

Reinforcement Learning

Temporal Difference - Prediction

TD - Introduction

DP
→ Bootstrapping.
Perfect model of
the env.

MC
→ Model free - Samples of
episodes.
 G_t

TD - Prediction

Monte - calro update rule (contant α MC):

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

new est \leftarrow old est + LR [Target - Old estimate]

$$R_{t+1} + \gamma V(S_{t+1})$$

$$\begin{aligned} V_{\pi}(S) &= E_{\pi}[G_t | S_t = S] \\ &= E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = S] \\ &= E_{\pi}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = S] \end{aligned}$$

TD(0) update rule:

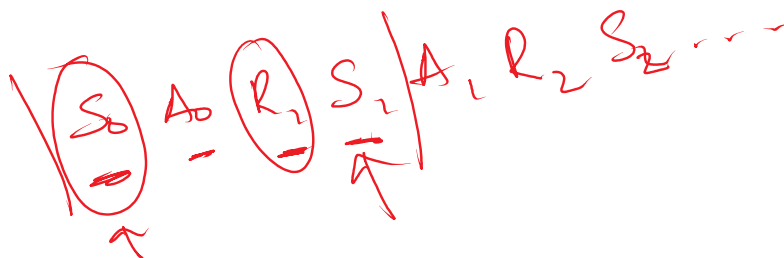
$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

TD-error

TD(0)

TD(0) - Prediction Algorithm

↓



Tabular TD(0) for estimating v_π

Input: the policy π to be evaluated

Algorithm parameter: step size $\alpha \in (0, 1]$

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop for each episode:

Initialize S

Loop for each step of episode:

$A \leftarrow$ action given by π for S

Take action A , observe R, S'

$V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

until S is terminal

Reinforcement Learning

On-Policy TD Control : SARSA



SARSA - Introduction

TD(0) update rule for action values:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

(S, A, R, S, A)

$(S, A) \quad R \quad (S, A)$

SARSA - Algorithm

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

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

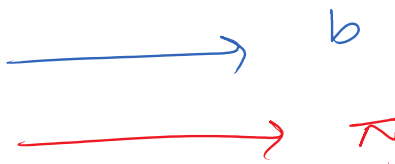
Reinforcement Learning

Off-Policy TD Control : Q-Learning

Q-Learning : Algorithm:

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

- Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
- Initialize $Q(s, a)$, for all $s \in S^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$
- Loop for each episode:
 - Initialize S
 - Loop for each step of episode:
 - Choose A from S using policy derived from Q (e.g., ε -greedy)
 - Take action A , observe R, S'
 - $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
 - $S \leftarrow S'$
 - until S is terminal



Comparison between Q-learning and SARSA:

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

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

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a')]$$

$$[R + \gamma E_{\pi}[Q(S_{t+1}, A_{t+1}) | S_t]]$$

$$(\sum_a \pi(a | S_{t+1}) \cdot Q(S_{t+1}, a))$$

Q-learning (off-policy TD control) for estimating $Q \approx q_*$

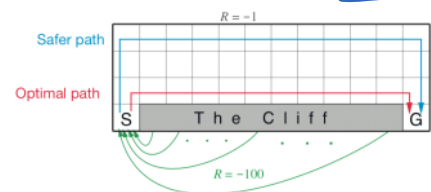
Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$
 Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$
 Loop for each episode:
 Initialize S
 Loop for each step of episode:
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 Take action A , observe R, S'
 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
 $S \leftarrow S'$
 until S is terminal

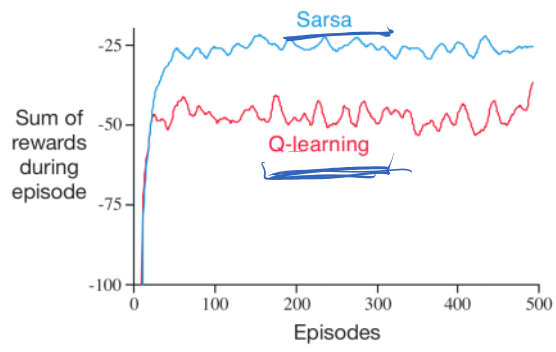
$$q_*(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]$$

but greedy
 $\pi \rightarrow$ greedy.

SARSA and Q-Learning on Cliff-Walking task:

1. Undiscounted, episodic task.
2. Task is to move from start state (S) to goal state (G)
3. Actions: Up, Down, Left and Right.
4. Deterministic environment.
5. Reward of -1 on all transitions, except those into the region marked as "The Cliff". Stepping into this region incurs a reward of -100 and sends the agent instantly to the start state.





Action selected according to ϵ -greedy ($\epsilon = 0.1$)

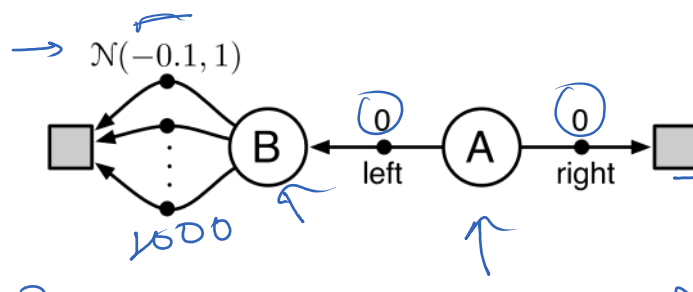


Reinforcement Learning Double Q-Learning

Maximization Bias:

$$SARSA: \quad Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

$$Q-learning: \quad Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$



1000 ↑

$$\begin{aligned} A_1 &\rightarrow -0.2 \\ A_2 &\rightarrow -0.1 \\ A_3 &\rightarrow +0.4 \end{aligned} \qquad \begin{aligned} A_1 &\rightarrow \underline{0.5} \end{aligned}$$

Double Learning:

Q → Best action
 Q → Estimated value for that action.

$$Q(a) \rightarrow \begin{matrix} Q_1(a) \\ \underline{Q_2(a)} \end{matrix}$$

$$A^* = \underset{a}{\operatorname{argmax}} Q_1(a)$$

$$Q_2(A^*)$$

$$\begin{matrix} Q_1(S_1, A_1) \\ Q_2(S_1, A_2) \end{matrix}$$

$$\begin{aligned} \rightarrow Q_1(S_t, A_t) &\leftarrow Q_1(S_t, A_t) + \alpha [R_{t+1} + \gamma \underline{Q_2(S_{t+1}, \operatorname{argmax}_a Q_1(S_{t+1}, a))} - Q_1(S_t, A_t)] \\ \rightarrow \underline{Q_2(S_t, A_t)} &\leftarrow Q_2(S_t, A_t) + \alpha [R_{t+1} + \gamma \underline{Q_1(S_{t+1}, \operatorname{argmax}_a Q_2(S_{t+1}, a))} - Q_2(S_t, A_t)] \end{aligned}$$

Double Q-Learning Algorithm:

Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, such that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

Initialize S

Loop for each step of episode:

Choose A from S using the policy ε -greedy in $Q_1 + Q_2$

Take action A , observe R, S'

With 0.5 probability:

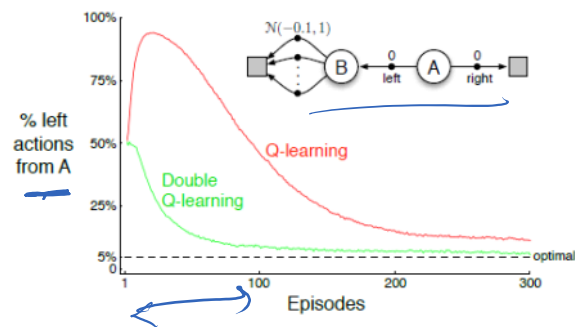
$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left(R + \gamma Q_2(S', \arg\max_a Q_1(S', a)) - Q_1(S, A) \right)$$

else:

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left(R + \gamma Q_1(S', \arg\max_a Q_2(S', a)) - Q_2(S, A) \right)$$

$S \leftarrow S'$

until S is terminal



Reinforcement Learning

n-step TD and TD(λ)

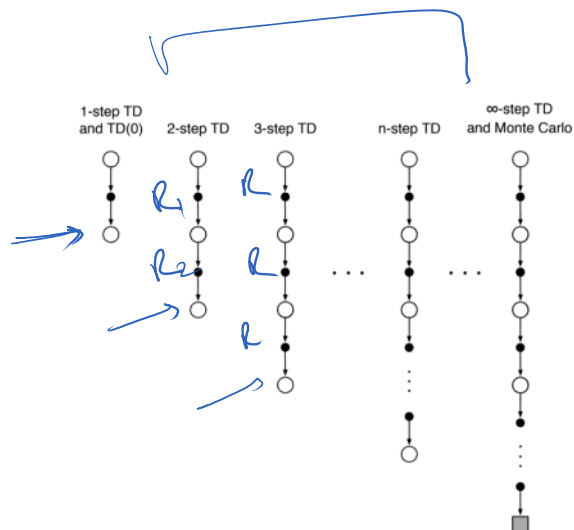
n-step Temporal Difference Learning:

Monte carlo update rule :

$$V(s_t) = V(s_t) + \alpha [G_t - V(s_t)]$$

TD(0) (one-step TD) update rule:

$$V(s_t) = V(s_t) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$



TD(0)

$$0.5 \times 4x^{(2)} + 0.5 \times 4x^{(5)}$$

n -step returns for $n = 1, 2, \dots, \infty$:

$$n = 1 : G_t^{(1)} = R_{t+1} + \gamma V(s_{t+1})$$

$$n = 2 : G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(s_{t+3})$$

\vdots

$$n = \infty : G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-1} R_T$$

General n -step return:

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

n -step TD update:

$$V(s_t) = V(s_t) + \alpha [G_t^{(n)} - V(s_t)]$$

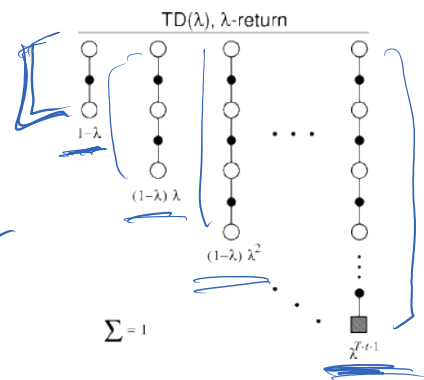
TD(λ):

Using a weight of $(1 - \lambda)\lambda^{n-1}$ for n^{th} return:

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

TD(λ) update rule:

$$V(s_t) = V(s_t) + \alpha [G_t^\lambda - V(s_t)]$$



200

$$1-\lambda + (1-\lambda)\lambda + (1-\lambda)\lambda^2 + \dots = 1$$

$$\cancel{1-\lambda} + \cancel{\lambda - \lambda^2} + \lambda^2 = 1$$

$$\Rightarrow \lambda^2 = 1$$