

# Sutton & Barto RL Cheatsheet

Pierluca D'Oro

## Monte Carlo Exploring Starts

Initialize, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$   
     $Q(s, a) \leftarrow$  arbitrary  
     $\pi(s) \leftarrow$  arbitrary  
     $Returns(s, a) \leftarrow$  empty list

Repeat forever:  
    Choose  $S_0 \in \mathcal{S}$  and  $A_0 \in \mathcal{A}(S_0)$  randomly  
    Generate episode starting from  $S_0, A_0$ , using  $\pi$   
    For each pair  $s, a$  appearing in episode:  
         $G \leftarrow$  return following first  $s, a$  occurrence  
        Append  $G$  to  $Returns(s, a)$   
         $Q(s, a) \leftarrow \text{average}(Returns(s, a))$   
    For each  $s$  in the episode:  
         $\pi(s) \leftarrow \text{argmax}_a Q(s, a)$

## On policy first-visit MC control

Initialize, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ :  
     $Q(s, a) \leftarrow$  arbitrary  
     $Returns(s, a) \leftarrow$  empty list  
     $\pi(a|s) \leftarrow$  arbitrary  $\epsilon$ -soft policy

Repeat forever:  
    (a) Generate an episode using  $\pi$   
    (b) For each pair  $s, a$  appearing in the episode:  
         $G \leftarrow$  return following first  $s, a$  occurrence  
        Append  $G$  to  $Returns(s, a)$   
         $Q(s, a) \leftarrow \text{average}(Returns(s, a))$   
    (c) For each  $s$  in the episode:  
         $A^* \leftarrow \text{argmax}_a Q(s, a)$   
        For all  $a \in \mathcal{A}(s)$ :

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

## Monte Carlo Policy Gradient (REINFORCE)

Given  $\pi_{\theta}(a|s)$ , initialize  $\theta \in \mathbb{R}^{d'}$   
Repeat forever:  
    Generate an episode following  $\pi$   
    For each step  $t$  to  $T$  in the episode:  
         $G \leftarrow$  return from step  $t$   
         $\theta \leftarrow \theta + \alpha \gamma^t G \nabla_{\theta} \ln \pi_{\theta}(A_t|S_t)$

## Q-learning (off-policy TD control)

Initialize  $Q(s, a)$  arbitrarily,  $Q(\text{terminal-state}, -) = 0$

Repeat (for each episode):  
    Initialize  $S$   
    Repeat (for each step of episode):  
        Choose  $A$  in  $S$  with  $\epsilon$ -greedy policy from  $Q$   
        Take action  $A$ , observe  $R, S'$   
         $Q(S, A) \leftarrow (1 - \alpha)Q(S, A) + \alpha(R + \gamma \max_a Q(S', a))$   
         $S \leftarrow S'$   
    Until  $S$  is terminal

## Sarsa (on-policy TD control)

Initialize  $Q(s, a)$  arbitrarily,  $Q(\text{terminal-state}, -) = 0$

Repeat (for each episode):  
    Initialize  $S$   
    Choose  $A$  in  $S$  with  $\epsilon$ -greedy policy from  $Q$   
    Repeat (for each step of episode):  
        Take action  $A$ , observe  $R, S'$   
        Choose  $A'$  in  $S'$  with  $\epsilon$ -greedy policy from  $Q$   
         $Q(S, A) \leftarrow (1 - \alpha)Q(S, A) + \alpha(R + \gamma Q(S', A'))$   
         $S \leftarrow S'$   
         $A \leftarrow A'$   
    Until  $S$  is terminal

## One-step Actor-Critic

Given  $\pi_{\theta}(a|s)$ ,  $\hat{v}_{\omega}(s)$ , initialize  $\theta \in \mathbb{R}^{d'}$ ,  $\omega \in \mathbb{R}^d$

Repeat forever:  
    Initialize  $S$   
     $I \leftarrow 1$   
    While  $S$  is not terminal:  
         $A \sim \pi_{\theta}(A|S)$   
        Take action  $A$ , observe  $S', R$   
         $\delta \leftarrow R + \gamma \hat{v}_{\omega}(S') - \hat{v}_{\omega}(S)$   
         $\omega \leftarrow \omega + \alpha I \delta \nabla_{\omega} \hat{v}_{\omega}(S)$   
         $\theta \leftarrow \theta + \beta I \delta \nabla_{\theta} \ln \pi_{\theta}(A|S)$   
         $I \leftarrow \gamma I$   
         $S \leftarrow S'$

## Value Iteration

Initialize array  $V$  arbitrarily

Repeat  
     $\Delta \leftarrow 0$   
    For each  $s \in \mathcal{S}$ :  
         $v \leftarrow V(s)$   
         $V(s) \leftarrow \max_a \sum_{s', r} p(s', r|s, a)[r + \gamma V(s')]$   
         $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
    until  $\Delta < \epsilon$

Output a deterministic policy  $\pi \approx \pi_*$ :  
     $\pi(s) = \text{argmax}_a \sum_{s', r} p(s', r|s, a)[r + \gamma V(s')]$

## Policy Iteration

**(1) Initialization**  
 $V(s) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in \mathcal{S}$

**(2) Policy Evaluation**  
Repeat  
     $\Delta \leftarrow 0$   
    For each  $s \in \mathcal{S}$ :  
         $v \leftarrow V(s)$   
         $V(s) \leftarrow \sum_{s', r} p(s', r|s, \pi(s))[r + \gamma V(s')]$   
         $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
    until  $\Delta < \epsilon$

**(3) Policy Improvement**  
 $\text{policy-stable} \leftarrow \text{true}$   
For each  $s \in \mathcal{S}$ :  
     $\text{old-action} \leftarrow \pi(s)$   
     $\pi(s) \leftarrow \text{argmax}_a \sum_{s', r} p(s', r|s, a)[r + \gamma V(s')]$   
    If  $\text{old-action} \neq \pi(s)$ , then  $\text{policy-stable} \leftarrow \text{false}$   
If  $\text{policy-stable}$ , then stop, else go to (2)