

Name: Alfred  
Student No.: 2020-12788

## Exercise 6

### Data Science & Reinforcement Learning

Spring 2023

1. Give pseudocode for semi-gradient one-step Expected Sarsa for control. You can build on the semi-gradient Sarsa code for this question.

**Episodic Semi-gradient Sarsa for Estimating  $\hat{q} \approx q_*$**

**Input:** a differentiable action-value function parameterization  $\hat{q}: \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$ , a policy  $\pi$

**Algorithm parameters:** step size  $\alpha > 0$ , small  $\epsilon > 0$

**Initialize** value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

**Loop for each episode:**

$S, A \leftarrow$  initial state and action of episode (e.g.,  $\epsilon$ -greedy)

**Loop for each step of episode:**

Take action  $A$ , observe  $R, S'$

If  $S'$  is terminal:

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$

Go to next episode

Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \mathbf{w})$  (e.g.,  $\epsilon$ -greedy)

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$

~~$S \leftarrow S'$~~

~~$A \leftarrow A'$~~

all again

for  $(S, A)$  pairs from ep.

Loop for every  $S, A$ :

take action  $A$ , observe  $R, S'$

If  $S'$  terminal:

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \sum_a \pi(a|S', \mathbf{w}) \hat{q}(S', a, \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$

goto next  $(S, A)$  pair

2. Show that tabular methods studied earlier in the course are a special case of linear function approximation. What would the feature vectors be?

$$v(\mathbf{s}_k, \vec{\mathbf{w}}) = \vec{\mathbf{x}}(\mathbf{s}_k) \cdot \vec{\mathbf{w}} = \sum_{i=1}^n x_i(\mathbf{s}_k) w_i$$

$$= \sum_{i=1}^n \delta_{ik} w_i = w_k$$

$$\vec{\mathbf{x}}(\mathbf{s}_k) \sim [0, \dots, 1, \dots, 0]$$

↑  
k-th position, only

$$\vec{\mathbf{x}}(\mathbf{s}, \mathbf{w}) = \mathbf{w}^T \mathbf{x}(\mathbf{s}) = \sum_{i=1}^d w_i x_i(\mathbf{s}) \leftarrow \mathbf{w}_n + 1 = \mathbf{w}_n$$

$(S, A)$  vectors are real feature vectors one-hot encoding, defining  $\mathbf{w}$  as vectors of same length.

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

becomes

$$\hat{q}(S, A) \leftarrow \hat{q}(S, A) + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})]$$

$\hat{q}(S, A, \mathbf{w})$   $\nabla \mathbf{w} = 1$

special case of linear function approx

3. For control with function approximation, we did not explicitly consider or give pseudocode for any Monte Carlo methods. Why is it reasonable to NOT consider Monte Carlo methods for function approximation?

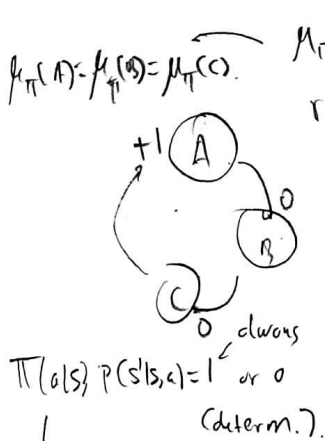
- Because it's too expensive computationally since the learning should go to terminal for any random policy. TD or SARSA would converge online faster.

- MC is variation/type of general TD's, won't be that different than

n-step SARSA but  $n=T$  (terminal).

Since MC updates value function as err. at ep. start - value but TD does from current value-estimation.

4. Consider a Markov reward process consisting of a ring of three states A, B, and C, with state transitions TD has loss going deterministically around the ring. A reward of +1 is received upon arrival in A and otherwise as otherwise as otherwise the reward is 0. What are the differential values of the three states, as defined below?



$$V_\pi(s) := \lim_{\gamma \rightarrow 1} \lim_{h \rightarrow \infty} \sum_{t=0}^h \gamma^t (\mathbb{E}_\pi[R_{t+1} | S_0 = s] - r(\pi))$$

$$\sum_{t=0}^{\infty} \gamma^t [\mathbb{E}_\pi[R_{t+1} | A] - 0] = 0 + \gamma + \gamma^2 + \dots = \frac{\gamma}{1-\gamma}$$

$$V_\pi(s=A) = \lim_{\gamma \rightarrow 1} \left( \frac{1}{1-\gamma} \right) \left( -\frac{1}{3} - \frac{1}{3}\gamma + \frac{2}{3}\gamma^2 \right) = -\frac{1}{3}$$

$$V_\pi(s=B) = \lim_{\gamma \rightarrow 1} \left( \frac{1}{1-\gamma} \right) \left( -\frac{1}{3} + \frac{2}{3}\gamma - \frac{1}{3}\gamma^2 \right) = 0$$

$$V_\pi(s=C) = \lim_{\gamma \rightarrow 1} \left( \frac{1}{1-\gamma} \right) \left( \frac{2}{3} - \frac{1}{3}\gamma - \frac{1}{3}\gamma^2 \right) = \frac{1}{3}$$

5. How would you use optimistic initial values for Sarsa with linear function approximation in episodic MDP with the episode length  $H$ ? Assume that your value function is parametrized with weights  $w \in \mathbb{R}^d$ , all the feature vectors  $\phi(s, a) \in \mathbb{R}^d$  are upper-bounded in L2 norm, i.e.,  $\|\phi(s, a)\|_2 \leq 1$  for all  $(s, a)$ , and the all rewards are bounded above by 1, i.e.,  $R(s, a) \leq 1$  for all  $(s, a)$ .

$$V_\pi(s) = \mu_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{a' \in \mathcal{A}} P(s', a'|s, a) V_\pi(s')$$

$$w = \begin{bmatrix} 1 \\ \vdots \end{bmatrix}_{d \times 1} \quad \phi(s, a) = \begin{bmatrix} 1 \\ \vdots \end{bmatrix}_{d \times 1}$$

$$\text{Target: } R(s, a) + \gamma w^T \phi(s', a')$$

- Initialize value-function weights  $w \in \mathbb{R}^d$  to optimistic values such as 1. to explore more initially.

- if other values, all of them should be less than

$$1 + r + \dots + r^{H-1} \leq \frac{1-r^H}{1-r} \text{ Since } G_{t:t+T} \text{ won't exceed } 1 + r + \dots + r^{H-1}$$

$$= \frac{1-r^H}{1-r} \text{ that's a consideration too}$$

2 of 2

But if n-step SARSA  $\rightarrow \frac{1-r^n}{1-r} + r^n \|w\|_{\text{initial}}$

can be used too.