

Reinforcement Learning

Exercise 3: Model-free Control, SARSA & Q-Learning

Nico Meyer

Exercise Sheet 3

Gymnasium Gridworld

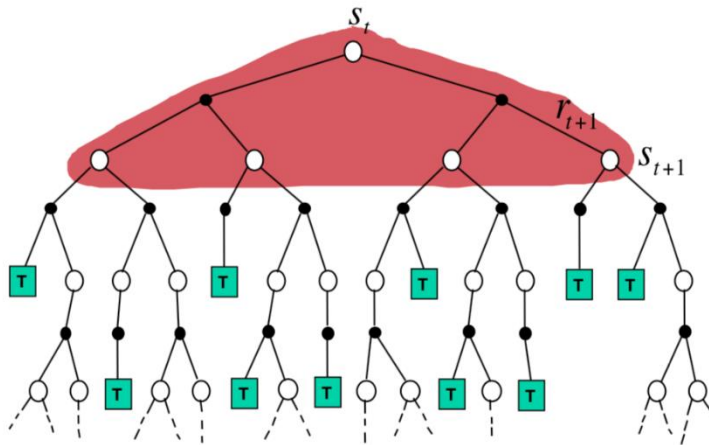


Exercise Sheet 3

Temporal Difference Learning

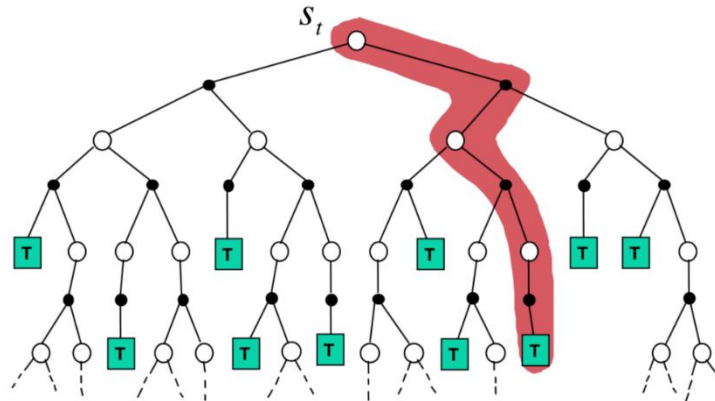
DP Backup

$$V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$$



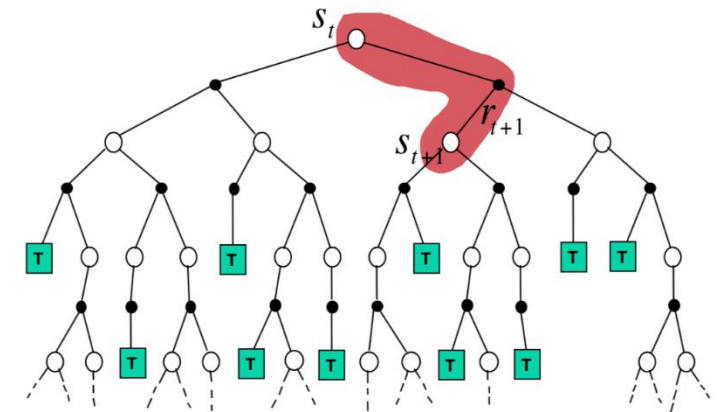
MC Backup

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



TD Backup

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



Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.

Model-free Control

TD Methods



Model-free Prediction vs. Control

Problem:

We do not know \mathcal{P} or \mathcal{R} or both of the MDP $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$

Solution:

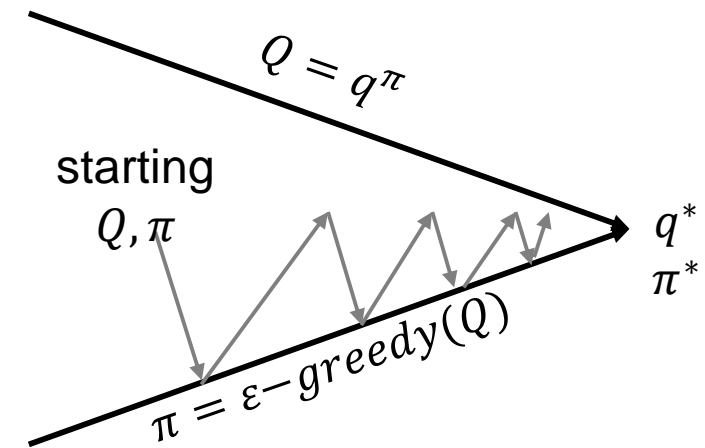
Model-free methods that use experience samples $s(s, a, r, s')$

In Exercise 3 we did:

Model-free Prediction: Evaluate the future, given the policy π .
(estimate the value function)

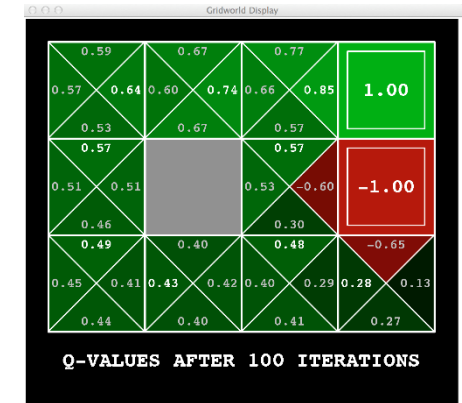
In Exercise 4 we will do:

Model-free Control: Optimize the future by finding the best policy π .
(optimize the value function)



State-action-value function

$$\begin{aligned}
 & S \xrightarrow{a, r_0} S_1 \xrightarrow{\pi(S_1), r_1} S_2 \xrightarrow{\pi(S_2), r_2} S_3 \dots S_{h-1} \xrightarrow{\pi(S_{h-1}), r_{h-1}} S_h \\
 & Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a \right]
 \end{aligned}$$



Greedy Policy Improvement over Q:

$$\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^\pi(s, a)$$

$$\forall s \in \mathcal{S}, \quad Q^{\pi'}(s, \pi'(s)) \geq Q^\pi(s, \pi(s))$$

SARSA

on-policy control

- Apply TD to $Q(s, a)$
- Use ε -greedy policy improvement
- Update at every time-step

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

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

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

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

 Loop for each step of episode:

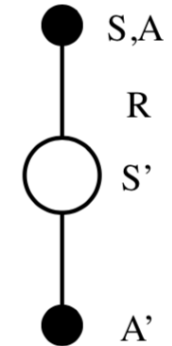
 Take action A , observe R, S'

 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



Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.

Q-Learning

off-policy control

- Evaluate one policy while following another
- Can re-use experience gathered from old policies

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 $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

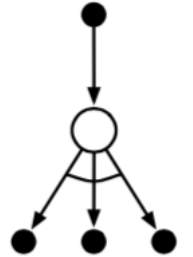
 Choose A from S using policy derived from Q (e.g., ε -greedy)

 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



Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.

Epsilon-greedy policy

Why should we follow an ϵ -greedy policy? Isn't this suboptimal?

Exercise Sheet 4



Thank you for your attention!