

Amina Tabassum

NUID: 002190127

Question 01:

≡

Amina Tabassum

NUID: 002190127

Q1:

- (a) NO. Because Dyna-Q will update every transition that it ever seen in planning phase at each step, but multi-step bootstrapping could only update 'n' states along trajectory. So, for an episode of length K, multistep bootstrapping will do 'K' updates whereas dyna-q will do  $Kn$  updates. i.e., K for episodes and 'n' for no. of planning steps in each episode.

Hence, if we compare n-step bootstrapping with 1-step bootstrapping, these definitely will be better performance for multistep bootstrapping. But, comparing it with Dyna-Q, it cannot reap efficiency gains that Dyna enjoys.

- (b) It might perform better, since, usually optimal step will be between 1-step method & monte-carlo method. So, if we use appropriate value of 'n', it will help algorithm converge faster. However, drawback, is that it is more computational consuming. There is not any method to select optimal value of n and it may slow down our algorithm as complexity will be  $O(n^2)$  rather than  $O(n)$ .

c)

I implemented n-step SARSA on blocking maze environment and results are shown below. It can be seen from plots that nstep SARSA does not perform as well as Dyna-Q. N-step bootstrapping is better than 1-step bootstrapping I.e, n-step SARSA is better than 1-step SARSA but when it comes to planning algorithms DynaQ, multistep bootstrapping is not better than Dyna-Q because Dyna-q performs number of episodes \* number of planning steps updates whereas multistep bootstrapping just performs number of episodes times steps updates.

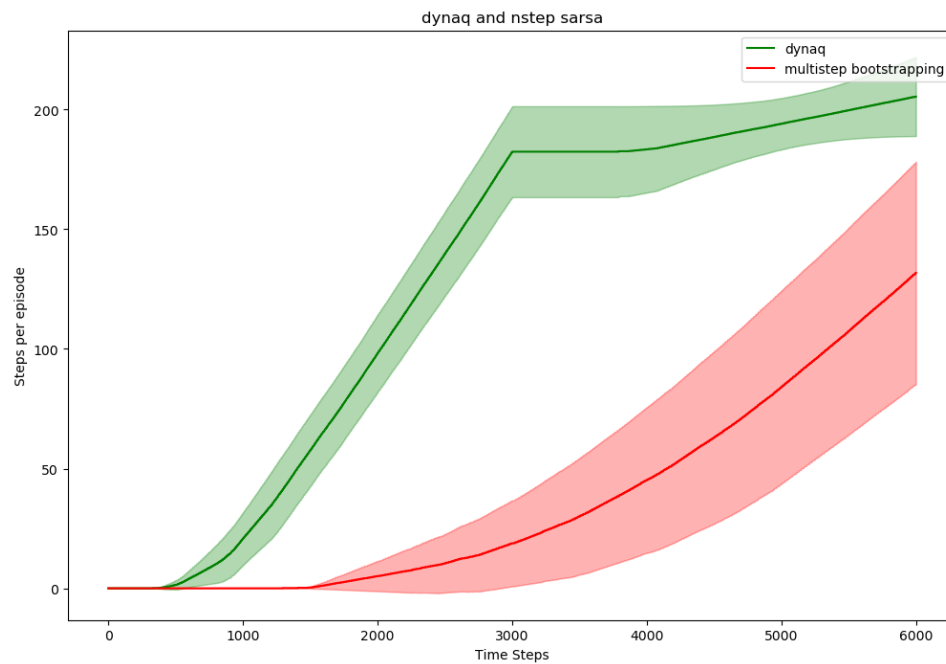


Figure1: Dyna-Q and Multistep Bootstrapping

Coming to experimental verification of part b) when n-step returns are implemented in planning step, the running complexity of algorithm increases by  $O(n^2)$  times number of episodes. In this case, algorithms perform both exploration and exploitation. It takes a lot of time in the planning phase. In this case, multistep bootstrapping updates are performed in every planning step. So, initially, the cumulative reward is slow. In other words, random sampling and multistep bootstrapping of already visited state-action pairs wastes a lot of computation on trivial updates because most of updates did not yet receive the backpropagated reward from goal state. However, when more states receive an update, the amount of wasted compute is reduced during subsequent episodes and increased timesteps. Hence, dyna-q with nstep sarsa in planning phase will outperform Dyna-q in long run but it is computationally expensive and time consuming.

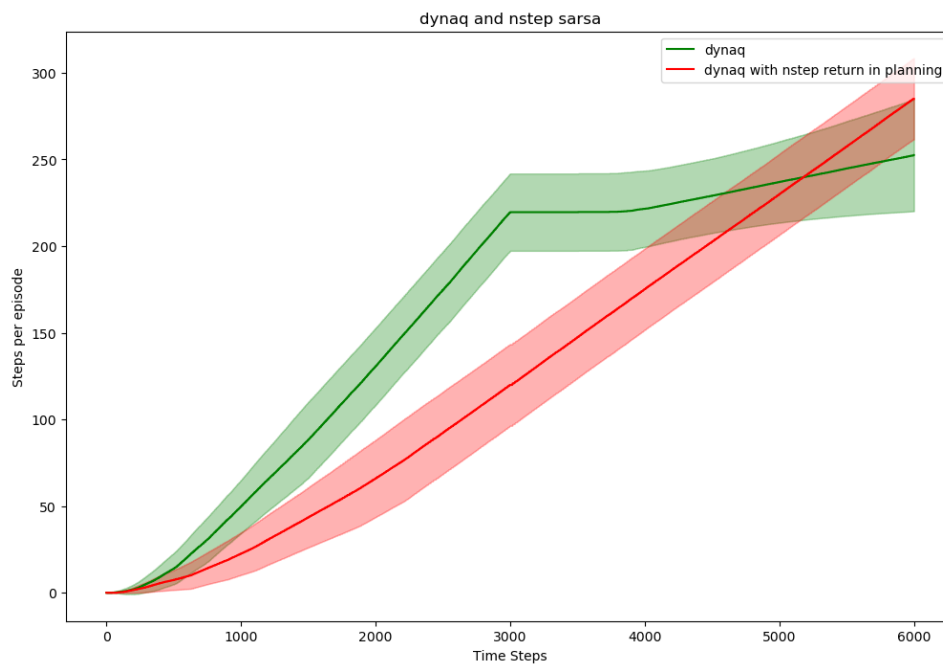


Figure 2: Dyna\_Q and Dyna-Q with n-step return in planning phase

## Question 02:

a)

Q2: (a)

Dyna-Q+ without footnote

Initialize  $Q(s,a)$  and  $\text{Model}(s,a)$  for all  $s \in S$  &  $a \in A(s)$

Loop forever:

a)  $S \leftarrow$  current (non-terminal) state

b)  $A \leftarrow \epsilon$ -greedy( $S, Q$ )

c) Take action  $A$ , observe reward  $R$  and state  $S'$

d)  $Q(S,A) = Q(S,A) + \alpha [R + \gamma \max_a Q(S',a) - Q(S,A)]$

(e)  $\text{Model}(S,A) = (R, S', \text{last-t})$

$\approx$  last-t = time step when transition occurred.

(f) Loop repeat  $n$  times:

$S \leftarrow$  randomly sample from previously observed state

$A \leftarrow$  randomly sample from previous taken actions.

$R, S', t \leftarrow \text{Model}(S,A)$  i.e. <sup>from</sup> updated model, take  $R, S'$  &  $t$  for  $S$  and  $A$ .

$\text{Reward} = \text{Reward} + K \sqrt{\text{last-t} - t}$   $\tau = \text{last-t} - t$

$\text{Reward} = \text{Reward} + K \sqrt{\tau}$

update Q-table as:

$Q(S,A) = Q(S,A) + \alpha [R + \gamma \max_a Q(S',a) - Q(S,A)]$

(b) Dyna-Q+ with footnote:

Initialize  $Q(s,a)$  and model  $\text{Model}(s,a)$  as

$Q(s,a) = 0$   $\text{Model}(s,a) = \{ \}$

env = class object

Loop forever:

a) Select current state  $S$  by  $s = \text{env.reset}$ .

b) Select action  $A$  using  $\epsilon$ -greedy policy on  $Q(s)$

as

$A = \epsilon$ -greedy-policy( $S$ )



3

c) Take action  $A$  and observe next state & reward:

$$R, s' = \text{step}(A)$$

d) update  $Q$ -table as:

$$Q(s, A) = Q(s, A) + \alpha [R + \gamma \max_a Q(s', a) - Q(s, A)]$$

e) update model:

$$\text{Model}(s, A) = (R, s', t')$$

f) loop for  $n$  times as:

for  $-$  in range( $n$ ):

$s$  = randomly selected state from previously visited states.

$A$  = randomly selected action from previous actions.

if  $A$  not in  $\text{Model}(s, A)$ :

$$R, s', t = 0, s, 0$$

else:

$$R, s', t = \text{Model}(s, A)$$

$$R = R + K \sqrt{t' - t}$$
$$= R + K \sqrt{T}$$

update  $Q$ -table as:

g)  $Q(s, A) = Q(s, A) + \alpha [R + \gamma \max_a Q(s', a) - Q(s, A)]$

b)

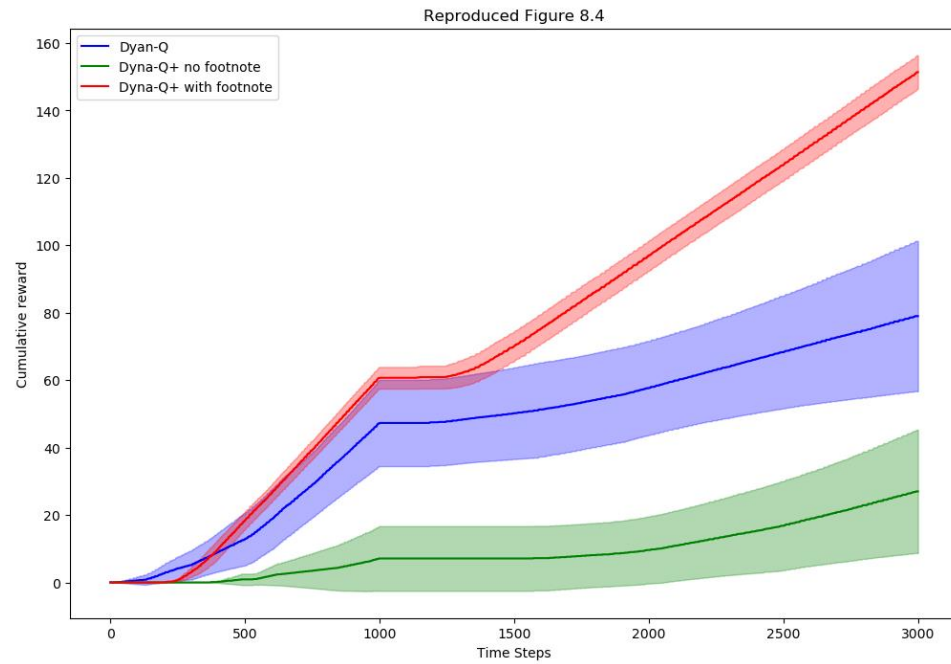


Figure 3: Dyna-Q and Dyna-Q + for Blocking Maze

In this case, both Dyna-Q and Dyna-Q+ are implemented for Blocking maze environment. It can be seen from plots that Dyna-Q+ with footnote performs better than both Dyna-Q and Dyna-Q+ without footnote. It has a high learning rate, higher cumulative reward and hence low variance.

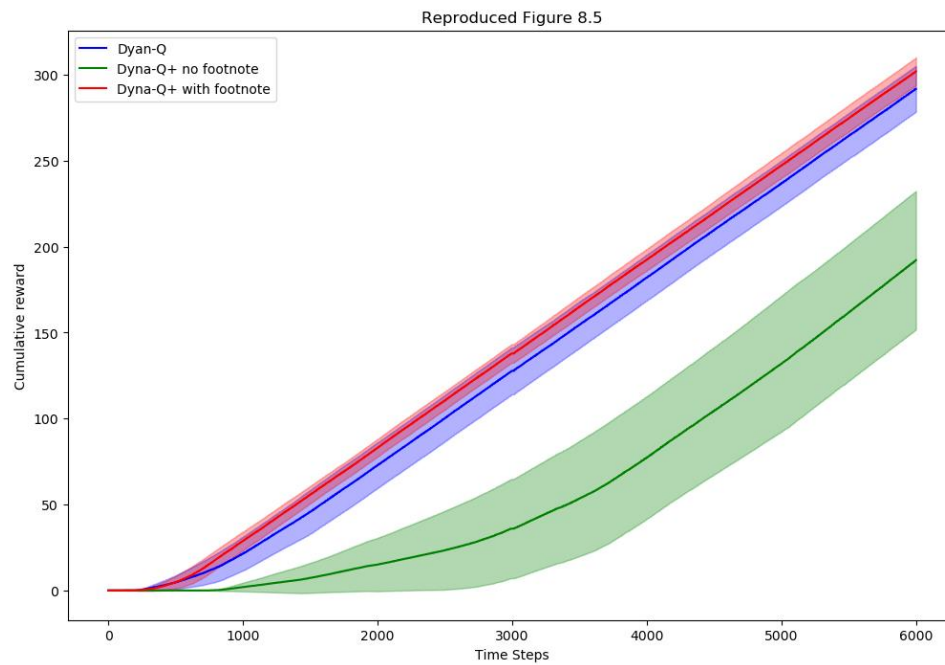


Figure 4: Dyna-Q and Dyna-Q + for Shortcut Maze

- c) The footnote does matter. It can be seen from the above results that learning curve with footnote rises rapidly and variance is small. In other words, DynaQ+ without footnote has high variance and learning curve rises very slowly. Without footnote, the actions that have not been taken will not be updates during planning phase and will appear in planning phase only when agent actually takes those actions in real environment and it depends on exploration factor epsilon which makes learning rate slow in new environments.

### Question No: 03

#### Question NO:03

In both methods, extra bonus either applied during action selection or to bonus reward will encourage the agent to choose those actions that have not been chosen for long time but UCB is used for k-arm bandit problem and it does not have transitions between different states, and it only need to pick action once to know how good is action now.

In Dyna-Q+, the bonus has actually been used in action value computation. If in the planning phase, the update time is large enough, the action value will be infinity, given learning target  $R + \text{Bonus} + \gamma \max_a Q(s', a)$ .

#### Dyna-Q+:

##### Pros:

Bootstrap update based on following state could handle state transition better than UCB.

##### Cons:

Extra bonus reward makes all action values keep increasing, which makes action value deviate from true value & not realistic.



UCB:

Advantages:

Since, we only need to record time that every state-action pair last appears.

Disadvantages:

Since, it does not update action value based on next state it reaches to, it cannot handle state transition well.

## Question No: 04

a)

Question No: 04  
Initialize  $Q(s,a)=0$ ,  $\text{Model}(s,a)=\{\}$   $\forall s \in S \ \& \ a \in A(s)$

- Select current state  $S$ :  
 $S = \text{reset}()$
- Select action using  $\epsilon$ -greedy policy:  
 $A = \epsilon\text{-greedy}[Q(S)]$
- Take action:  $R, S' = \text{step}(A)$   
Observe next-state reward  $R$  and next step  $S'$ .
- $Q(S,A) = Q(S,A) + \alpha [R + \gamma \max_a Q(S',a) - Q(S,A)]$
- $\text{model}(S,A) = [(R,S')]$   $\leftarrow$  if 'A' has not been previously taken.  
i.e, if  $A$  not in  $\text{model}[S].\text{keys}()$ :  
 $\text{model}[S][a] = [(R,S')]$   
else:  
 $\text{model}[S][a].\text{append}((R,S'))$

So, basically here, I am appending reward 'R' and next state  $S'$  for every (state, action) pair in list. I will then randomly sample  $(R,S')$  for  $(S,A)$  in planning stage.

- Loop for  $n$ -times:  
 $S \leftarrow$  Randomly select from previously visited states.  
 $A \leftarrow$  Randomly select from previously visited actions  
 $R', S' \leftarrow$  Randomly sample from list of  $(R,S')$  stored for every state/action pair.
- $Q(S,A) = Q(S,A) + \alpha [R + \gamma \max_a Q(S',a) - Q(S,A)]$

So, my logic is to store <sup>next state</sup> (reward, next-state) for every (state, action) pair. So, greater this action is taken for certain state, greater will be its probability of being selected in planning stage. Since, I am randomly selecting  $(R,S')$  for  $(S,A)$ , so, it is stochastic.

b)

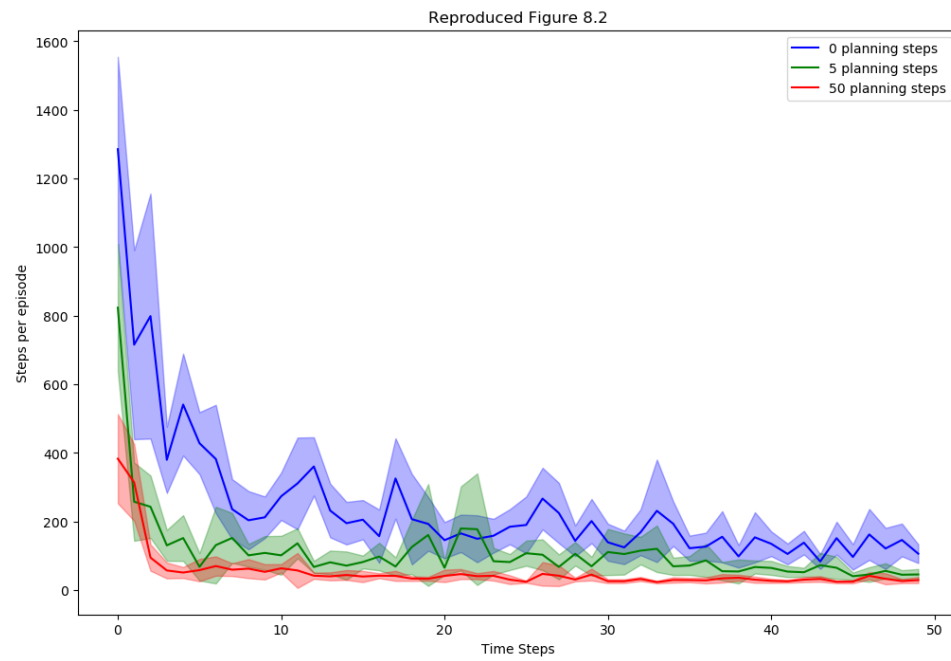


Figure 5: Dyna-Q and Dyna-Q + for Stochastic Windy Gridworld

## Question No: 05

Question No: 05

Iteration 01:

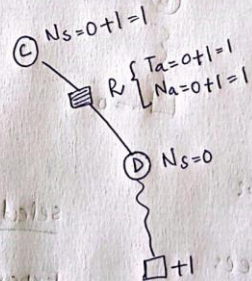
since Right = 1  
Left = 0

Selection: C

Expansion: 1 (Use random list)  $\rightarrow$  D

Simulation:  $\textcircled{C} \xrightarrow[R]{1} \textcircled{E} \xrightarrow[R]{1} \text{Terminal state}$   
Reward = +1

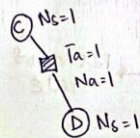
Backpropagation:



Iteration 02:

Right = 1  
Left = 0

Previous tree:



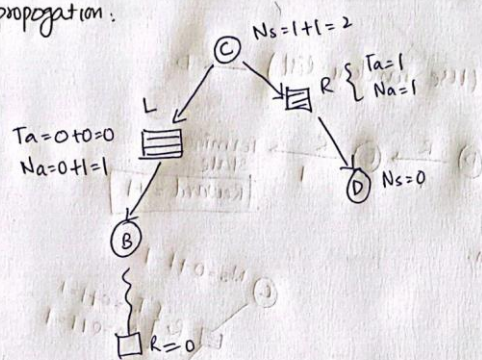


Selection: C

Expansion: 0 (no use of random list)  $\rightarrow$  B

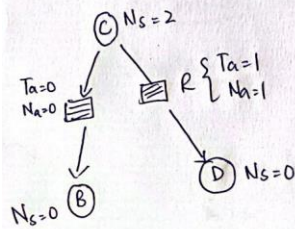
Simulation: B  $\xrightarrow{L}$  A  $\xrightarrow{L}$  Terminal state

Backpropagation:



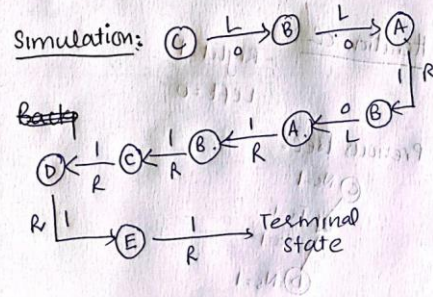
Iteration 3:

Previous tree:



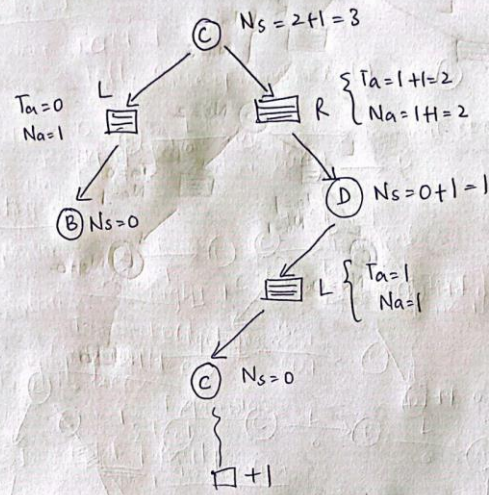
selection: C  $\xrightarrow{R}$  D

Expansion: 0  $\xrightarrow{\text{Use random list}}$  C

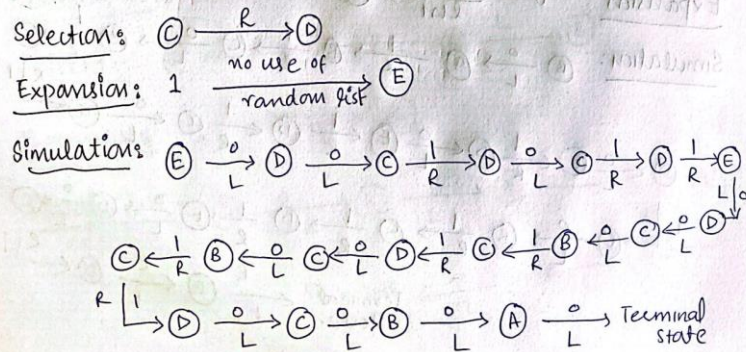




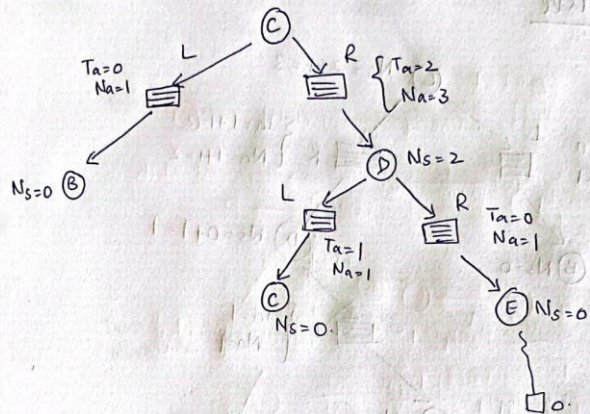
### Final tree



### Iteration 04:



### Final tree



### Iteration 05

Selection  $C \xrightarrow{R, 1} D \xrightarrow{L, 0} C$

Expansion  $C \xrightarrow{\text{Use random list}} B$

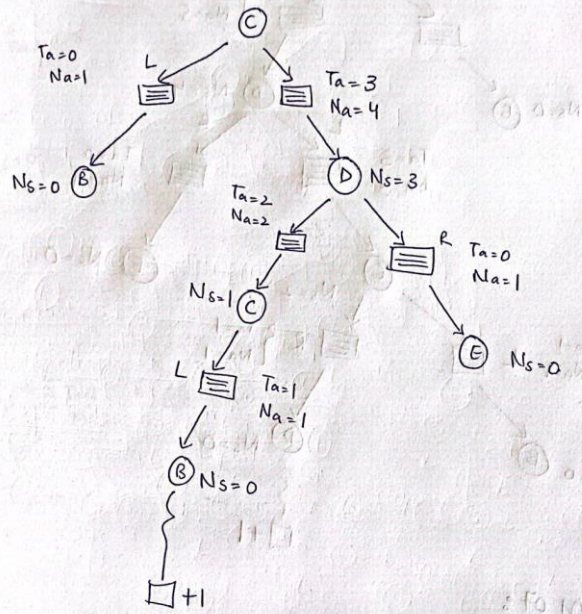
Simulation:  $B \xrightarrow{0, L} A \xrightarrow{1, R} B \xrightarrow{1, R} C \xrightarrow{0, L} B \xrightarrow{1, R} C$

$D \xrightarrow{0, L} E \xrightarrow{1, R} D \xrightarrow{0, L} E \xrightarrow{1, R} D \xrightarrow{1, R} C \xrightarrow{0, L} D \xrightarrow{1, R} C$

$L \xrightarrow{0, L} C \xrightarrow{0, L} B \xrightarrow{1, R} C \xrightarrow{1, R} D \xrightarrow{1, R} E \xrightarrow{0, L} D \xrightarrow{0, L} C$

Terminal state.

# Final tree



## Iteration 6:

