

Exercise 5

Data Science & Reinforcement Learning

Spring 2023

1. Recall the n -step value function is defined as

$$V_{t:t+n}(S_t) = V_{t+n-1}(S_t) + \alpha [G_{t:t+n} - V_{t+n-1}(S_t)], \quad 0 \leq t \leq T$$

where $G_{t:t+n}$ is the n -step return

$$G_{t:t+n} := R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n}) \quad V_{(S_t)} = V_{(S_{t+1})} = \dots = V$$

for $0 \leq t < T-n$, and $G_{t:t+n} := G_t$ for $t \geq T-n$. Show that the n -step TD error, $G_{t:t+n} - V_{t+n-1}(S_t)$, can be written as a sum of TD errors if the value estimates don't change from step to step.

$$\begin{aligned} \sum_{k=t}^{\min(t+n, T)-1} \delta_k &= (R_{t+1} + \gamma V_{(S_{t+1})} - V_{(S_t)}) + (R_{t+2} + \gamma V_{(S_{t+2})} - V_{(S_{t+1})}) + \dots + (R_{t+n} + \gamma V_{(S_{t+n})} - V_{(S_{t+n-1})}) \\ &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{(S_{t+n})} - V_{(S_t)} \\ &= G_{t:t+n} - V_{t+n-1}(S_t) \end{aligned}$$

2. Recall the n -step return of Sarsa defined as

$$G_{t:t+n} := R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n})$$

for $0 \leq t < T-n$, and $G_{t:t+n} := G_t$ for $t \geq T-n$. Prove that the n -step return of Sarsa can be written exactly in terms of a novel TD error, as

$$\begin{aligned} G_{t:t+n} &= Q_{t-1}(S_t, A_t) + \sum_{k=t}^{\min(t+n, T)-1} \gamma^{k-t} [R_{k+1} + \gamma Q_k(S_{k+1}, A_{k+1}) - Q_{k-1}(S_k, A_k)] \\ \text{RHS} &= \sum_{k=t}^{\min(t+n, T)-1} \gamma^{k-t} R_{k+1} + \sum_{k=t}^{\min(t+n, T)-1} \gamma^{k-t} (Q_k(S_{k+1}, A_{k+1}) - Q_{k-1}(S_k, A_k)) \\ &\quad \left\{ \begin{array}{l} R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} \quad t < T-n \\ G_T \quad t \geq T-n \end{array} \right. \\ &\quad \left\{ \begin{array}{l} \gamma (Q_t(S_{t+1}, A_{t+1}) - Q_{t-1}(S_t, A_t)) \\ + \gamma^2 (Q_{t+2}(S_{t+2}, A_{t+2}) - Q_t(S_{t+1}, A_{t+1})) \\ \vdots \\ \gamma^n (Q_{t+n-1}(S_{t+n}, A_{t+n}) - Q_{t+n-2}(S_{t+n-1}, A_{t+n-1})) \end{array} \right. \\ &\quad \left\{ \begin{array}{l} 0 - Q_{t-1}(S_t, A_t) \quad (t \geq T-n) \\ (-Q_{\text{terminal}} = 0) \end{array} \right. \end{aligned}$$

1 of 4

then equal sides.

3. An agent observes the following two episodes from an MDP,

$$\begin{aligned} & \neg S_0 = 0, A_0 = 1, R_1 = 1, \underline{S_1 = 1}, A_1 = 1, R_2 = 1 \\ & S_0 = 0, A_0 = 0, R_1 = 0, \underline{S_1 = 0}, A_1 = 1, R_2 = 1, S_2 = 1, A_2 = 1, R_3 = 1 \end{aligned}$$

and updates its deterministic model accordingly. What would the model output for the following queries:

- (a) Model($S = 0, A = 0$): ~~$S = \text{initial}$~~ , ~~$R = 0$~~ . $S = 0, R = 0$
- (b) Model($S = 0, A = 1$): ~~$S = \text{initial}$~~ , ~~$R = 1$~~ . $S = 1, R = 1$
- (c) Model($S = 1, A = 0$): ~~$S = 1$~~ , ~~$R =$~~ initial values
- (d) Model($S = 1, A = 1$): ~~$S = 1$~~ , ~~$R =$~~ $R = 1, S = \text{terminal}$

4. Why did the Dyna agent with exploration bonus, Dyna-Q+, perform better in the first phase as well as in the second phase of the blocking experiment in the figure below.

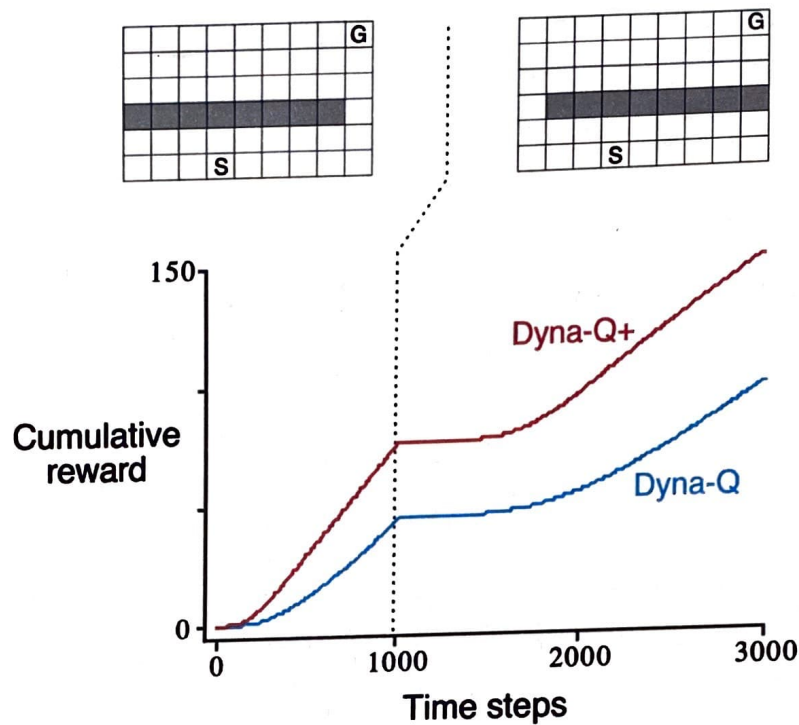


Figure 1: Average performance of Dyna agents on a blocking task. The left environment was used for the first 1000 steps, the right environment for the rest. Dyna-Q+ is Dyna-Q with an exploration bonus that encourages exploration

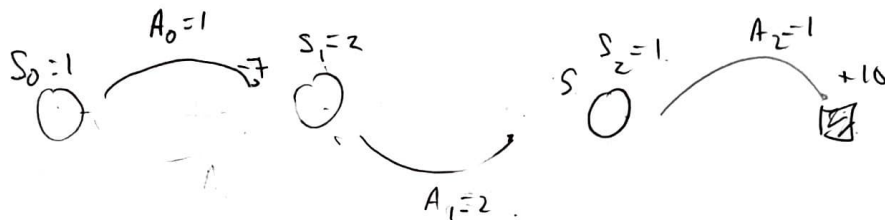
First half: early beginnings after Dyna-Q⁺ tend to give higher R,
 for second phase Dyna-Q's poorer than Dyna-Q⁺ because
 it explores more, finds optimal policy faster as well as
 model's difference too. Phase 1 & 2 (after 1000 time steps)
 the dynamics change too.

because that's
 bonus
 reward

5. Consider an MDP with three states $S = \{1, 2, 3\}$, where each state has two possible actions $A = \{1, 2\}$ and a discount rate $\gamma = 0.5$. Suppose estimates of $Q(S, A)$ are initialized to 0 and you observed the following episode according to an unknown behaviour policy where S_3 is the terminal state.

$$S_0 = 1, A_0 = 1, R_1 = -7, S_1 = 2, A_1 = 2, R_2 = 5, S_2 = 1, A_2 = 1, R_3 = 10$$

- (a) Suppose you used Q-learning with the above trajectory to estimate $Q(S, A)$, what are your new estimates for $Q(S = 1, A = 1)$ using $\alpha = 0.1$?
- (b) Suppose in the planning loop, after search control, we would like to update $Q(S = 1, A = 1)$ with Q-planning. What are the possible outputs of $\text{Model}(S = 1, A = 1)$?



a)

$$Q^{\text{new}}(S=1, A=1) = 0 + 0.1(-7 + 0.5(0) - 0) \quad \text{1st iteration}$$

2nd iteration.

$$Q^{\text{new}}(S=1, A=1) = -0.7$$

$$Q(1,1) \leftarrow Q(1,1) + \alpha(10 + 0.5 \times 0 - Q(1,1)) = 1 - 0.63 = 0.37$$

b)

$$R, S = (-7, 2) \quad N(\text{10, terminate}) \text{ both true.}$$

or if stochastic $\rightarrow \{(-7, 2), (10, \text{terminate})\}$ true too

50% 50%