lambda \in (0, 1)$	$\lambda = 1$
Backward view	TD(0) 	TD(λ) 	TD(1)
Forward view	TD(0) 	Forward TD(λ) 	MC
Exact Online	TD(0)	Exact Online TD(λ)	Exact Online TD(1)

① In offline setting, Backward view and Forward view are equal.

$$MC = TD(\lambda)$$



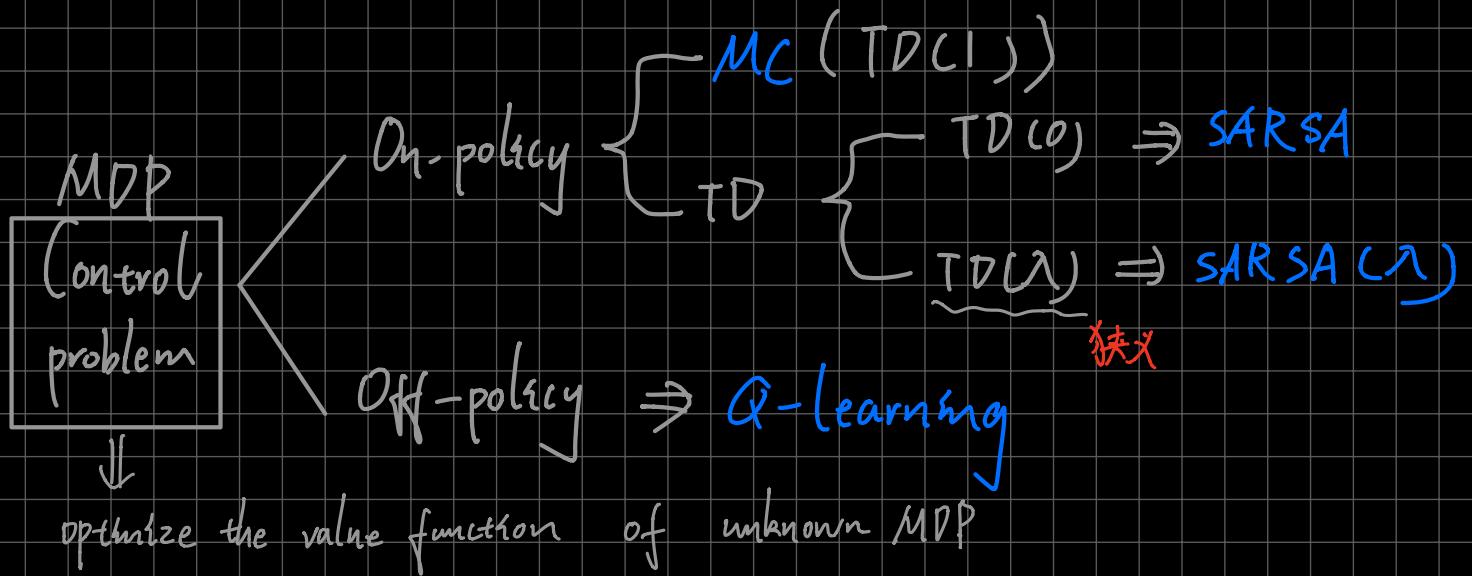
② In online setting, Backward view and Forward view are not equal (except $TD(0)$, obviously).

$$\underbrace{\text{Forward } TD(1)}_{\text{Exact Online } TD(1)} = MC$$

$$\underbrace{\text{Forward } TD(\lambda)}_{\text{Exact Online } TD(\lambda)} = MC$$

Forward means "online"

Lecture 6 MDP (Control)



MDP Prediction problem :

(lecture 4, 5) \Rightarrow estimate the value function of unknown MDP
 \Rightarrow ~~fix TD(λ)~~

(lecture 3) \Rightarrow estimate & optimize the value function of known MDP
 policy evaluation policy improvement

MC =

On-policy first-visit MC control (for ϵ -soft policies), estimates $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$:

$Q(s, a) \leftarrow$ arbitrary

$Returns(s, a) \leftarrow$ empty list

$\pi(a|s) \leftarrow$ an arbitrary ϵ -soft policy

Repeat forever:

(a) Generate an episode using π

(b) For each pair s, a appearing in the episode:

$\underbrace{\text{policy evaluation}}_{G} \leftarrow$ return following the first occurrence of s, a

Append G to $Returns(s, a)$

$Q(s, a) \leftarrow \text{average}(Returns(s, a))$

(c) For each s in the episode:

$\underbrace{\text{policy improvement}}_{A^* \leftarrow \arg \max_a Q(s, a)}$

For all $a \in \mathcal{A}(s)$:

$$\pi(a|s) \leftarrow \begin{cases} 1 - \epsilon + \epsilon / |\mathcal{A}(s)| & \text{if } a = A^* \\ \epsilon / |\mathcal{A}(s)| & \text{if } a \neq A^* \end{cases}$$

why use ϵ -greedy, not use greedy?

Greedy strategy is model-free, so in $\underbrace{\text{on-policy}}$ situation can't ensure the best result. ϵ -greedy is the simplest way to ensure continual exploration.

$\Rightarrow \epsilon$ -greedy, with ϵ prob randomly choose an action, $1 - \epsilon$ prob greedy choose.

SARSA:

Initialize $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
 Repeat (for each episode):

Initialize S

Choose A from S using policy derived from Q (e.g., ε -greedy)

Repeat (for each step of episode):

Take action A , observe R, S'

Choose A' from S' using policy derived from Q (e.g., ε -greedy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

$$S \leftarrow S'; A \leftarrow A'; \Rightarrow \text{on-policy}$$

until S is terminal

Here we can directly choose from V -table, or in-directly choose from π -table.

But for latter option, we need to perform policy improvement for π -table.
 can be skip if we use former method.

SARSA(λ):

Initialize $Q(s, a)$ arbitrarily, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat (for each episode):

$E(s, a) = 0$, for all $s \in \mathcal{S}, a \in \mathcal{A}(s) \Rightarrow$ init the eligibility trace

Initialize $S, A \Rightarrow$ compare to SARSA, which use ε -greedy choose A .

Repeat (for each step of episode): But in SARSA(λ), each episode reaches the end, so we just need to init A , not choose A .

Take action A , observe R, S'

$$\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$$

$$E(S, A) \leftarrow E(S, A) + 1$$

For all $s \in \mathcal{S}, a \in \mathcal{A}(s) \Rightarrow$ for each step of episode, update all the state,

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a) \text{ action item in } Q\text{-table to get } \lambda\text{-returns.}$$

$$E(s, a) \leftarrow \gamma \lambda E(s, a)$$

$S \leftarrow S'; A \leftarrow A' \Rightarrow$ on-policy, go to next step.

until S is terminal

$$\text{Forward: } Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \cdot (q_t^\lambda - Q(S_t, A_t))$$

$$q_t^\lambda = (1-\lambda) \cdot \sum_{n=1}^{\infty} \lambda^{(n-1)} \cdot \underbrace{q_t^{(n)}}_{n \text{ step } Q \text{ returns}}$$

$$q_t^{(n)} = R_{t+1} + \gamma \cdot R_{t+2} + \dots + \gamma^{n-1} \cdot R_{t+n} + Q(S_{t+n})$$

$$\text{Backward: } Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \cdot \delta_t \cdot E_t(s, a)$$

$$\delta_t = R_{t+1} + \gamma \cdot Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$

$$E_t(s, a) = \gamma \cdot \pi \cdot E_{t-1}(s, a) + \underbrace{\left((S_t = s, A_t = a) \right)}$$

Q-learning:

Initialize $Q(s, a), \forall s \in S, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
 Repeat (for each episode):

Initialize S

behavior policy (π) is ϵ -greedy

Repeat (for each step of episode):

Choose A from S using policy derived from Q (e.g., ϵ -greedy)

Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$; \Rightarrow off-policy

until S is terminal

target policy (π) is greedy

On-policy vs Off-policy:

- On policy is generally simple.

- Off policy is more powerful and general.

But it has greater variance and slower convergence.

Relationship Between DP and TD

	<i>Full Backup (DP)</i>	<i>Sample Backup (TD)</i>
Bellman Expectation Equation for $v_\pi(s)$ <i>predict problem</i>	<p>$v_\pi(s) \leftrightarrow s$</p> <p>$v_\pi(s') \leftrightarrow s'$</p>	<p>TD Learning</p>
Bellman Expectation Equation for $q_\pi(s, a)$	<p>$q_\pi(s, a) \leftrightarrow s, a$</p> <p>$q_\pi(s', a') \leftrightarrow a'$</p>	<p>S, A</p> <p>R</p> <p>S'</p> <p>A'</p> <p>Sarsa</p>
Bellman Optimality Equation for $q_*(s, a)$ <i>control problem</i>	<p>$q_*(s, a) \leftrightarrow s, a$</p> <p>$q_*(s', a') \leftrightarrow a'$</p>	<p>Q-Learning</p>

value iteration = policy iteration + policy improvement

Relationship Between DP and TD (2)

Predict Problem

Full Backup (DP)

Iterative Policy Evaluation

$$V(s) \leftarrow \mathbb{E}[R + \gamma V(S') | s]$$

Q-Policy Iteration

$$Q(s, a) \leftarrow \mathbb{E}[R + \gamma Q(S', A') | s, a]$$

Q-Value Iteration (沒有 a ~~re~~ Value Iteration)

$$Q(s, a) \leftarrow \mathbb{E} \left[R + \gamma \max_{a' \in \mathcal{A}} Q(S', a') | s, a \right]$$

Sample Backup (TD)

TD Learning

$$V(S) \xleftarrow{\alpha} R + \gamma V(S') \Rightarrow \text{TD}(\lambda)$$

Sarsa

$$Q(S, A) \xleftarrow{\alpha} R + \gamma Q(S', A') \Rightarrow \text{SARSAC}$$

Q-Learning

$$Q(S, A) \xleftarrow{\alpha} R + \gamma \max_{a' \in \mathcal{A}} Q(S', a')$$

where $x \xleftarrow{\alpha} y \equiv x \leftarrow x + \alpha(y - x)$

Control problems