NAME: KARMA TARAP
CSCI E-89c Deep Reinforcement Learning
Part I of Assignment 11

Suppose each state $s \in \mathcal{S}$ of the Markov Decision Process can be represented by a vector of 2 real-valued features: $\mathbf{x}(s) = (x_1(s), x_2(s))^T$.

Given some policy $\pi$, suppose we model the state value function $v_\pi(s)$ with a *fully connected feedforward neural network* (please see the table below) which has two inputs ($x_1(s)$ and $x_2(s)$), one hidden layer that consists of two neurons ($u_1$ and $u_2$) with ReLU activation functions, and one output ($\hat{v}(s, \mathbf{w})$) with the ReLU activation function.

The explicit representation of this network is

| input layer | hidden layer | output layer |
|---|---|---|
| $x_1$ | $u_1 = f(w_{01}^{(1)} + w_{11}^{(1)} x_1 + w_{21}^{(1)} x_2)$ | $\hat{v} = f(w_0^{(2)} + w_1^{(2)} u_1 + w_2^{(2)} u_2)$ |
| $x_2$ | $u_2 = f(w_{02}^{(1)} + w_{12}^{(1)} x_1 + w_{22}^{(1)} x_2)$ | |

Here, $f(x)$ denotes the rectified linear unit (ReLU) defined as follows:

$$f(x) = \begin{cases} x, & \text{if } x \geq 0, \\ 0, & \text{if } x < 0. \end{cases}$$

Assume that the weights,

$$\mathbf{w} = \big( \underbrace{w_{01}^{(1)}, w_{11}^{(1)}, w_{21}^{(1)}, w_{02}^{(1)}, w_{12}^{(1)}, w_{22}^{(1)}}_{\text{hidden layer}}, \underbrace{w_0^{(2)}, w_1^{(2)}, w_2^{(2)}}_{\text{output layer}} \big)^T,$$

are currently estimated as follows:

| hidden layer | output layer |
|---|---|
| $w_{01}^{(1)} = -1.2, w_{11}^{(1)} = 0.1, w_{21}^{(1)} = 0.5$ | $w_0^{(2)} = 0.2, w_1^{(2)} = -0.8, w_2^{(2)} = 1.2$ |
| $w_{02}^{(1)} = 0.9, w_{12}^{(1)} = 0.8, w_{22}^{(1)} = -0.3$ | |

Assume the agent minimizes the mean squared error loss function,

$$L \doteq \frac{1}{2} \left( \hat{v}(S_t, \mathbf{w}) - v_\pi(S_t) \right)^2,$$

using Stochastic Gradient Descent (SGD), i.e. the Neural Network is trained in mini-batches of size 1.

If for current state $S_t$, the features are $x_1(S_t) = 1.3$ and $x_2(S_t) = 0.7$; and the agent "observes" $v_\pi(S_t)$ (this, of course, means the agent uses MC return, 1-step TD return, etc. as a "measurement" of $v_\pi(S_t)$) to be 4.1, please find

(a) Error associated with the output layer:

$$\varepsilon^{(2)} \doteq \frac{\partial L}{\partial \hat{v}}.$$

(b) Errors associated with the hidden layer:

$$\varepsilon_h^{(1)} \doteq \frac{\partial L}{\partial u_h}, \quad h = 1, 2.$$

(c) Partial derivatives of the loss function with respect to weights in the output layer:

$$\frac{\partial L}{\partial w_h^{(2)}}, \quad h = 0, 1, 2.$$

(d) Partial derivatives of the loss function with respect to weights in the hidden layer:

$$\frac{\partial L}{\partial w_{jh}^{(1)}}, \quad j = 0, 1, 2 \text{ and } h = 1, 2.$$

(e) The next SGD update of the weights using $\alpha = 0.1$:

$$\mathbf{w} - \alpha \nabla L,$$

where $\nabla L \doteq \Big( \underbrace{\frac{\partial L}{\partial w_{01}^{(1)}}, \frac{\partial L}{\partial w_{11}^{(1)}}, \frac{\partial L}{\partial w_{21}^{(1)}}, \frac{\partial L}{\partial w_{02}^{(1)}}, \frac{\partial L}{\partial w_{12}^{(1)}}, \frac{\partial L}{\partial w_{22}^{(1)}}}_{\text{hidden layer}}, \underbrace{\frac{\partial L}{\partial w_0^{(2)}}, \frac{\partial L}{\partial w_1^{(2)}}, \frac{\partial L}{\partial w_2^{(2)}}}_{\text{output layer}} \Big)^T.$

Please notice that the "measurement" of the state-value $v_\pi(S_t)$ here is considered to be independent of $\mathbf{w}$ (please see, for example, the Semi-gradient 1-step Temporal-Difference (TD) prediction).

SOLUTION:

a)

$$\varepsilon^{(2)} \doteq \frac{\partial L}{\partial \hat{v}}$$

$$= (\hat{v}(S_t, \mathbf{w}) - v_\pi(S_t))$$

$$= f(w_0^{(2)} + w_1^{(2)} u_1 + w_2^{(2)} u_2) - 4.1$$

$$= f(w_0^{(2)} + w_1^{(2)} f(w_{01}^{(1)} + w_{11}^{(1)} x_1 + w_{21}^{(1)} x_2) + w_2^{(2)} f(w_{01}^{(1)} + w_{11}^{(1)} x_1 + w_{21}^{(1)} x_2)) - 4.1$$

$$= f(w_0^{(2)} + w_1^{(2)}(0) + w_2^{(2)}(1.73)) - 4.1$$

$$= 2.276 - 4.1$$

$$= -1.824$$

b)

$$\varepsilon_1^{(1)} \doteq \frac{\partial L}{\partial u_1} = \frac{\partial L}{\partial \hat{v}}\frac{\partial \hat{v}}{\partial u_1} = \varepsilon^{(2)}u_1 w_1^{(2)} = 0$$

$$\varepsilon_2^{(1)} \doteq \frac{\partial L}{\partial u_2} = \frac{\partial L}{\partial \hat{v}}\frac{\partial \hat{v}}{\partial u_2} = \varepsilon^{(2)}u_2 w_2^{(2)} = -3.787$$

c)

$$\frac{\partial L}{\partial w_0^{(2)}} = \frac{\partial \hat{v}}{\partial w_0^2}\frac{\partial L}{\partial \hat{v}} = \epsilon^{(2)}\hat{v}(1) = -4.15$$

$$\frac{\partial L}{\partial w_1^{(2)}} = \frac{\partial \hat{v}}{\partial w_1^2}\frac{\partial L}{\partial \hat{v}} = \epsilon^{(2)}\hat{v}(u_1) = 0$$

$$\frac{\partial L}{\partial w_2^{(2)}} = \frac{\partial \hat{v}}{\partial w_2^2}\frac{\partial L}{\partial \hat{v}} = \epsilon^{(2)}\hat{v}(u_2) = -7.18$$

d)

$$\frac{\partial L}{\partial w_{01}^{(1)}} = \frac{\partial u_1}{\partial w_{01}^{(1)}}\frac{\partial \hat{v}}{\partial u_1}\frac{\partial L}{\partial \hat{v}} = (1)\epsilon_1^{(1)}\epsilon^{(2)} = 0$$

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \frac{\partial u_1}{\partial w_{11}^{(1)}}\frac{\partial \hat{v}}{\partial u_1}\frac{\partial L}{\partial \hat{v}} = (x_1)\epsilon_1^{(1)}\epsilon^{(2)} = 0$$

$$\frac{\partial L}{\partial w_{21}^{(1)}} = \frac{\partial u_1}{\partial w_{21}^{(1)}}\frac{\partial \hat{v}}{\partial u_1}\frac{\partial L}{\partial \hat{v}} = (x_2)\epsilon_1^{(1)}\epsilon^{(2)} = 0$$

$$\frac{\partial L}{\partial w_{02}^{(1)}} = \frac{\partial u_1}{\partial w_{02}^{(1)}}\frac{\partial \hat{v}}{\partial u_2}\frac{\partial L}{\partial \hat{v}} = (x_2)\epsilon_2^{(1)}\epsilon^{(2)} = 6.9$$

$$\frac{\partial L}{\partial w_{12}^{(1)}} = \frac{\partial u_1}{\partial w_{12}^{(1)}}\frac{\partial \hat{v}}{\partial u_2}\frac{\partial L}{\partial \hat{v}} = (x_2)\epsilon_2^{(1)}\epsilon^{(2)} = 8.98$$

$$\frac{\partial L}{\partial w_{22}^{(1)}} = \frac{\partial u_1}{\partial w_{22}^{(1)}}\frac{\partial \hat{v}}{\partial u_2}\frac{\partial L}{\partial \hat{v}} = (x_2)\epsilon_2^{(1)}\epsilon^{(2)} = 4.8$$

e)

$$= (-1.2, 0.1, 0.5, 0.21, -0.098, -0.783, 0.615, -0.8, 2.018)^T$$