

# RL HW EX6:

1. Planning vs. n-step returns.
  - a. The nonplanning method looks particularly poor in Figure 8.3 because it is a one-step method; a method using multi-step bootstrapping would do better. Do you think one of the multi-step bootstrapping methods from Chapter 7 could do as well as the Dyna method? Explain why or why not.
    - i. No, I wouldn't expect a multi-step bootstrapping (MSB) method from Chapter 7 to perform as well as Dyna. MSB methods rely solely on real experience and require more iterations to learn effectively. In contrast, Dyna's planning component allows it to propagate value updates quickly by generating simulated experiences, significantly reducing the sample complexity. This ability to update values without needing additional real experiences gives Dyna a distinct advantage in learning efficiency.
  - b. See the pseudocode on page 164 in RL2e. Consider using both n-step returns and Dyna (specifically, using n-step returns instead of one-step returns in tabular Dyna-Q). We can clearly do this for the learning phase (d). Can we also do this for the planning phase (f)? What are the advantages and disadvantages of using n-step returns in the planning phase (f)?
    - i. Yes, n-step returns can be used in the planning phase (f) of Dyna-Q, allowing rewards to propagate over longer distances, which could speed up learning in tasks with delayed rewards. However, this approach also could also hamper the planning phase as it increases computational cost and introduces higher variance, which may affect stability and slow down the algorithm.
2. Implementing Dyna-Q and Dyna-Q+.
  - a. For some unknown reason, the textbook is quite vague about Dyna-Q+ and does not provide pseudocode for the modification. Read Section 8.3 carefully (specifically p. 168) and reconstruct the pseudocode for Dyna-

Q+. In particular, please write your pseudocode that also contains the suggestion from the footnote.

i.

```
Dyna-Q+
Initialize Q(s,a) and Model(s, a) for all s in S and a in A

# Store T values for the bonus
Initialize LastTried(s, a) to 0 for all s, a

# small k for exploration bonus
Set k > 0
Loop forever:
(a) S = current (nonterminal) state
(b) A = e-greedy(S,Q)
(c) Take action A; observe resultant reward, R, and state, S'
(d) Q(S, A) = Q(S, A) + alpha[R + gamma * maxa Q(S', a) - Q(S, A)]
(e) Model(S, A) = R, S'(assuming deterministic environment)

# Just tried to reset T
(new step) LastTried(s, a) = 0
(f) Loop repeat n times:
    S = random previously observed state

    # action not yet tried allowed to be considered in this step
    A = random action from A(S)

    # initial model for untried action leads back to same state
    if (S, A) has never been tried:
        Model(S, A) = (0, S)

    R, S' = Model (S, A)

    # Exploration bonus
    bonus = k * sqrt(TimeSinceLastVisit(S, A))
```

```

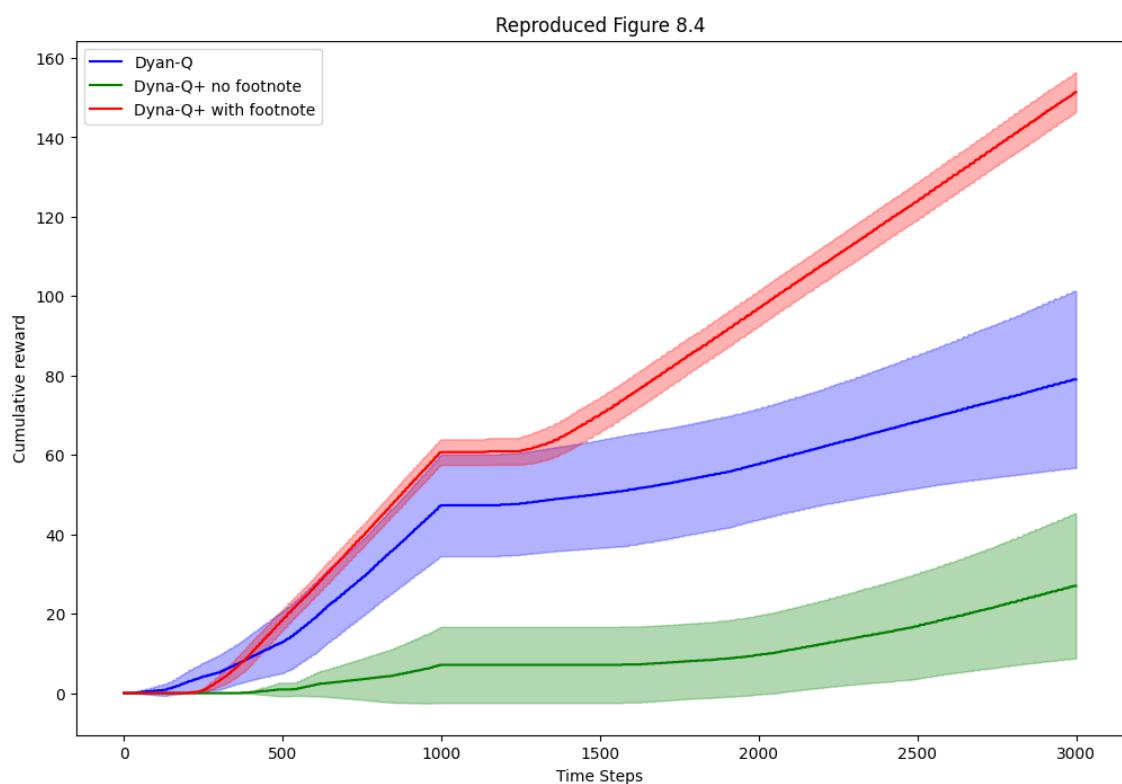
Q(S, A) = Q(S, A) + alpha[(R + bonus) + gamma * maxa (S', A')]

# Just tried to reset T
LastTried(S, A) = 0

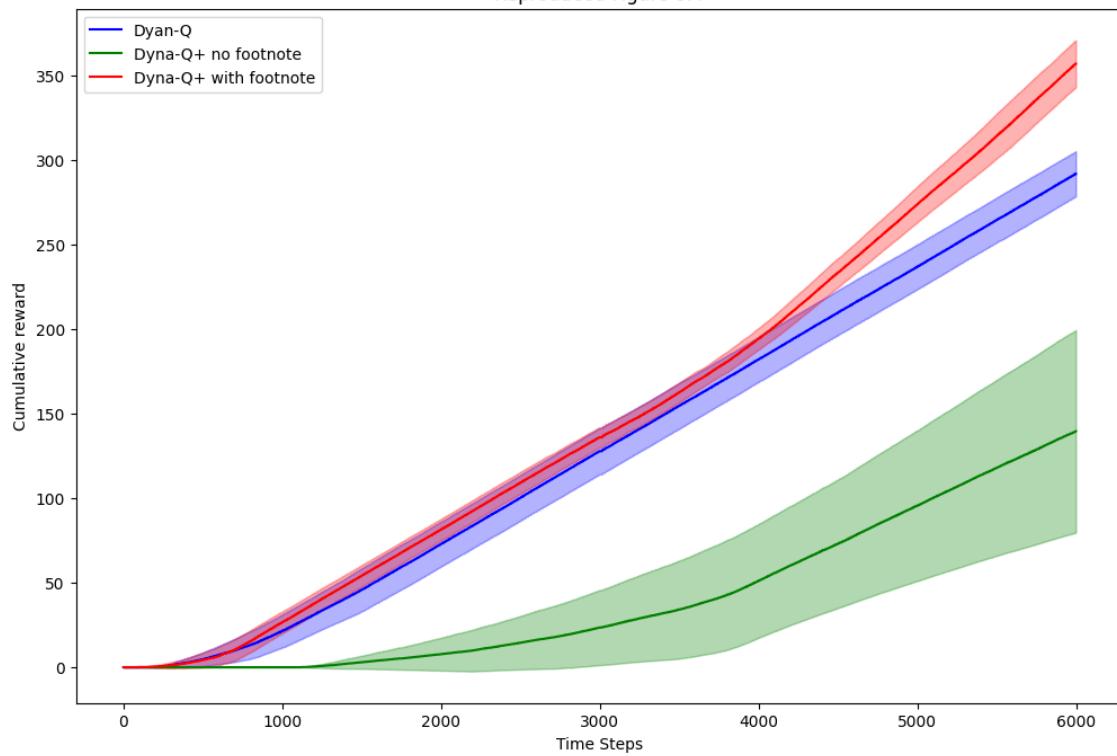
# Took a step so add 1 to all
(new step) LastTried += 1

```

b. Plot:



Reproduced Figure 8.4



c. Did the footnote matter? Why or why not?

i. Yes the footnote does matter as with the dyna-q+ we add a bonus reward that promotes exploration. Specifically, it promotes choosing untried actions. This is where the footnote comes in handy as without the footnote, untried actions in a state are ignored. This footnote makes a major difference in the changing environments we were looking at, as it prevents the agent from getting stuck in familiar but suboptimal actions, leading to faster discovery of optimal actions and improved performance.

3. Think carefully about the implications of applying bonuses to the reward function ("pseudo-reward") versus during action selection. Make some predictions:

a. What would generally happen to the policy's behavior during the training under each approach? How do they differ from each other?

i. When using exploration bonuses as a pseudo-reward (Dyna-Q+ style), the bonus directly impacts the Q-values by integrating into the simulated rewards, leading the learned policy to inherently favor

persistent exploration. This approach results in a policy that continuously balances exploration and exploitation, as actions with higher bonuses (even if they yield lower actual rewards) are prioritized. In contrast, when bonuses are applied only during action selection (UCB style), they temporarily influence which actions are chosen without altering the Q-values. This setup drives short-term exploration but allows the agent to gradually shift towards exploitation as it becomes more confident in the true reward structure, leading to reduced exploration once an optimal policy emerges. Thus, the key difference is that pseudo-reward bonuses embed exploration into the learned policy itself, while action selection bonuses encourage early exploration but ultimately favor exploitation.

b. What are the advantages and disadvantages of each approach?

i. The pseudo-reward approach has the advantage of promoting ongoing exploration by embedding the bonus into Q-values, making it particularly effective in dynamic environments where things change, like in the blocking and shortcut maze examples. This lets the agent continually explore and adapt to new paths or obstacles. The downside to this is that it can distort the true value estimates, leading the agent to sometimes over-prioritize exploration over finding the actual optimal path in stable settings. In contrast, the action selection approach has an advantage in that it keeps Q-values accurate, encouraging short-term exploration that naturally shifts to exploitation, which is advantageous in stationary environments. However, in complex or changing settings, this approach can lead to exploration tapering off too quickly, causing the agent to miss out on learning better routes when the environment shifts.

5.

## Iteration 4

Selection:  $\textcircled{C} \xrightarrow{R} \textcircled{D}$

Expansion: 1 (no use md)

$$\frac{2.05}{2} + \frac{\sqrt{2 \ln 3}}{2} + \frac{0}{1} + \frac{\sqrt{2 \ln 3}}{1}$$

$$\max\left(\frac{T_r}{N_r} + \sqrt{\frac{2 \ln N_e}{N_r}}, \frac{T_L}{N_L} + \sqrt{\frac{2 \ln N_e}{N_L}}\right)$$

## Iteration 4

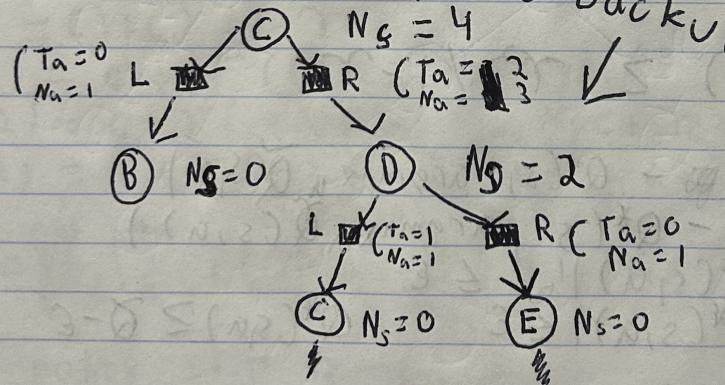
Selection:  $\max_a\left(\frac{T_a}{N_a} + \sqrt{\frac{2 \ln N_e}{N_a}}\right)$

$\max(R: (2/2) + \sqrt{2 \ln(3)})/2, L: (0/1) + \sqrt{2 \ln(3)/1}\right)$

$\max(R: 2.05, L: 1.48) = R$

$\textcircled{C} \xrightarrow{R} \textcircled{D}$

Expansion: 1 (no use md)  $\rightarrow \textcircled{E}$  Backup



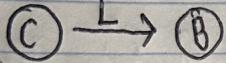
Simulation:  $\textcircled{B} \xrightarrow{L} \textcircled{D} \xrightarrow{L} \textcircled{C} \xrightarrow{R} \textcircled{D} \xrightarrow{L} \textcircled{C} \xrightarrow{R} \textcircled{D} \xrightarrow{R} \textcircled{E} \xrightarrow{L} \textcircled{D} \xrightarrow{L} \textcircled{C} \xrightarrow{L} \textcircled{B}$

$\hookrightarrow$  Term reward 0

Iteration 5:

Selection: C:  $\max (R: \frac{2}{3} + \sqrt{2 \ln(4)/3}, L: \sqrt{2 \ln(4)})$

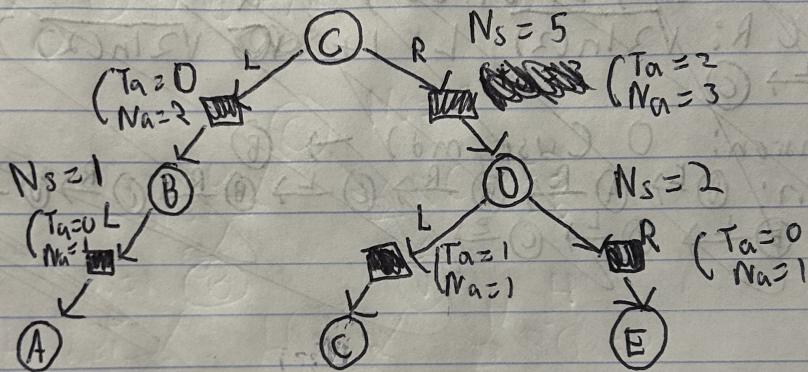
$$\max (R: 1.63, L: 1.66) = L$$



Expansion: (B) (use  $n_d$ )  $\rightarrow$  (A)

Simulation: (A)  $\xrightarrow{L}$  Terminal reward 0

Backup



## Iteration 6

Selection: C:  $\max CR: 2/3 + \sqrt{2\ln(5)/3}$ , L:  $\sqrt{2\ln(5)/2}$

$\max(CR: 1.70, L: 1.26) = R$  ~~C~~  $\xrightarrow{C \rightarrow D}$

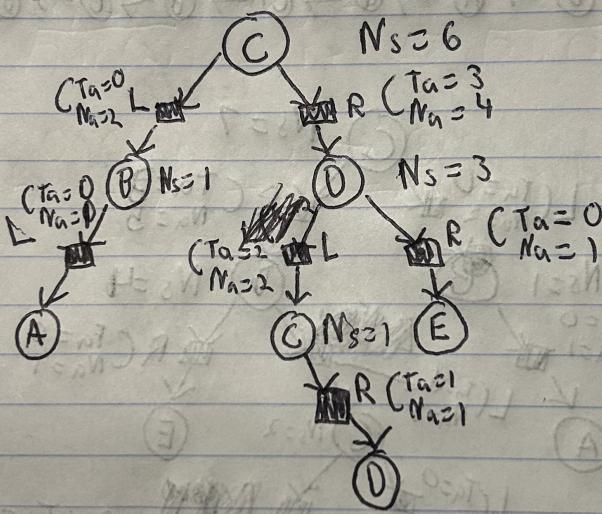
D:  $\max CR: \sqrt{2\ln(2)}$ , L:  $1 + \sqrt{2\ln(2)} = L$

$\circlearrowleft \rightarrow \textcircled{1} \xrightarrow{\textcircled{1}} \textcircled{2}$

Expansion: 1 (use rnd)  $\textcircled{2} \xrightarrow{R} \textcircled{1}$

Simulation:  $\textcircled{1} \xrightarrow{R} \textcircled{2} \xrightarrow{L} \textcircled{1} \xrightarrow{R} \textcircled{2} \xrightarrow{R} \textcircled{3}$  Terminal reward = 1

Backup



# Iteration 7

Selection: C:  $\max(C, R: 3/4 + \sqrt{2 \ln(6)/4}, L: \sqrt{2 \ln(6)/2})$

$$\max(C: 1.70, L: 1.26) = R \quad C \xrightarrow{R} D$$

~~$$D: \max(C: \sqrt{2 \ln(3)}, L: 2/2 + \sqrt{2 \ln(3)/2})$$~~

~~$$\max(C: 1.48, L: 2.05) = L$$~~

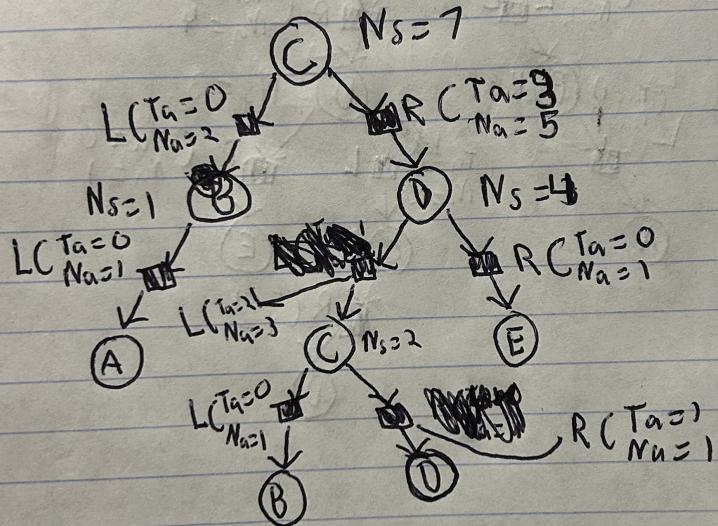
~~$$① \xrightarrow{R} ② \xrightarrow{L} ③$$~~

Expansion: ① (no rnd)

Simulation: ③  $\xrightarrow{A} ④ \xrightarrow{R} ⑤ \xrightarrow{L} ⑥ \xrightarrow{B} ⑦ \xrightarrow{R} ⑧ \xrightarrow{L} ⑨ \xrightarrow{A} \text{Term}$

reward = 0

Backup



## Iteration 8

Selection: C:  $\max(R: 3.15 + \sqrt{2 \ln(7)}/5)$ , L:  $\sqrt{2 \ln(7)}/12$

$$\max(R: 1.48, L: 1.39) = R \quad C \xrightarrow{R} D$$

$$D: \max(R: \sqrt{2 \ln(4)}) / 3, L: 2/3 + \sqrt{2 \ln(4)}/3$$

$$\max(R: 1.67, L: 1.63) = R \quad D \xrightarrow{R} E \quad C \xrightarrow{R} D \xrightarrow{R} E$$

Expansion: I ~~rnd~~ (use rnd)  $E \xrightarrow{!} \text{Terminal}$

