

## Assignment 1 Report

### Q1. What methods have you tried for async DP? Compare their performance.

#### 1. Comparison of three asynchronous DP methods

In-place DP	Prioritized Sweeping	Real-time DP
<b>1144 steps</b>	<b>1288 steps</b>	<b>1567 steps</b>

- **In-place dynamic programming:** This method converged the fastest, taking only 1144 steps. It avoids additional memory use by updating values directly in place, but might update the same state multiple times.
- **Prioritized sweeping:** This method is more efficient in focusing updates on high-priority states. It balances computation and prioritization but requires maintaining a priority queue.
- **Real-time dynamic programming (RTDP):** This method took the longest, with 1567 steps. RTDP optimizes decisions in real-time but requires significant interaction with the environment, leading to slower convergence.

#### 2. Method Pseudocode

##### (a) In-place dynamic programming

```

1: Initialize values for all states arbitrarily
2: repeat
3:   Set  $\delta \leftarrow 0$ 
4:   for each state in the state space do
5:      $v \leftarrow \text{values}[\text{state}]$ 
6:     Update  $\text{values}[\text{state}] \leftarrow \max(\text{Q-value}(\text{state}, \text{action}))$  for all possible actions
7:     Update  $\delta \leftarrow \max(\delta, |v - \text{values}[\text{state}]|)$ 
8:   end for
9: until  $\delta < \text{threshold}$ 
10: for each state in the state space do
11:    $\text{policy}[\text{state}] \leftarrow \arg \max(\text{Q-value}(\text{state}, \text{action}))$ 
12: end for
```

##### (b) Prioritized sweeping

```

1: Initialize values for all states arbitrarily
2: Create an empty priority queue
3: for each state in the state space do
4:   Compute the error:  $\text{error} \leftarrow |\text{get\_state\_value}(\text{state}) - \text{values}[\text{state}]|$ 
5:   if  $\text{error} > \theta$  then
6:     Push state into priority queue with negative priority
7:   end if
8: end for
9: while priority queue is not empty do
10:   if  $\text{policy\_evaluation}() < \text{threshold}$  then
11:     exit loop
12:   end if
13:   Pop the state with the highest priority from the queue
14:   Update  $\text{values}[\text{state}] \leftarrow \text{get\_state\_value}(\text{state})$ 
15:   for each predecessor of the state do
16:     Store the old value of the predecessor
```

```

17:     Update values[predecessor] ← get_state_value(predecessor)
18:     Compute the new error for the predecessor: pre_error ← |new_value − old_value|
19:     if pre_error >  $\theta$  then
20:         Push predecessor into the priority queue
21:     end if
22: end for
23: end while
24: Perform policy improvement
(c) Real-time dynamic programming
1: for each state in the state space do
2:     if policy_evaluation() < threshold then
3:         exit loop
4:     end if
5:     while True do
6:         Select action: action ← arg max(Q-value(state, action))
7:         Update values[state] ← get_state_value(state)
8:         Take a step: (next_state, reward, done) ← grid_world.step(state, action)
9:         if done then
10:            break
11:        end if
12:        for each predecessor of the state do
13:            Update values[predecessor] ← get_state_value(predecessor)
14:        end for
15:        state ← next_state
16:    end while
17: end for
18: policy_improvement()

```

**Q2. What is your final method? How is it better than other methods you’ve tried?**

1. My final method is In-place Dynamic Programming. This method is simple since it avoids the need for extra data structures like priority queues and that for tracking predecessors as well as future states continuously during interaction with the environment. While this method directly updates the value of each state based on the immediate best action for that state, it quickly reduces the difference between continuous value estimates. With all of the reasons mentioned, In-place Dynamic Programming is better than other methods I’ve tried.
2. Tried a novel method aside from the 3 method mentioned in class, but the step number 2193 is much larger than the others.
  - 1: Sort the states based on values in descending order
  - 2: **for** each sorted state **do**
  - 3: **if** policy\_evaluation() < threshold **then**
  - 4: **exit loop**
  - 5: **end if**
  - 6: **while** True **do**
  - 7: Select action: action ← arg max(Q-value(state, action))
  - 8: Update values[state] ← get\_state\_value(state)
  - 9: Take a step: (next\_state, reward, done) ← grid\_world.step(state, action)
  - 10: **if** done **then**
  - 11: **break**
  - 12: **end if**
  - 13: **for** each predecessor of the state **do**
  - 14: Update values[predecessor] ← get\_state\_value(predecessor)
  - 15: **end for**
  - 16: state ← next\_state

```
17:   end while  
18: end for  
19: policy_improvement()
```