

Reinforcement Learning with Tensorflow

In exercise 11, we made use of the quadratic TD error as our loss function:

$$\begin{aligned}\mathcal{L} &= \delta^2 \\ &= \underbrace{(r_{k+1} + \gamma \hat{q}(\mathbf{x}_{k+1}, u_{k+1}, \mathbf{w}))}_{\text{target}} - \underbrace{\hat{q}(\mathbf{x}_k, u_k, \mathbf{w})}_{\text{prediction}})^2.\end{aligned}$$

The semi-gradient of this function is then be given by:

$$\nabla_{\mathbf{w}} \delta^2 = -2 \delta \nabla_{\mathbf{w}} \hat{q}(\mathbf{x}_k, u_k, \mathbf{w}).$$

The true gradient would of course be $\nabla_{\mathbf{w}} \delta^2 = 2 \delta \nabla_{\mathbf{w}} [\gamma \hat{q}(\mathbf{x}_{k+1}, u_{k+1}, \mathbf{w}) - \hat{q}(\mathbf{x}_k, u_k, \mathbf{w})]$. However, Tensorflow will not be able to compute this gradient, unless the prediction $\hat{q}(\mathbf{x}_{k+1}, u_{k+1}, \mathbf{w})$ happens within the same gradient tape as the prediction of $\hat{q}(\mathbf{x}_k, u_k, \mathbf{w})$. If we estimate the target without the use of a gradient tape it will be considered as a constant number, whose derivative will then be zero.

Tensorflow will inherently minimize loss, so it will perform a gradient **descent** with this gradient when the `apply_gradients` function is used:

$$\begin{aligned}\mathbf{w}_{\text{new}} &= \mathbf{w}_{\text{old}} - \alpha \nabla_{\mathbf{w}} \delta^2 \\ &= \mathbf{w}_{\text{old}} + 2 \alpha \delta \nabla_{\mathbf{w}} \hat{q}(\mathbf{x}_k, u_k, \mathbf{w}).\end{aligned}$$

Which is the desired operation for Sarsa(0), when assuming the presence of the factor 2 to be only of minor importance, or by defining $\tilde{\alpha} = 2\alpha$.

Algorithmic Implementation: Semi-Gradient Sarsa

input: a differentiable function $\hat{q} : \mathbb{R}^{\kappa} \times \mathbb{R}^{\zeta} \rightarrow \mathbb{R}$
input: a policy π (only if estimating q_{π})
parameter: step size $\alpha \in \{\mathbb{R} | 0 < \alpha < 1\}$, $\varepsilon \in \{\mathbb{R} | 0 < \varepsilon < 1\}$
init: parameter vector $\mathbf{w} \in \mathbb{R}^{\zeta}$ arbitrarily
for $j = 1, 2, \dots$ *episodes* **do**
 initialize \mathbf{x}_0 ;
 for $k = 0, 1, 2 \dots$ *time steps* **do**
 $u_k \leftarrow$ apply action from $\pi(\mathbf{x}_k)$ or ε -greedy on $\hat{q}(\mathbf{x}_k, \cdot, \mathbf{w})$;
 observe \mathbf{x}_{k+1} and r_{k+1} ;
 if \mathbf{x}_{k+1} *is terminal* **then**
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha [r_{k+1} - \hat{q}(\mathbf{x}_k, u_k, \mathbf{w})] \nabla_{\mathbf{w}} \hat{q}(\mathbf{x}_k, u_k, \mathbf{w})$;
 go to next episode;
 choose u' from $\pi(\mathbf{x}_{k+1})$ or ε -greedy on $\hat{q}(\mathbf{x}_{k+1}, \cdot, \mathbf{w})$;
 $\mathbf{w} \leftarrow$
 $\mathbf{w} + \alpha [r_{k+1} + \gamma \hat{q}(\mathbf{x}_{k+1}, u', \mathbf{w}) - \hat{q}(\mathbf{x}_k, u_k, \mathbf{w})] \nabla_{\mathbf{w}} \hat{q}(\mathbf{x}_k, u_k, \mathbf{w})$;

Algo. 10.2: Semi-gradient Sarsa action-value estimation (output: parameter vector \mathbf{w} for \hat{q}_{π} or \hat{q}^*)

However, for Sarsa(λ), we want a different gradient:

Algorithmic Implementation: Semi-Gradient Sarsa(λ)

```

input: a differentiable function  $\hat{q} : \mathbb{R}^\kappa \times \mathbb{R}^\zeta \rightarrow \mathbb{R}$ 
input: a policy  $\pi$  (only if estimating  $q_\pi$ )
parameter:  $\alpha \in \{\mathbb{R} | 0 < \alpha < 1\}$ ,  $\varepsilon \in \{\mathbb{R} | 0 < \varepsilon < 1\}$ ,  $\lambda \in \{\mathbb{R} | 0 \leq \lambda \leq 1\}$ 
init: parameter vector  $\mathbf{w} \in \mathbb{R}^\zeta$  arbitrarily
for  $j = 1, 2, \dots$  episodes do
    initialize  $\mathbf{x}_0$  and set  $\mathbf{z} = 0$ ;
     $u_0 \leftarrow$  choose action from  $\pi(\mathbf{x}_0)$  or  $\varepsilon$ -greedy on  $\hat{q}(\mathbf{x}_0, \cdot, \mathbf{w})$ ;
    for  $k = 0, 1, 2 \dots$  time steps do
        apply action  $u_k$ , observe  $\mathbf{x}_{k+1}$  and  $r_{k+1}$ ;
        if  $\mathbf{x}_{k+1}$  is terminal then  $\delta \leftarrow r_{k+1} - \hat{q}(\mathbf{x}_k, u_k, \mathbf{w})$ ;
        else
             $u_{k+1} \leftarrow \pi(\mathbf{x}_{k+1})$  or  $\varepsilon$ -greedy on  $\hat{q}(\mathbf{x}_{k+1}, \cdot, \mathbf{w})$ ;
             $\delta \leftarrow r_{k+1} + \gamma \hat{q}(\mathbf{x}_{k+1}, u_{k+1}, \mathbf{w}) - \hat{q}(\mathbf{x}_k, u_k, \mathbf{w})$ ;
         $\mathbf{z} \leftarrow \gamma \lambda \mathbf{z} + \nabla_{\mathbf{w}} \hat{q}(\mathbf{x}_k, u_k, \mathbf{w})$ ;
         $\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta \mathbf{z}$ ;
        exit loop if  $\mathbf{x}_{k+1}$  is terminal;

```

Algo. 11.3: Semi-gradient Sarsa(λ) (output: parameter vector \mathbf{w} for \hat{q}_π or \hat{q}^*)

Oliver Wallscheid

RL Lecture 11

27

In this case, it is easier to directly calculate the derivative of the action-value approximator:

$$\mathcal{L} = \hat{q}(\mathbf{x}_k, u_k, \mathbf{w}).$$

This way, we receive $\nabla_{\mathbf{w}} \hat{q}(\mathbf{x}_k, u_k, \mathbf{w})$ without any scaling factors.