

## **Exercise 6**

Gregory Attra

03.10.2022

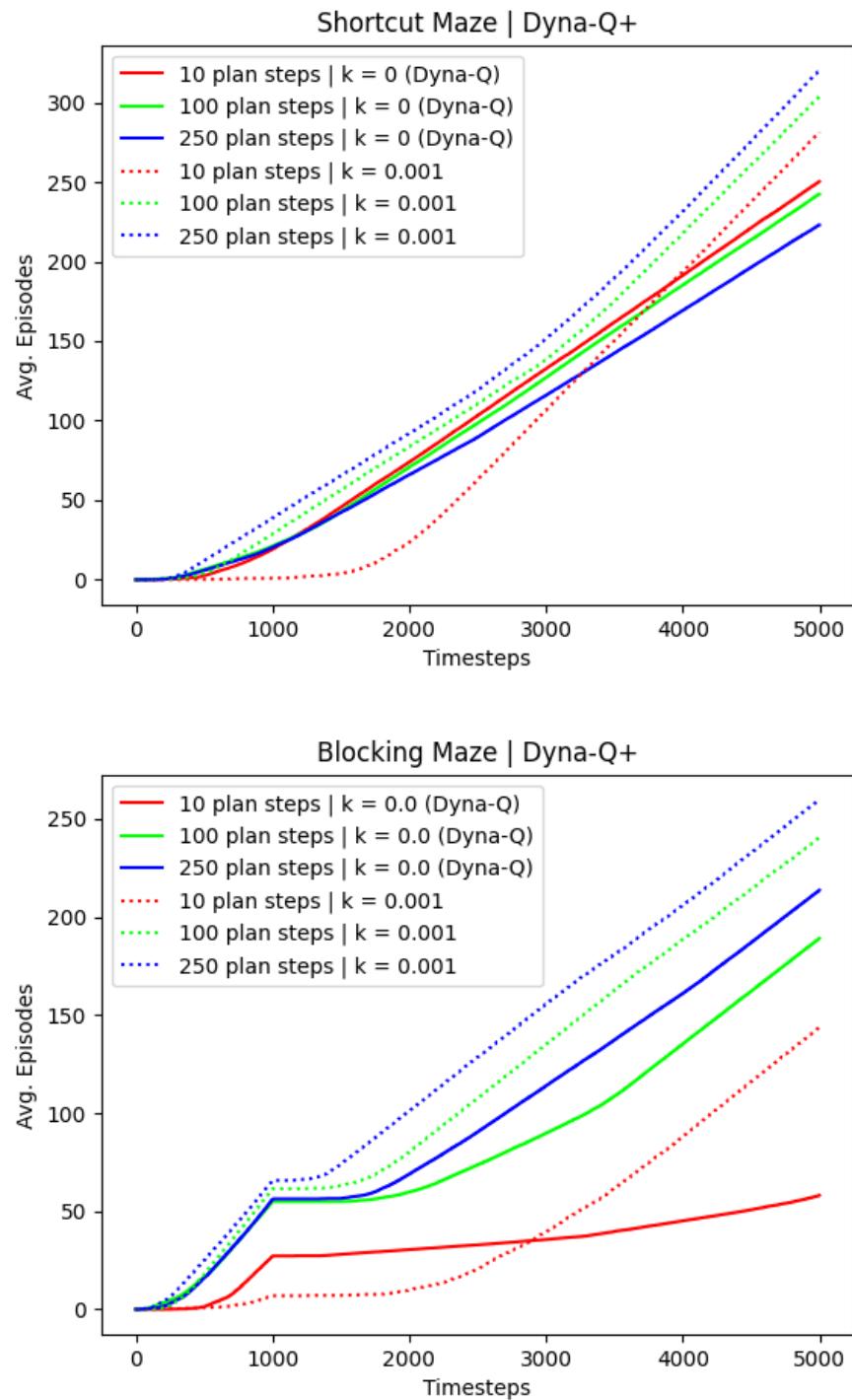
CS 7180 - Prof. Amato

### **Code Setup:**

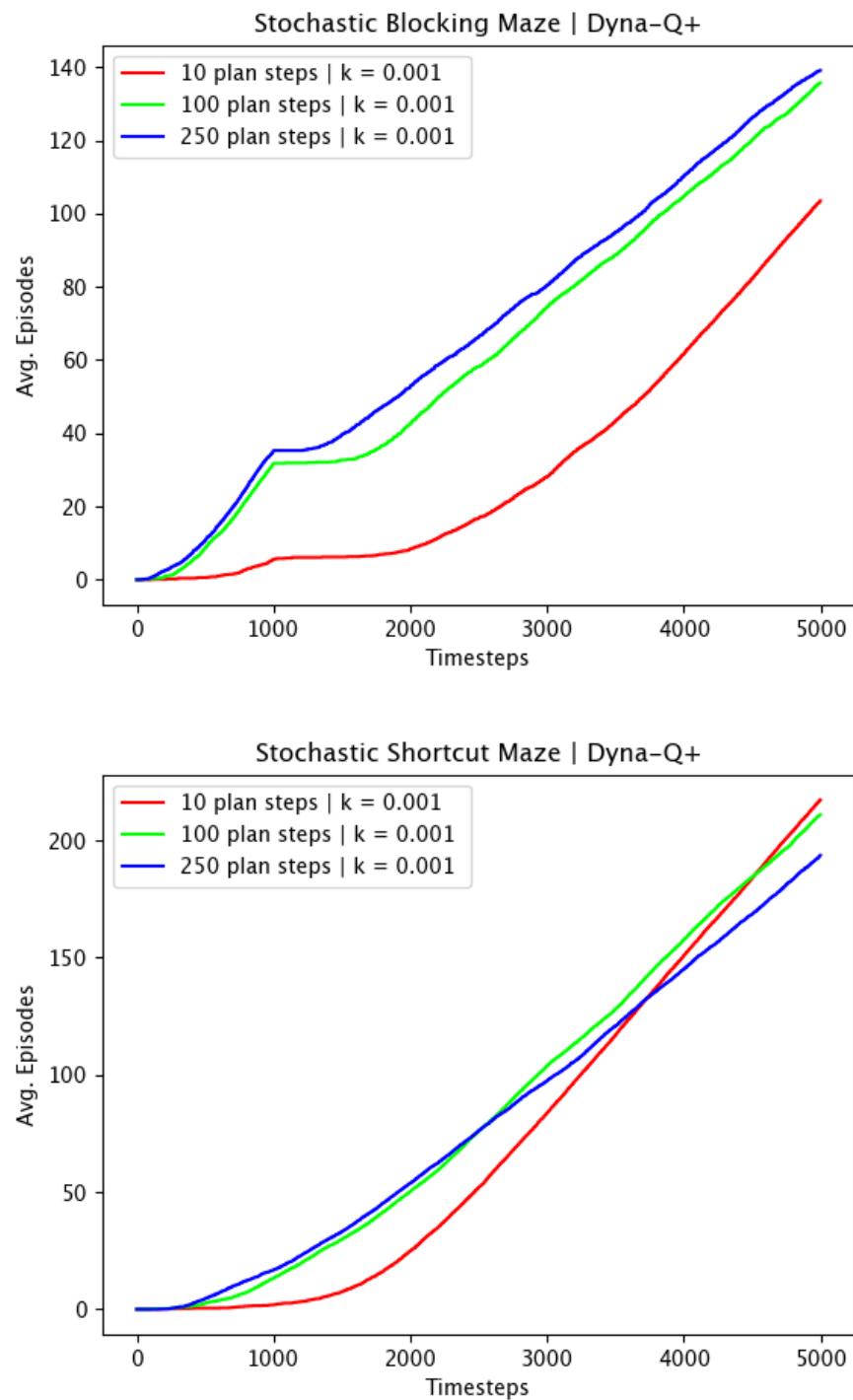
1. Unzipped the ‘ex6.zip’ file
2. ‘cd‘ into the ‘code‘ directory
3. run ‘source ./init‘ to setup the pyenv

### **Questions:**

1. (a) It would perform slightly better as, after the first episode, the final n-steps would be updated. But it would still not perform as well as Dyna-Q as we would still not make as many updates. In Dyna-Q, we could sample  $n + 1$  states and have better results, especially if we used prioritized sampling. We would also make updates to states not necessarily in our trajectory, unlike with n-step SARSA.  
(b) i. Advantage: We have a higher likelihood of making meaningful updates; in other words we have a higher likelihood of sampling a trajectory where  $Q(s, a)$  was updated from learning or from a prior planning update  
ii. Disadvantage: Computationally more expensive:  $M \times N$  samples each step.
2. The benefit of Dyna-Q+ is that it incentivizes exploring, but this can also mean taking suboptimal actions. In this case, Dyna-Q+ learned a more optimal policy initially, but overtime began to explore less-visited states, which brought down the average reward.
3. (a) The modifications to make to Dyna-Q are:
  - For each state/action pair, track the most-recent timestep at which that state/action pair was taken in the experience phase.
  - In the planning phase, compute  $\tau = t - T[s][a]$  where  $t$  is the current timestep and  $T[s][a]$  is the most recent timestep at which the state action pair  $(s, a)$  was executed in the experience phase.
  - Add  $k\sqrt{\tau}$  to the model’s reward for the state action pair:  $r = M[s][a]$ ; where  $k$  is some exploration balancing constant.  
(b) To run the Dyna-Q+ implementation:
  - i. Follow the code setup instructions at the top of this document
  - ii. Run ‘python src/planning/run\_shortcut\_maze.py‘
  - iii. Run ‘python src/planning/run\_blocking\_maze.py‘

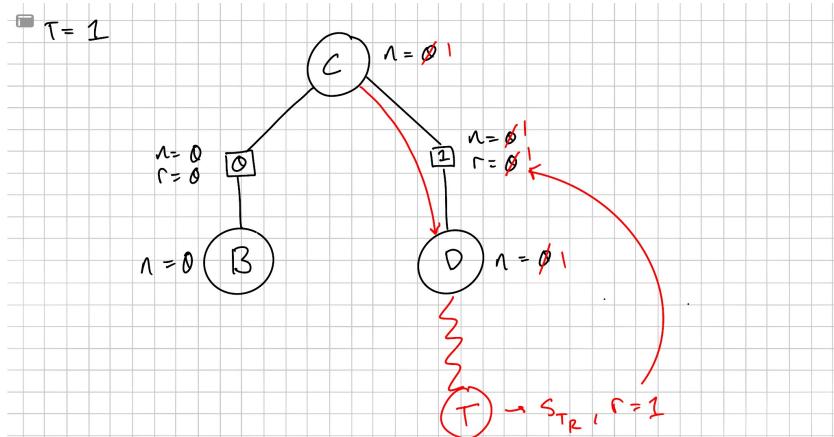


4. (a) In action selection, the bonus is ephemeral; that is: it is used only after a one-time computation and then thrown out. It does not integrate into the value function.  
 When the bonus is applied to the Q-function, it becomes a "permanent" part of the system and can only be "erased" through more experience at that state/action pair. In this case, we are making our Q-function inaccurate, whereas in action selection we are overriding our Q-function.
- (b)
  - i. Advantage of bonus applied to action selection: Not a permanent change to our system. We aren't making our Q-function inaccurate.
  - ii. Disadvantage of bonus applied to action selection: Only doable in small-dim state spaces as we must iterate over all  $\{state, action\}$  pairs and compute the bonus for each.
  - iii. Advantage of bonus applied to Q-function: can be done for high-dim state space as we embed the bonus in the Q-value.
  - iv. Disadvantage of bonus applied to Q-function: bonus becomes embedded in the system and must be overcome through experience. This leads to taking more suboptimal actions overtime.
5. (a) In stochastic environments, the transitions for a given  $\{state, action\}$  pair are not stochastic, and our model updates would have high variance and make the planning phase very inconsistent and unreliable. Further, if a bad sample is used to update the model, then our subsequent planning phase will push our Q-function in a suboptimal direction.
- (b) To fix this, our model could be an expectation of rewards for a given  $\{state, action\}$  pair. That is: we could track  $\{s, a, s', r, t\}$ ; where  $\{s, a\}$  is the current state/action pair,  $s'$  is the experienced next state,  $r$  is the reward at that next state, and  $t$  is the number of times  $\{s', r\}$  was experienced when taking action  $a$  at state  $s$ . When we execute the planning phase, we compute the reward for the samples  $\{state, action\}$  pair as the sum of rewards experienced at each  $s'$  from that  $\{state, action\}$  pair, multiplied by the probability of entering  $s'$  when taking action  $a$  at state  $s$ .
- (c) To deal with non-stationary environments, we can integrate Dyna-Q+ with the above mentioned changes to encourage exploration.
- (d) To run the code implementation of these changes:
  - i. Follow the code setup instructions at the top of this document
  - ii. Run 'python src/planning/run\_stochastic\_dynaq.py'
- (e) We see the average number of episodes completed by timestep 5000 in the stochastic environment is much lower than in the deterministic environment (plotted on page 2). This could be due to slippage



causing the agent to take sub-optimal routes. Additionally, the policy may be accounting for the additional risk caused by stochasticity. The optimal path may bring the agent dangerously close to low-value states. Since the domain has slippage, there is the risk of accidentally entering one of these low-value states. The policy will learn to cushion itself from low-value states but in doing so take a safer, less-optimal, route.

6. Manual MCTS (pictures on pages below):

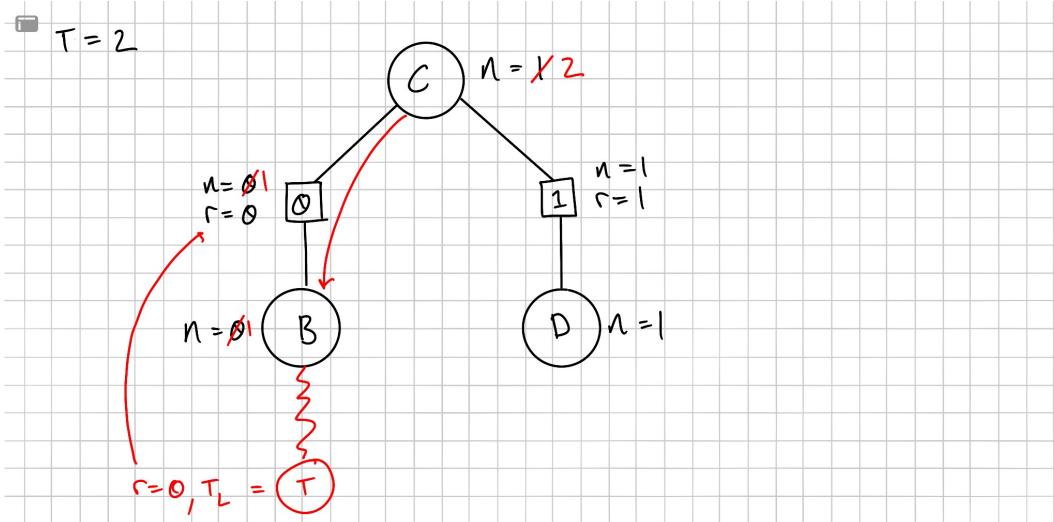


action selection

$$s_c, a_1 \rightarrow \frac{\phi}{\phi} + \sqrt{\frac{2 \ln \phi}{\phi}} = \phi \rightarrow a_1 \text{ chosen @ } s_c$$

rollout sequence:

$$\{s_D, a_1, s'_E, r=0\}, \{s_E, a_1, s'_{T_R}, r=1\}$$

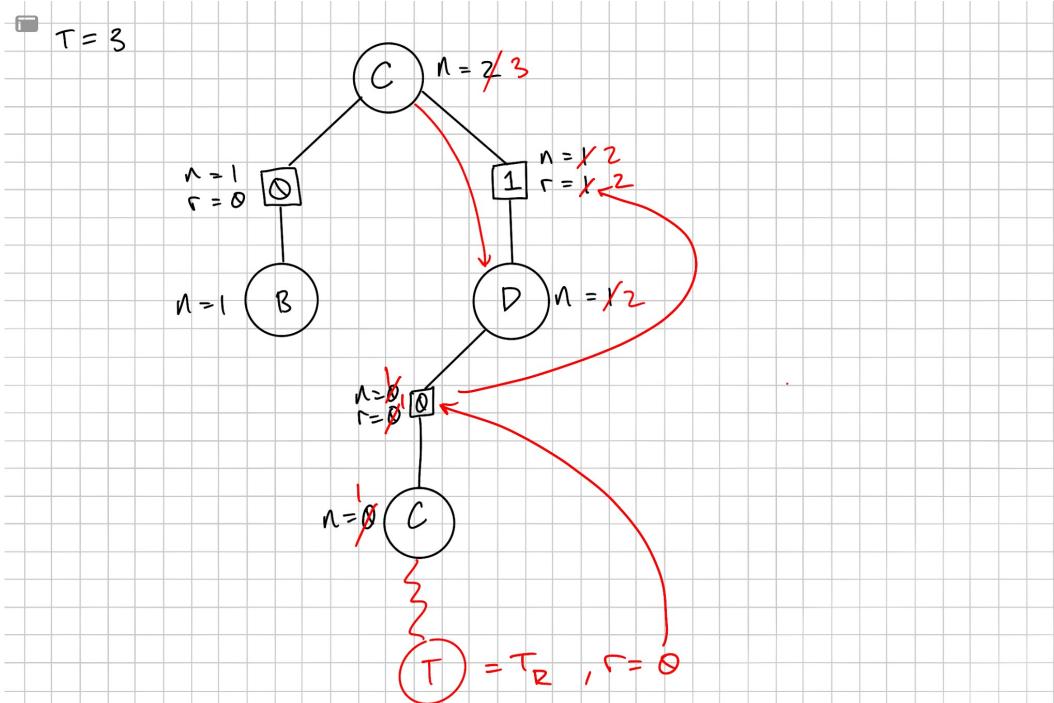


action selection:

$$s_c, a_0 = \frac{0}{0} + \sqrt{\frac{2 \ln 1}{0}} = \infty \rightarrow s_c, a_0 \text{ chosen}$$

rollout sequence:

$$\{s_B, a_0, s_A, r=0\}, \{s_A, a_0, s_{T_L}, r=0\}$$



action selection

$$s_C, a_0 = \frac{0}{1} + \sqrt{\frac{2 \ln 2}{1}} = 1.18$$

$$s_C, a_1 = \frac{1}{1} + \sqrt{\frac{2 \ln 2}{1}} = 1 + 1.18 = 2.18 \rightarrow \text{chosen action}$$

$$s_D, a_0 = \sqrt{\frac{2 \ln 1}{0}} = \infty \rightarrow \text{chosen action}$$

rollout sequence

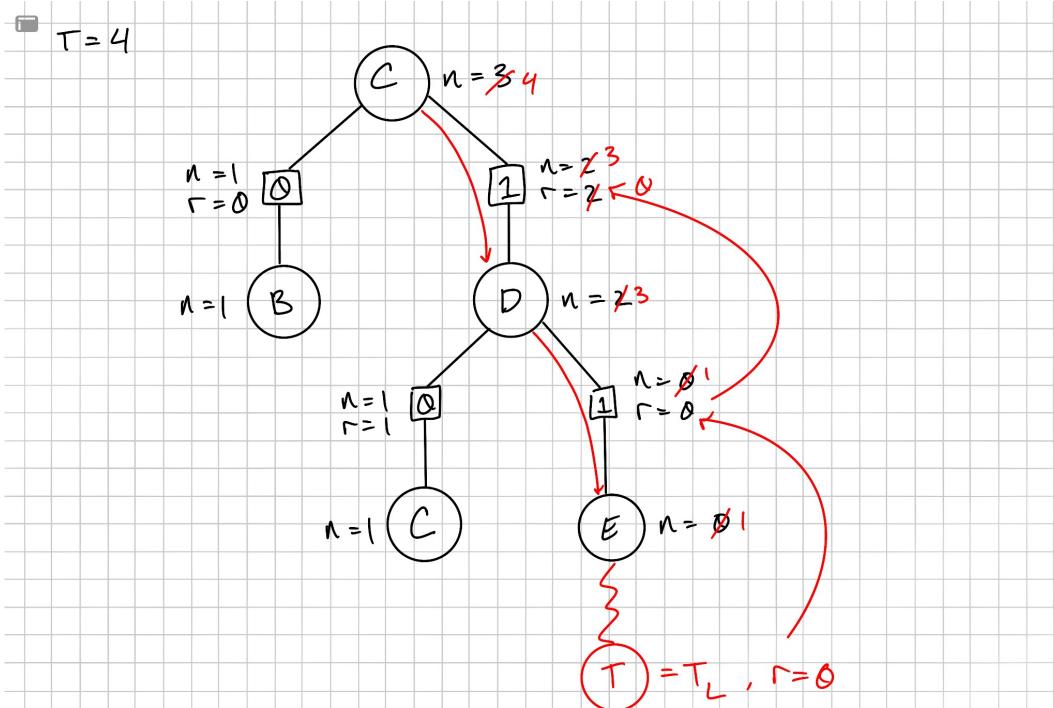
$$\{s_C, a_0, s'_B, r=0\} \quad \{s_B, a_0, s'_A, r=0\}$$

$$\{s_A, a_1, s'_B, r=0\} \quad \{s_B, a_0, s'_A, r=0\}$$

$$\{s_A, a_1, s'_B, r=0\} \quad \{s_B, a_1, s'_C, r=0\}$$

$$\{s_C, a_1, s'_D, r=0\} \quad \{s_D, a_1, s'_E, r=0\}$$

$$\{s_E, a_1, s'_{T_R}, r=0\}$$



action selection:

$$s_{C,a_0} = \frac{0}{1} + \sqrt{\frac{2 \ln 3}{1}} = 1.48$$

$$s_{C,a_1} = \frac{2}{2} + \sqrt{\frac{2 \ln 3}{2}} = 1 + 1.05 = 2.05 \rightarrow \text{chosen action}$$

$$s_{D,a_1} = \frac{0}{0} + \sqrt{\frac{2 \ln 2}{0}} = \infty \rightarrow \text{chosen action}$$

rollout sequence:

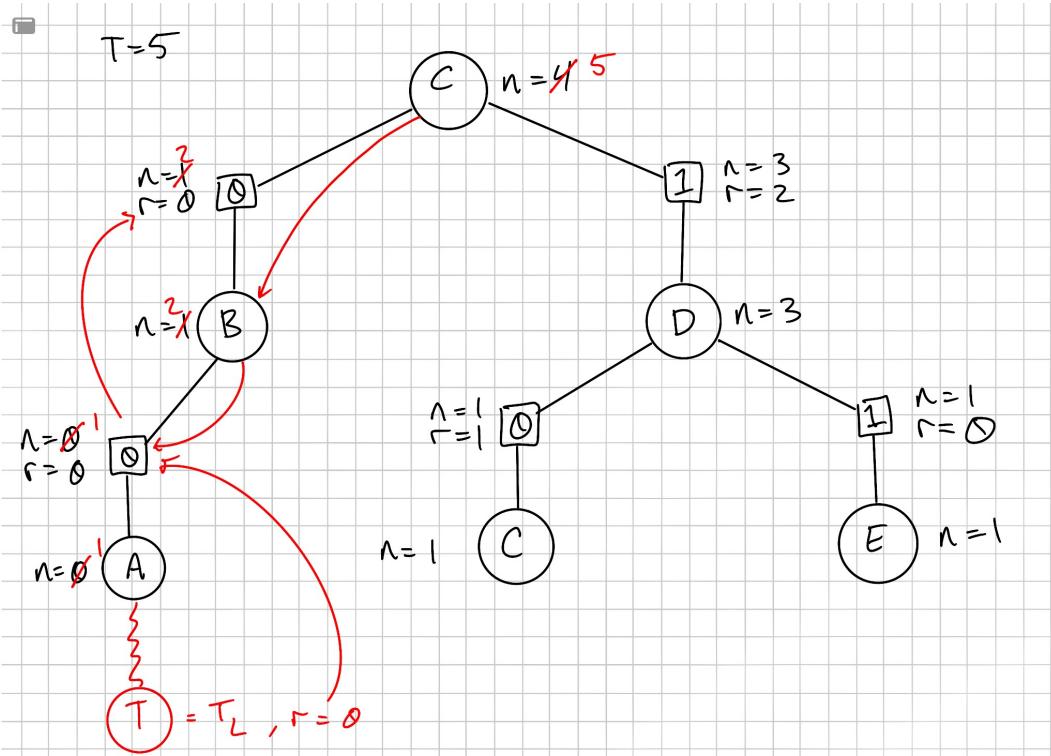
$\{s_D, a_0, s'_C\}, \{s_C, a_0, s'_B\}, \{s_B, a_1, s'_C\}, \{s_C, a_0, s'_B\}$

$\{s_B, a_1, s'_C\}, \{s_C, a_1, s'_D\}, \{s_D, a_0, s'_C\}, \{s_C, a_0, s'_B\}$

$\{s_B, a_0, s'_A\}, \{s_A, a_1, s'_B\}, \{s_B, a_1, s'_C\}, \{s_C, a_0, s'_B\}$

$\{s_B, a_0, s'_A\}, \{s_A, a_1, s'_C\}, \{s_C, a_1, s'_D\}, \{s_D, a_0, s'_C\}$

$\{s_C, a_0, s'_B\}, \{s_B, a_0, s'_A\}, \{s_A, a_0, s'_C\}, \{s_C, a_0, s'_B\}$



action selection:

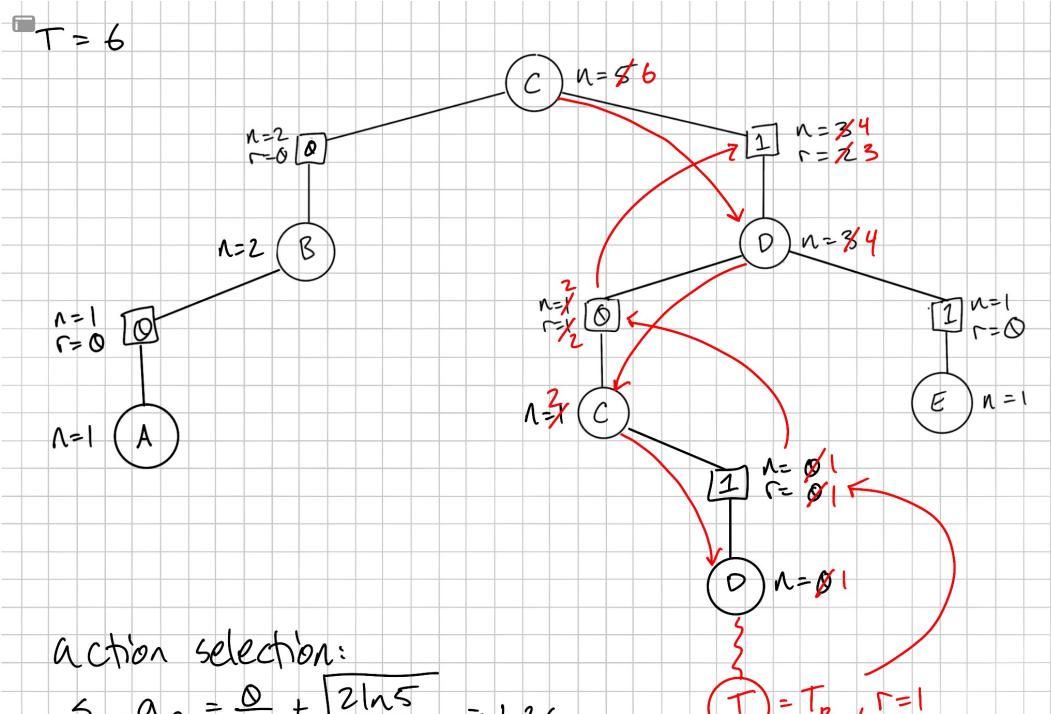
$$s_C, a_0 = \frac{0}{1} + \sqrt{\frac{2 \ln 4}{1}} = 1.66 \rightarrow \text{action chosen}$$

$$s_C, a_1 = \frac{2}{3} + \sqrt{\frac{2 \ln 4}{3}} = 0.96 + 0.66 = 1.62$$

$$s_B, a_0 = \frac{0}{0} + \sqrt{\frac{2 \ln 4}{0}} = \infty \rightarrow \text{action chosen}$$

rollout sequence:

$$\{s_A, a_0, s_{T_L}, r=0\}$$



Action selection:

$$S_C, a_0 = \frac{0}{2} + \sqrt{\frac{2 \ln 5}{2}} = 1.26$$

$$S_C, a_1 = \frac{2}{3} + \sqrt{\frac{2 \ln 5}{3}} = 1.7 \rightarrow \text{action chosen}$$

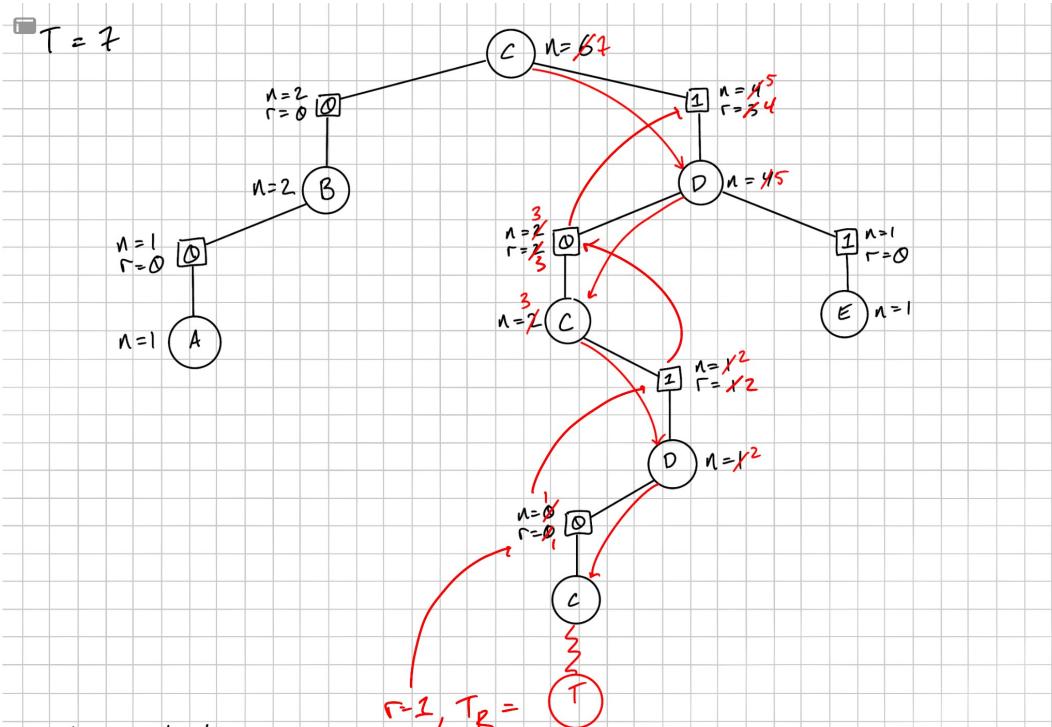
$$S_D, a_0 = \frac{1}{1} + \sqrt{\frac{2 \ln 3}{1}} = 2.5 \rightarrow \text{chosen action}$$

$$S_D, a_1 = \frac{0}{1} + \sqrt{\frac{2 \ln 3}{1}} = 1.5$$

$$S_{C'}, a_1 = \frac{0}{0} + \sqrt{\frac{2 \ln 1}{0}} = \infty \rightarrow \text{chosen action}$$

rollout sequence:

$$\{s_0, a_1, s'_E\} \{s_E, a_0, s'_D\} \{s_D, a_1, s'_E\} \{s_E, a_1, s_{T_R}, r=1\}$$



action selection:

$$s_c, a_0 = \frac{0}{2} + \sqrt{\frac{2 \ln 6}{2}} = 1.34$$

$$s_c, a_1 = \frac{3}{4} + \sqrt{\frac{2 \ln 6}{4}} = 1.7 \rightarrow \text{action chosen}$$

$$s_d, a_0 = \frac{2}{2} + \sqrt{\frac{2 \ln 4}{2}} = 2.17 \rightarrow \text{action chosen}$$

$$s_d, a_1 = \frac{0}{1} + \sqrt{\frac{2 \ln 4}{2}} = 1.17$$

$$s_{d_2}, a_0 = \frac{0}{0} + \sqrt{\frac{2 \ln 2}{0}} = \infty \rightarrow \text{action chosen}$$

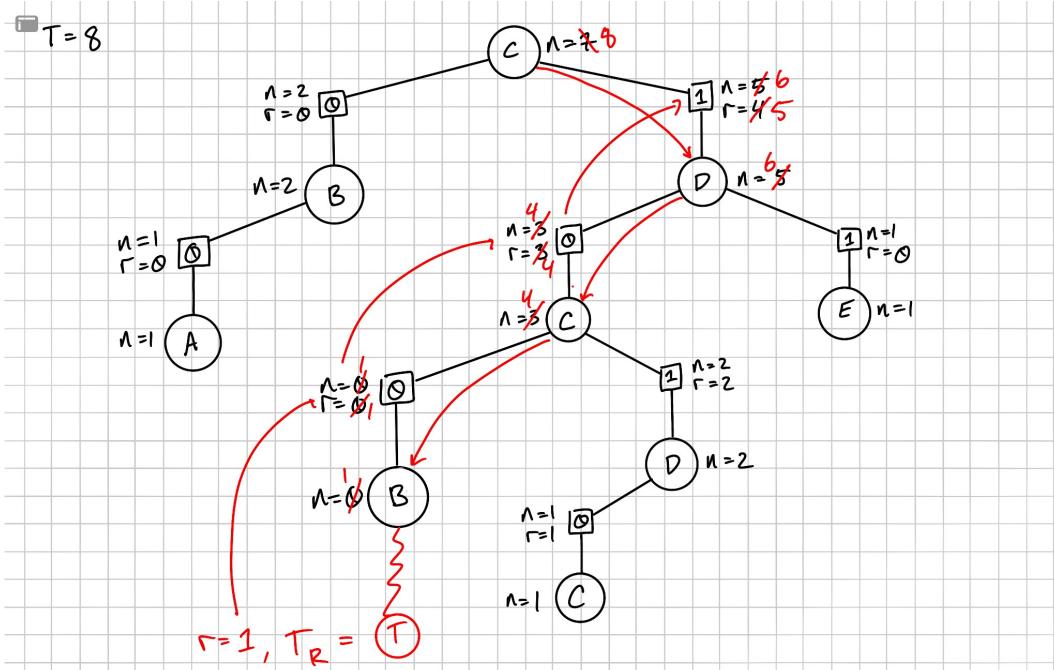
$$s_{d_2}, a_0 = \frac{0}{0} + \sqrt{\frac{2 \ln 1}{0}} = \infty \rightarrow \text{action chosen}$$

rollout sequence

$$\{s_c, a_1, s'_d\} \{s_d, a_1, s'_e\} \{s_e, a_0, s'_d\} \{s_d, a_1, s'_e\} \{s_e, a_0, s'_d\}$$

$$\{s_d, a_0, s'_c\} \{s_c, a_1, s'_b\} \{s_b, a_1, s'_c\} \{s_c, a_1, s'_d\} \{s_d, a_1, s'_e\} \{s_e, a_0, s'_d\}$$

$$\{s_d, a_0, s'_c\} \{s_c, a_1, s'_d\} \{s_d, a_1, s'_e\} \{s_e, a_1, s_{T_R}\}, r=1\}$$



action selection

$$s_c, a_0 = \frac{0}{2} + \sqrt{\frac{2 \ln 7}{2}} = 1.4$$

$$s_c, a_1 = \frac{4}{5} + \sqrt{\frac{2 \ln 7}{5}} = 1.7 \rightarrow \text{action chosen}$$

$$s_d, a_0 = \frac{3}{3} + \sqrt{\frac{2 \ln 5}{3}} = 2.03 \rightarrow \text{action chosen}$$

$$s_d, a_1 = \frac{0}{1} + \sqrt{\frac{2 \ln 5}{1}} = 1.8$$

$$s_e, a_0 = \frac{0}{0} + \sqrt{\frac{2 \ln 3}{0}} = \infty \rightarrow \text{action chosen}$$

rollout seq:

$\{s_B, a_0, s'_B\} \{s_A, a_1, s'_A\} \{s_B, a_1, s'_B\} \{s_B, a_1, s'_C\} \{s_C, a_1, s'_D\} \{s_D, a_1, s'_E\}$

$\{s_E, a_0, s'_D\} \{s_D, a_0, s'_C\} \{s_C, a_1, s'_D\} \{s_D, a_1, s'_E\} \{s_E, a_0, s'_D\}$

$\{s_D, a_0, s'_C\} \{s_C, a_0, s'_B\} \{s_B, a_1, s'_C\} \{s_C, a_0, s'_B\} \{s_B, a_1, s'_C\}$

$\{s_C, a_1, s'_D\} \{s_D, a_1, s'_E\} \{s_E, a_1, s_{TR}\}, r=13$

$T=9$  /Final Tree

