Name: Alfred
Student No.: 2020-12788

# Exercise 4
## Data Science & Reinforcement Learning
### Spring 2023

$$E\left[\left(\gamma \max_{a'} Q_t(s',a') - (\bar{T}Q_t)(s,a)\right)^2\right]$$

1. Modify the Tabular TD(0) algorithm for estimating $V_\pi$ below to estimate $Q_\pi$.

---

**Tabular TD(0) for estimating $v_\pi$**

Input: the policy $\pi$ to be evaluated
Algorithm parameter: step size $\alpha \in (0,1]$
Initialize $V(s)$, for all $s \in S^+$, arbitrarily except that $V(terminal) = 0$

Loop for each episode:
  Initialize $S$
  Loop for each step of episode:
    $A \leftarrow$ action given by $\pi$ for $S$
    Take action $A$, observe $R$, $S'$
    $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$
    $S \leftarrow S'$
  until $S$ is terminal

---

Algorithm parameters: step size $\alpha \in (0,1]$, small $\varepsilon > 0$
Initialize $Q(s,a)$, for all $s \in S^+$, $a \in A(s)$ arbitrarily except $Q(terminal, \cdot) = 0$.

Loop for each episode:

  Initialize S.

  Loop for each episode:

    Take action A, observe $R, s'$

    choose $A'$ from $s'$ using policy from $Q$

    $Q(S,A) \leftarrow Q(S,A) + \alpha[R + \gamma Q(s',A') - Q(S,A)]$

    $S \leftarrow s'; A \leftarrow A'$

  until $S$ is terminal.

---

5) $\quad Var(F_t(s,a) \mid H_t) = \mathbb{E}\left[\left(r_t(s,a) + \gamma \max_{a'} Q_t(s',a') - Q^*(s,a) - (\bar{T}Q_t)(s,a) + Q^*(s,a)\right)^2\right]$

$\quad = Var\left[ r_t(s,a) + \gamma \max_{a'} Q_t(s',a') \mid H_t \right] = Var\left[ r_t(s') + \gamma \max_a Q_t(s,a)\right]$

$\quad$ reward $\wedge \max_a Q_t(s',c)$ bounded $\qquad$ bounded $\qquad$ then $Var\left[F_t(s,a)\mid H_t\right] \le C(1 + \|\Delta_t\|_\infty^2)$

then by lemma 2.

$\quad \Delta_t(e) = Q_t(s,a) - Q^*(s-a) \to 0.$

$\quad$ convergence: $Q_t(s,a) = Q^*(s-c)$ w.p. 1

$\max_a Q_t(s',a) - \max_{s,a} Q^*(s,c)$
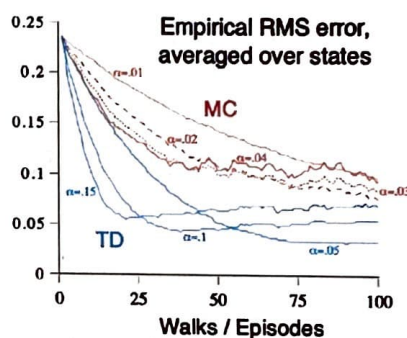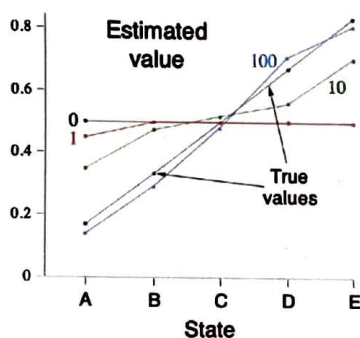$\le \max_{s,a} \|Q_t - Q^*\| \le \|\Delta_t\|_\infty.$

2. Consider the following *Random Walk* example shown in the textbook.

---

**Example 6.2　Random Walk**

In this example we empirically compare the prediction abilities of TD(0) and constant-$\alpha$ MC when applied to the following Markov reward process:



A *Markov reward process*, or MRP, is a Markov decision process without actions. We will often use MRPs when focusing on the prediction problem, in which there is no need to distinguish the dynamics due to the environment from those due to the agent. In this MRP, all episodes start in the center state, C, then proceed either left or right by one state on each step, with equal probability. Episodes terminate either on the extreme left or the extreme right. When an episode terminates on the right, a reward of +1 occurs; all other rewards are zero. For example, a typical episode might consist of the following state-and-reward sequence: C, 0, B, 0, C, 0, D, 0, E, 1. Because this task is undiscounted, the true value of each state is the probability of terminating on the right if starting from that state. Thus, the true value of the center state is $v_\pi(C) = 0.5$. The true values of all the states, A through E, are $\frac{1}{6}, \frac{2}{6}, \frac{3}{6}, \frac{4}{6},$ and $\frac{5}{6}$.



The left graph above shows the values learned after various numbers of episodes on a single run of TD(0). The estimates after 100 episodes are about as close as they ever come to the true values—with a constant step-size parameter ($\alpha = 0.1$ in this example), the values fluctuate indefinitely in response to the outcomes of the most recent episodes. The right graph shows learning curves for the two methods for various values of $\alpha$. The performance measure shown is the root mean square (RMS) error between the value function learned and the true value function, averaged over the five states, then averaged over 100 runs. In all cases the approximate value function was initialized to the intermediate value $V(s) = 0.5$, for all $s$. The TD method was consistently better than the MC method on this task.

---

On the left plot for the estimated values, it appears that the first episode results in a change in only $V(A)$. What does this tell you about what happened in the first episode? Why was only the estimate for $V(A)$ changed? By exactly how much was it changed?

- first episode is untill state A without knowing the transition
  (must have termitated at State A)

$$V_A \leftarrow V_A + \alpha(R + \gamma V_+ - V_A) = 0.5 - 0.1(0 + 1 - 0.5) = 0.45$$

- only V(A) changed because rewards are 0 ∧ $\gamma = 1$ since initialized at V(s) = 0.5. TD update doesn't update some V.

For A, E states only non-zero changes are possible.

$$\Delta V_A = -0.05$$

3. In the *Random Walk* example of Question 2, the curves on the right plot are dependent on the value of the step-size parameter $\alpha$. Based on these results, can you come up with a strategy for the choice of the step-size parameter $\alpha$? Justify your answer.

Small parameter steps can definitely reduce error for both MC and TD.

But when at limit it will reach the minimum lowest error.

Then arbitrary/strategical choice won't improve the plot.

4. Q-learning is an off-policy algorithm, but there does not seem to be any importance sampling ratio used in the update rule. Why is it the case?

target policy is greedy w.r.t $Q(S,A)$ while b is another different policy.

(Importance)
Sampling

It's not needed because $Q(S,A) \leftarrow Q(S,A) + \alpha (r + \sigma \max_a (Q(s',a) - Q(S,A))$

learns from taking various actions (even random) which don't need

policy. We use the max Q instead.   first time only.

from behavior policy, A is sampled the

later weren't so no need of importance sampling.

$a$ are all actions probed in state $s'$

5. Prove the convergence of Q-learning.

**Lemma 1 (Contraction)**

the operator $T$ defined for generic function $q: S \times A \to \mathbb{R}$ as

$$(\bar{T}q)(s,a) := \sum_{s' \in S} P(s'|s,a)\left[r(s,a) + \gamma \max_{a'} q(s',a')\right]$$

this operator is a contraction in sup-norm, i.e. $\sim \infty$.

$$\|\bar{T}q_1 - \bar{T}q_2\|_\infty \leq \gamma \|q_1 - q_2\|_\infty$$

**Proof**

$$\|\bar{T}q_1 - \bar{T}q_2\|_\infty = \max_{s,a}\left|(\bar{T}q_1)(s,a) - (\bar{T}q_2)(s,a)\right|$$

$$= \max_{s,a}\left|\sum_{s'\in S} P(s'|s,a)\left[r(s,a) + \gamma \max_{a'} q_1(s'-a') - r(s,a) - \gamma \max_{a'} q_2(s',a')\right]\right|$$

$$= \max_{s,a} \gamma\left|\sum_{s'\in S} P(s'|s,a)\left[\max_{a'} q_1(s',a') - \max_{a'} q_2(s',a')\right]\right|$$

$$\leq \max_{s,a} \gamma \sum_{s'} P(s'|s,a) \max_{s',a'}\left|q_1(s',a') - q_2(s',a')\right|$$

$$\underbrace{\phantom{XXXXXX}}_{\|q_1 - q_2\|_\infty}$$

$$= \gamma\|q_1 - q_2\|_\infty \underbrace{\sum_{s'} P(s'|s,a)}_{1}$$

$$\boxed{\mathbb{E}_x[f(x)] \leq \max_x f(x)}$$
$$\boxed{\bar{T}Q^* = Q^*}$$

**Lemma 2 (Stochastic approx.):** A random process $\{\Delta_t\}_t$ taking values in $\mathbb{R}^n$ defined as

Jaakkola, Jordan, Siegl, 1993

$$\Delta_{t+1}(x) := (1 - \alpha_t(x))\Delta_t(x) + \alpha_t(x) F_t(x) \; ; \quad \Delta_t(x) \to 0 \text{ w. P. } 1 \text{ under:}$$

i) $0 \leq \alpha_t \leq 1, \sum_t \alpha_t(x) = \infty$

$\sum_t \alpha_t^2(x) < \infty$

**Q-learning**

$$Q_{t+1}(s_T, a_t) = Q_t(s_T, a_t) + \alpha_t(s_t, a_t)\left[r_t + \gamma \max_{a'} Q_t(s_{t+1}, a') - Q(s_t, a_t)\right]$$

ii) $\max_x \left|\mathbb{E}[F_t(x)|H_t]\right|$

$\leq \gamma \max_x |\Delta_t(x)|$

WTS $Q_t \to Q^*$ w.P.1

under $\alpha (\sum \alpha = \infty \wedge \sum \alpha^2 < \infty)$

$= \gamma\|\Delta_t\|_\infty$

**[Proof]** Define $\Delta_{t,a} \, \overset{df}{=} \Delta_t(s,a) := Q_t(s,a) - Q^*(s,a)$

iii) $\text{Var}(F_t(x)|H_t) \leq C(1 + \|\Delta_t\|_\infty^2)$

$$\to \Delta_{T+1}(s,a) = (1 - \alpha_t(s,a))\Delta_t(s,a) + \alpha_t(s,a)\left[r_T + \gamma \max_{a'} Q_t(s',a) - Q^*(s,a)\right]$$   for $x, C > 0$

Define $F_t(s,a) = r_t(s,a) + \gamma \max_{a'} Q_t(x,a) - Q^*(s,a)$, $X$: random sample from $= P(\cdot|s,a)$ rMDP

$$\mathbb{E}[F_t(s,a)|H_t] = \sum_{s'} P(s'|s,a)\left[r(s,a) + \gamma \max_{a'} Q_t(s',a') - Q^*(s,a)\right] \quad \text{4 of 4}$$

$H_t = \{\Delta_T, \Delta_{T-1}, \ldots F_{t-1}, \alpha_{t-1}\}$

$$= (\bar{T}Q_t)(s,a) - Q^*(s,a) = (\bar{T}Q_t)(s,a) - (\bar{T}Q^*)(s,a)$$

$$\|\mathbb{E}[F_t|H_t]\|_\infty = \|\bar{T}Q_t - \bar{T}Q^*\|_\infty \leq \gamma\|Q_t - Q^*\|_\infty \text{ (by Lemma 1)} = \gamma\|\Delta_t\|_\infty.$$

$\boxed{\text{PAGE I}}$