

**Problem 1**

[30 points] Consider the following stochastic approximation with a fixed step size  $\epsilon \in (0, 1)$

$$\theta_{k+1} = (1 - \epsilon)\theta_k + \epsilon X_k,$$

where  $X_k$  are i.i.d random variables with mean  $\mu$  and variance  $\sigma^2$ . In addition, given any constant  $a > 0$ , the Chebyshev's inequality implies

$$\mathbb{P}(|X - E(X)| \geq a) \leq \frac{\text{Var}(X)}{a^2},$$

where  $X$  is a random variable with  $E(X)$  and  $\text{Var}(X)$  are its expected value and variance, respectively. Show that for any  $c > 0$

$$\limsup_{k \rightarrow \infty} \mathbb{P}(|\theta_k - \mu| \geq c\sqrt{\epsilon}) \leq \frac{\sigma^2}{c^2}.$$

Step 1: Showing that  $\lim_{k \rightarrow \infty} E(\theta_k) = \mu$

$$\theta_{k+1} = (1 - \epsilon)\theta_k + \epsilon X_k$$

$$\therefore \theta_1 = (1 - \epsilon)\theta_0 + \epsilon X_0$$

$$\begin{aligned} \theta_2 &= (1 - \epsilon)((1 - \epsilon)\theta_0 + \epsilon X_0) + \epsilon X_1 \\ &= (1 - \epsilon)^2 \theta_0 + (1 - \epsilon)\epsilon X_0 + \epsilon X_1 \end{aligned}$$

$$\theta_3 = (1 - \epsilon)^3 \theta_0 + (1 - \epsilon)^2 \epsilon X_0 + (1 - \epsilon)\epsilon X_1 + \epsilon X_2$$

$$\vdots$$

$$\theta_k = (1 - \epsilon)^k \theta_0 + \sum_{j=0}^{k-1} (1 - \epsilon)^j \epsilon X_{k-j}$$

$$E(\theta_k) = E\left((1 - \epsilon)^k \theta_0 + \sum_{j=0}^{k-1} (1 - \epsilon)^j \epsilon X_{k-j}\right)$$

$$= E((1 - \epsilon)^k \theta_0) + E\left(\sum_{j=0}^{k-1} (1 - \epsilon)^j \epsilon X_{k-j}\right)$$

$$\begin{aligned} \lim_{k \rightarrow \infty} E(\theta_k) &= \lim_{k \rightarrow \infty} E((1 - \epsilon)^k \theta_0) + \lim_{k \rightarrow \infty} E\left(\sum_{j=0}^{k-1} (1 - \epsilon)^j \epsilon X_{k-j}\right) \\ &= 0 + \lim_{k \rightarrow \infty} \left[ \sum_{j=0}^{k-1} E((1 - \epsilon)^j \epsilon X_{k-j}) \right] \end{aligned}$$

$$= \lim_{k \rightarrow \infty} \left[ \sum_{j=0}^{k-1} E(\epsilon(1 - \epsilon)^j) E(X_{k-j}) \right]$$

$$= \lim_{k \rightarrow \infty} \left[ \sum_{j=0}^k \varepsilon (1-\varepsilon)^j \mu \right]$$

as  $E(X_k) = \mu$   
given in question

$$= \varepsilon \mu \sum_{j=0}^{\infty} (1-\varepsilon)^j$$

$$= \frac{\varepsilon \mu}{1 - (1-\varepsilon)} = \frac{\varepsilon \mu}{\varepsilon} = \mu$$

$$\therefore \lim_{k \rightarrow \infty} E(\Theta_k) = \mu$$


---

Step 2: Solving  $\lim_{k \rightarrow \infty} \text{Var}(\Theta_k)$

---

$$\Theta_k = (1-\varepsilon)^k \Theta_0 + \sum_{j=0}^k (1-\varepsilon)^j \varepsilon X_{(k-j)} \quad (\text{from step 1 above})$$

$$= (1-\varepsilon)^k \varepsilon X_0 + (1-\varepsilon)^{k-1} \varepsilon X_1 + \dots + (1-\varepsilon) \varepsilon X_{k-1} + \varepsilon X_k$$

$$\text{as } X \text{ is iid: } \text{Var}(aX_1 + bX_2) = a^2 \text{Var}(X_1) + b^2 \text{Var}(X_2)$$

$$\therefore \text{Var}(\Theta_k) = \text{Var}[(1-\varepsilon)^k \varepsilon X_0] + \text{Var}[(1-\varepsilon)^{k-1} \varepsilon X_1] \\ \dots + \text{Var}[(1-\varepsilon) \varepsilon X_{k-1}] + \text{Var}[\varepsilon X_k]$$

$$= (\varepsilon (1-\varepsilon)^k)^2 \text{Var}(X_0) \\ + (\varepsilon (1-\varepsilon)^{k-1})^2 \text{Var}(X_1) \\ + \dots \\ + (\varepsilon (1-\varepsilon)^1)^2 \text{Var}(X_{k-1}) \\ + (\varepsilon (1-\varepsilon)^0)^2 \text{Var}(X_k)$$

$$= \sigma^2 \varepsilon^2 [(1-\varepsilon)^0 + (1-\varepsilon)^2 + \dots + (1-\varepsilon)^{2k-2} + (1-\varepsilon)^{2k}]$$

as  $\text{Var}(X_k) = \sigma^2$   
given in question

$$\text{let } S = \lim_{k \rightarrow \infty} \left[ (1-\varepsilon)^0 + (1-\varepsilon)^2 + \dots + (1-\varepsilon)^{2k-2} + (1-\varepsilon)^{2k} \right]$$

$$S = 1 + (1-\varepsilon)^2 S$$

$$\Rightarrow S - S(1-\varepsilon)^2 = 1$$

$$S(1 - (1-\varepsilon)^2) = 1$$

$$\Rightarrow S = \frac{1}{1 - (1-\varepsilon)^2} = \frac{1}{2\varepsilon - \varepsilon^2}$$

$$\therefore \lim_{k \rightarrow \infty} \text{Var}(\Theta_k) = S \lim_{k \rightarrow \infty} \sigma^2 \varepsilon^2 = \frac{\sigma^2 \varepsilon^2}{2\varepsilon - \varepsilon^2}$$

Step 3: Chebyshev's equation

$$P(|x - E(x)| \geq a) \leq \frac{\text{Var}(x)}{a^2}$$

$$\therefore \lim_{k \rightarrow \infty} P(|\Theta_k - E(\Theta_k)| \geq c\sqrt{\varepsilon})$$

$$= \lim_{k \rightarrow \infty} P(|\Theta_k - \mu| \geq c\sqrt{\varepsilon})$$

$$\leq \lim_{k \rightarrow \infty} \frac{\text{Var}(\Theta_k)}{c^2 \varepsilon}$$

$$= \frac{\sigma^2 \varepsilon^2}{c^2 \varepsilon (2\varepsilon - \varepsilon^2)}$$

$$= \frac{\sigma^2}{c^2} \left( \frac{\varepsilon}{2\varepsilon - \varepsilon^2} \right) = \frac{\sigma^2}{c^2(2-\varepsilon)}$$

$$\leq \frac{\sigma^2}{c^2}$$

$$\therefore \limsup_{k \rightarrow \infty} P(|\Theta_k - \mu| \geq c\sqrt{\varepsilon}) \leq \frac{\sigma^2}{c^2}$$

$$\left| \begin{array}{l} a) \lim_{k \rightarrow \infty} E(\Theta_k) = \mu \\ a) \text{ shown in step 1} \end{array} \right.$$

$$\left| \begin{array}{l} a) \lim_{k \rightarrow \infty} \text{Var}(\Theta_k) = \frac{\sigma^2 \varepsilon^2}{2\varepsilon - \varepsilon^2} \\ a) \text{ shown in step 2} \end{array} \right.$$

$$\left| \begin{array}{l} a) \frac{1}{2-\varepsilon} \in \left(\frac{1}{2}, 1\right) \\ \text{since } \varepsilon \in (0, 1) \end{array} \right.$$

We consider a discounted MDP problem with finite state space  $\mathcal{S}$  and finite action space  $\mathcal{A}$ . For any stationary policy  $\mu$  define the value function  $V_\mu$

and let  $V_\mu(s) \leq V^*(s) = V_{\mu^*}(s)$  for the optimal policy  $\mu^*$  and for all  $s \in \mathcal{S}$ . Given any value function  $V_\mu$  we denote by  $Q_\mu$  the state-action value function

Similarly, the optimal state-action value function  $Q^*$  is defined as

Given a stationary policy  $\mu$ , defined the Bellman operator  $T_\mu$  as

**Questions:** Let  $\mu$  be a stationary policy and any state action-value functions  $Q, Q'$ .

- 1)  $\|V_{Q,\mu} - V_{Q',\mu}\|_\infty \leq \|Q - Q'\|_\infty$
- 2)  $\|V_Q - V_{Q'}\|_\infty \leq \|Q - Q'\|_\infty$

- given any  $s \in S$ :

$\mu(s)$  subset of all possible actions, a

$$\therefore \|U_{Q,n} - U_{Q',n}\|_{\infty} \leq \|Q - Q'\|_{\infty}$$

---

## Part 2

note: given  $x = \{x_1, x_2, \dots, x_n\}$ , w.l.o.g:  $\max_i x_i \geq \max_j y_j$   
 $y = \{y_1, y_2, \dots, y_n\}$ .

$$\begin{aligned} \Rightarrow \quad | \max_i x_i - \max_j y_j | &= \max_i x_i - \max_j y_j \\ &= x_{i_2} - \max_j y_j \quad \left| \begin{array}{l} \text{letting} \\ i_2 = \arg \max x_i \end{array} \right. \\ &\leq x_{i_2} - y_{i_2} \\ &= |x_{i_2} - y_{i_2}| \\ &\leq \max_i |x_i - y_i| \quad (\text{eq 1}) \end{aligned}$$

$$\therefore | \max_i x_i - \max_j y_j | \leq \max_i |x_i - y_i|$$

given any  $s \in S$ :

$$\begin{aligned} |V_Q(s) - V_{Q'}(s)| &= \left| \max_{a_1} Q(s, a_1) - \max_{a_2} Q'(s, a_2) \right| \quad (a_1, a_2 \in A) \\ &\leq \max_a |Q(s, a) - Q'(s, a)| \quad \left| \begin{array}{l} \text{using} \\ \text{eq (1) from} \\ \text{above.} \end{array} \right. \\ &= \|Q(s, a) - Q'(s, a)\|_\infty \end{aligned}$$

$$\therefore \underline{\underline{\|V_Q - V_{Q'}\|_\infty \leq \|Q - Q'\|_\infty}}$$

2. Show that [10 points]

$$1) \|T_\mu Q - T_\mu Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

$$2) \|TQ - TQ'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

Part 1

given any  $(s, a)$ :

$$\begin{aligned} (T_\mu Q - T_\mu Q')(s, a) &= \gamma \sum_{s' \in S} P_{ss'}(a) [Q(s', u) - Q'(s', u)] \Big|_{\substack{u \in A \\ u \in A}} \\ &\leq \gamma \sum_{s' \in S} P_{ss'}(a) |Q(s', u) - Q'(s', u)| \\ &\leq \gamma \sum_{s' \in S} P_{ss'}(a) \max_{u, u'} |Q(s', u) - Q'(s', u)| \\ &\leq \gamma \sum_{s' \in S} P_{ss'}(a) \max_v |Q(s', v) - Q'(s', v)| \\ &\leq \gamma \sum_{s' \in S} P_{ss'}(a) \max_{v, y} |Q(y, v) - Q'(y, v)| \Big|_{y \in S} \\ &= \gamma \max_{v, y} |Q(y, v) - Q'(y, v)| \\ &= \gamma \|Q - Q'\|_\infty \end{aligned}$$

$$\therefore \|T_\mu Q - T_\mu Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

## Part 2

given any  $(s, a)$ :

$$TQ(s, a) - TQ'(s, a) = \gamma \sum_{s' \in S} P_{ss'}(a) \left[ \max_v Q(s', v) - \max_u Q'(s', u) \right] \quad \left| \begin{array}{l} v \in A \\ u \in A \end{array} \right.$$

$$\Rightarrow |TQ(s, a) - TQ'(s, a)| = \left| \gamma \sum_{s' \in S} P_{ss'}(a) \left[ \max_v Q(s', v) - \max_u Q'(s', u) \right] \right|$$
$$\leq \gamma \sum_{s' \in S} P_{ss'}(a) \left| \left[ \max_v Q(s', v) - \max_u Q'(s', u) \right] \right| \quad \left| \begin{array}{l} \text{via triangle} \\ \text{inequality.} \end{array} \right.$$

note: given  $x = \{x_1, x_2, \dots, x_n\}$   
 $y = \{y_1, y_2, \dots, y_n\}$ , w.l.o.g:  $\max_i x_i \geq \max_j y_j$

$$\Rightarrow \left| \max_i x_i - \max_j y_j \right| = \max_i x_i - \max_j y_j$$
$$= x_{i_2} - \max_j y_j \quad \left| \begin{array}{l} \text{letting} \\ i_2 = \arg \max x_i \end{array} \right.$$
$$\leq x_{i_2} - y_{i_2}$$
$$= |x_{i_2} - y_{i_2}|$$
$$\leq \max_i |x_i - y_i| \quad (\text{eq 1})$$

• applying eq 1 from above:

$$\gamma \sum_{s' \in S} P_{ss'}(a) \left| \left[ \max_v Q(s', v) - \max_u Q'(s', u) \right] \right|$$
$$\leq \gamma \sum_{s' \in S} P_{ss'}(a) \max_v |Q(s', v) - Q'(s', v)| \quad \left| \begin{array}{l} \text{via (eq 1)} \\ v \in A \end{array} \right.$$
$$\leq \gamma \sum_{s' \in S} P_{ss'}(a) \max_{v, g} |Q(g, v) - Q'(g, v)| \quad \left| \begin{array}{l} g \in S \end{array} \right.$$
$$= \gamma \max_{v, g} |Q(g, v) - Q'(g, v)|$$
$$= \gamma \|Q - Q'\|_\infty$$

$$\therefore \|TQ - TQ'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

3. Show that [10 points]

$$i) \quad \|Q - Q_\mu\|_\infty \leq \frac{\|Q - T_\mu Q\|_\infty}{1 - \gamma}$$

$$2) \quad \|Q - Q^*\|_\infty \leq \frac{\|Q - TQ\|_\infty}{1 - \gamma}$$

## Part 1

$$Q(s,a) - T_M Q(s,a) \leq \|Q - T_M Q\|_\infty$$

let  $\|Q - T_n Q\|_\infty = r$ .

$$\Rightarrow T_M Q \leq Q + r1$$

$$\Rightarrow T_M^2 Q = T_M(T_M Q) \leq T_M(Q + r1) = T_M Q + \gamma r 1$$

$$\Rightarrow T_M^3 Q \leq T_M^2 Q + \gamma^2 r 1$$

$$\Rightarrow T_M^L Q \leq T_M^{L-1} Q + \gamma^{L-1} r_1$$

$$\Rightarrow T_M^L Q - T_M^{L^{-1}} Q \leq \gamma^{L^{-1}} r 1$$

note that:  $T_n^m Q - Q = \sum_{L=1}^m (T_n^L Q - T_n^{L-1} Q)$

letting  $n \rightarrow \infty$ :

$$\bullet \lim_{m \rightarrow \infty} (T_m^m Q - Q) = Q_m - Q$$

$$\lim_{n \rightarrow \infty} \left( \sum_{k=1}^n \gamma^{k-1} \cdot 1 \right) = \frac{r1}{1-\gamma}$$

$$\therefore \lim_{m \rightarrow \infty} (T_m^n Q - Q) \leq \lim_{m \rightarrow \infty} \left( \sum_{L=1}^m \gamma^{L-1} \cdot 1 \right)$$

$$\Rightarrow Q_n - Q \leq \frac{r1}{1-\gamma}$$

$$\Rightarrow Q_n - Q \leq \frac{\|Q - T_n Q\|_\infty}{1 - \delta}$$

$$\Rightarrow \|Q_n - Q\|_\infty \leq \frac{\|Q - T_n Q\|_\infty}{1 - \gamma}$$

---



## Part 2

$$Q(s,a) - TQ(s,a) \leq \|Q - TQ\|_{\infty}$$

$$\text{let } \|Q - TQ\|_{\infty} = r.$$

$$\Rightarrow TQ \leq Q + r1$$

$$\Rightarrow T^2 Q = T(TQ) \leq T(Q + r1) = TQ + \gamma r1$$

$$\Rightarrow T^3 Q \leq T^2 Q + \gamma^2 r1$$

$$\vdots$$
$$\Rightarrow T^L Q \leq T^{L-1} Q + \gamma^{L-1} r1$$

$$\Rightarrow T^L Q - T^{L-1} Q \leq \gamma^{L-1} r1$$

$$\text{note that: } T^m Q - Q = \sum_{L=1}^m (T^L Q - T^{L-1} Q)$$
$$\leq \sum_{L=1}^m \gamma^{L-1} r1$$

letting  $m \rightarrow \infty$ :

$$\bullet \lim_{m \rightarrow \infty} (T^m Q - Q) = Q^* - Q$$

$$\bullet \lim_{m \rightarrow \infty} \left( \sum_{L=1}^m \gamma^{L-1} r1 \right) = \frac{r1}{1-\gamma}$$

$$\therefore \lim_{m \rightarrow \infty} (T^m Q - Q) \leq \lim_{m \rightarrow \infty} \left( \sum_{L=1}^m \gamma^{L-1} r1 \right)$$

$$\Rightarrow Q^* - Q \leq \frac{r1}{1-\gamma}$$

$$\Rightarrow Q^* - Q \leq \frac{\|Q - TQ\|_{\infty} 1}{1-\gamma}$$

$$\Rightarrow \|Q - Q^*\|_{\infty} \leq \frac{\|Q - TQ\|_{\infty}}{1-\gamma}$$

///

4. Let  $\mu$  be the greedy policy for any state-action value function  $Q$ , i.e.,

$$\mu(s) = \arg \max_a Q(s, a)$$

Define the Bellman error for  $Q$  as  $\beta = \|TQ - Q\|_\infty$ . Let  $V_\mu$  be the value function associated with the greedy policy  $\mu$ . Show that [10 points]

$$V_\mu(s) \geq V^*(s) - \frac{2\beta}{1-\gamma}, \quad \forall s \in \mathcal{S}.$$

• from Q2.3:  $\|Q - Q_\mu\|_\infty \leq \frac{\|Q - TQ\|_\infty}{1-\gamma}$  (eq 1)

$$\|Q - Q^*\|_\infty \leq \frac{\|Q - TQ\|_\infty}{1-\gamma} \quad (\text{eq 2})$$

$$\|Q^* - Q_\mu\|_\infty = \|(Q^* - Q) + (Q - Q_\mu)\|_\infty$$

$$\leq \|Q^* - Q\|_\infty + \|Q - Q_\mu\|_\infty$$

| via triangle inequality.

$$= \|Q - Q^*\|_\infty + \|Q - Q_\mu\|_\infty$$

$$\leq \frac{\|Q - TQ\|_\infty}{1-\gamma} + \frac{\|Q - TQ\|_\infty}{1-\gamma}$$

| using eq (1) and eq (2) above

$$= \frac{2\|Q - TQ\|_\infty}{1-\gamma}$$

$$\therefore \|Q^* - Q_\mu\|_\infty \leq \frac{2\|Q - TQ\|_\infty}{1-\gamma}$$

$$\Rightarrow |Q^*(s, \pi^*(s)) - Q(s, \mu(s))| \leq \frac{2\|Q - TQ\|_\infty}{1-\gamma} \quad \left| \begin{array}{l} \text{for a given} \\ s \in \mathcal{S} \end{array} \right.$$

$$\Rightarrow V^*(s) - V_\mu(s) \leq \frac{2\|Q - TQ\|_\infty}{1-\gamma}$$

$$\Rightarrow V_\mu(s) \geq V^*(s) - \frac{2\|Q - TQ\|_\infty}{1-\gamma}, \quad \forall s \in \mathcal{S}$$

=====

### Problem 3 - Coding question [30 points]

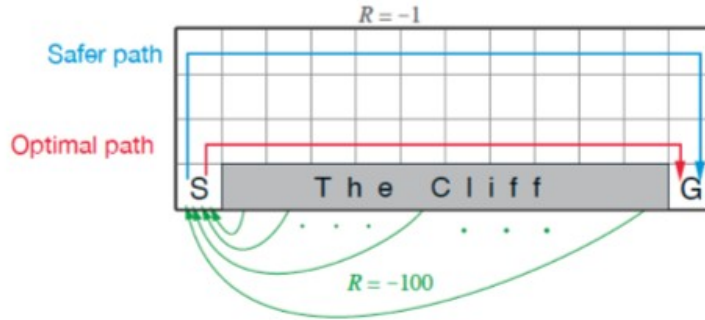


Figure 1: Cliff-walking or gridworld problem (Example 6.6 in Sutton and Barto's book)

Consider the gridworld shown in Fig. 1. This is a standard undiscounted ( $\gamma = 1$ ), episodic task, with start ( $S$ ) and goal ( $G$ ) states, and the usual actions causing movement up, down, right, and left. Reward is  $-1$  on all transitions except those into the region marked “The Cliff.” Stepping into this region incurs a reward of  $-100$  and sends the agent instantly back to the start. There are 48 states (positions in the grid) and 4 actions. This environment can be found here <https://github.com/caburu/gym-cliffwalking>.

Let  $\mu$  be a fixed stochastic policy, which assigns uniform distribution on  $\mathcal{A}$ , i.e., given any state  $i$ , we have the probability of taking action  $a$  is  $\mu(a|i) = 1/4$  for all  $a \in \mathcal{A}$ . Your job is to implement  $\text{TD}(\lambda)$  for finding  $V_\mu$ . You can either implement  $\text{TD}(\lambda)$  with tabular settings or with linear function approximations [1]. Note that this is an undiscounted and episodic problem.

[1.] Value Function Approximation in Reinforcement Learning using the Fourier Basis. George Konidaris and Sarah Osentoski and Philip Thomas.

**Questions:** Write a simulation program to compute  $V_\mu$  for the cliff-walking problem using  $\text{TD}(\lambda)$ . In your simulation, consider 1000 episodes, where each episode runs 20 steps of  $\text{TD}(\lambda)$  given in class. For each episode, compute the norm of the expected TD update (NEU), the average of temporal difference, i.e.,

$$\text{NEU} = \frac{1}{20} \sum_{k=1}^{20} (d_k z_k)^2,$$

where  $z_k$  is the trace vector. Then for every 10 episode, you take the average of the NEU values and plot this average as a function of the number of episodes. Note that for each episode, you should initialize your function values  $V_\mu$  as the values returned by the previous step.

You are asked to submit a pseudo code to explain your simulation and a plot which shows 5 curves of the average of NEU values as a function of the number of episodes for  $\lambda = 0, .3, .5, .7, 1$ . Finally, briefly explain the impacts of  $\lambda$  on the performance of TD learning.

## Pseudo Code:

### On-line Tabular TD( $\lambda$ )

Input: a policy  $\pi$ , trace decay rate  $\lambda \in [0, 1]$

Output: NEU\_episodes (List of NEU values corresponding to each episode)

Algorithm parameters: step size  $\alpha > 0$ , learning rate  $\gamma \in [0, 1]$ , number of episodes, length of episodes

Initialize:  $V(s) = 0$ ,  $Z(s) = 0$ ,  $\forall s \in S$ , NEU\_episodes = empty list, NEU\_steps = empty list

Loop for each episode:

$s \leftarrow 0$

$Z(s) \leftarrow 0 \forall s$

NEU\_steps  $\leftarrow$  empty list

Loop for each step of episode:

$a \leftarrow \pi(s)$

Take action  $a$ , observe reward  $r$ , and next state  $s'$

$\delta \leftarrow r + \gamma V(s') - V(s)$

$Z(s) \leftarrow Z(s) + 1$

For all  $s$ :

$V(s) \leftarrow V(s) + \alpha \delta Z(s)$

$Z(s) \leftarrow \gamma \lambda Z(s)$

NEU\_steps.append( $(\delta Z(s)^2)$ )

$s \leftarrow s'$

If  $s$  is terminal or end of episode

NEU\_episodes.append(sum(NEU\_steps)/number of steps in episode)

Go to next episode

### HW3 Generate Plots

Algorithm parameters: lambda\_list = [0, 0.3, 0.5, 0.7, 1]

For lambda in lambda\_list:

NEU\_episodes  $\leftarrow$  TD( $\pi$  = random policy,  $\lambda$ =lambda)

NEU\_averaged  $\leftarrow$  Place entries of NEU\_episodes into bins of length 10, average each bin

Plot NEU\_averaged vs episode counts for current lambda

Display plot

Briefly explain the impacts of  $\lambda$  on the performance of TD learning:

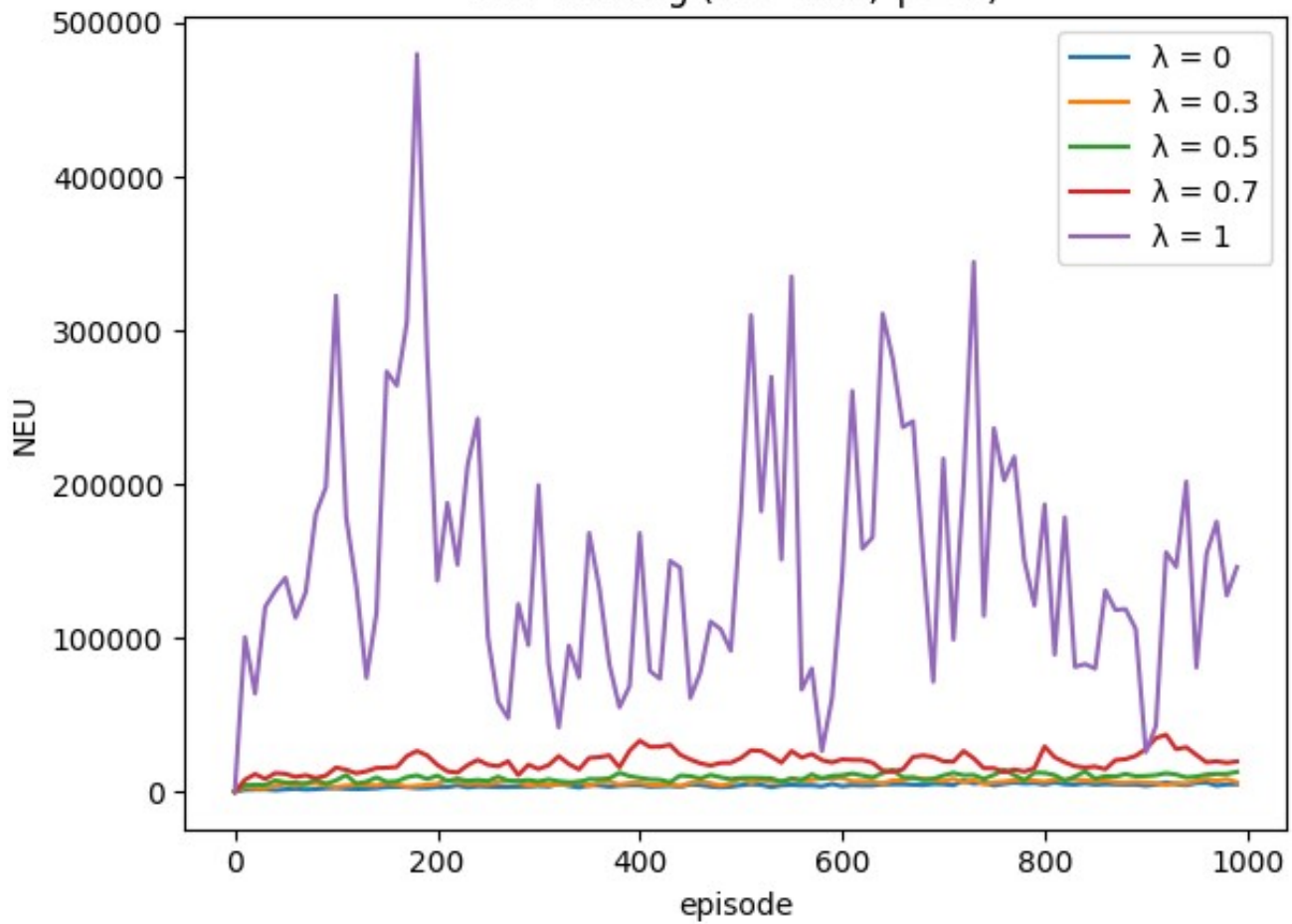
The value of lambda ( $\lambda \in [0, 1]$ ) governs the decay rate of the trace vector. These trace vectors determine the degree to which previous predictions at a given state are eligible for updates. Values of  $\lambda$  closer to zero cause the trace vector to decay faster, whereas  $\lambda$  values closer to one create a longer lasting trace. As such, high values of  $\lambda$  result in value predictions of states visited further into the past getting assigned credit/blame (based on currently observed errors) to a greater extent than if  $\lambda$  were lower.

When  $\lambda = 0$  you have Temporary Difference Learning (TD(0)), and when  $\lambda = 1$  you effectively have Monte Carlo (MC). TD(0) and MC have different characteristics, with one being preferable over the other for a given problem. Lambda provides the ability to pick a point between these two extremes, maximizing the cumulative benefit.

For example, MC has high variance, zero bias, is not very sensitive to the initial values chosen, and does not exploit the Markov property (and therefore usually more effective in non-Markov environments). Whereas TD(0) has low variance, some bias, is sensitive to the initial value, and exploits the Markov property via bootstrapping (and therefore usually more effective in Markov environments). Tuning the value of  $\lambda$  for a given problem and environment can therefore speed up convergence through utilizing these respective properties.

Note that the increasing variance can be observed within the following graphs.

Cliff walking ( $\alpha = 0.01, \gamma = 1$ )



Cliff walking without  $\lambda=1$  ( $\alpha = 0.01, \gamma = 1$ )

