

RL_HW1_gridworld

ShengHsun Chang

September 29, 2025

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

1.1 Summary

I've tried three asynchronous DP methods mentioned in class: In-place DP, Prioritized Sweeping, and Real-Time DP. The implementation details and performance are discussed in the following subsections.

1.2 Method1: In-place DP

I implement In-place DP base on value iteration method. The difference is that I update the value function $V(s)$ right after computing the new value, instead of waiting until the end of the episode. It is better than the original value iteration because it uses the up-to-date value. However, it still waste steps in states with little residual. The pseudo code is shown below.

Algorithm 1: In-place Dynamic Programming

Input : GridWorld \mathcal{G} , discount γ ; threshold θ
Output : Updated $V(\cdot)$ and greedy $\pi(\cdot)$

```
1 Initialization;  
2  $V(s) \leftarrow 0$  for all  $s \in S$  ; // initial value  
3  $\Delta \leftarrow 2\theta$  ; // ensure at least one sweep  
4 while  $\Delta > \theta$  do  
5    $\Delta_{max} \leftarrow 0$ ;  
6   foreach  $s \in S$  in a fixed order do  
7      $v_{old} \leftarrow V(s)$ ;  
8      $q_{max} \leftarrow -\infty$ ;  
9     foreach  $a \in \{0, \dots, |A| - 1\}$  do  
10       $(s', r, done) \leftarrow \text{step}(s, a)$ ;  
11      if done then  $q \leftarrow r$  ;  
12      else  $q \leftarrow r + \gamma \cdot V(s')$  ;  
13       $q_{max} \leftarrow \max(q_{max}, q)$ ;  
14       $V(s) \leftarrow q_{max}$  ; // in-place overwrite  
15       $\Delta_{max} \leftarrow \max(\Delta_{max}, |V(s) - v_{old}|)$ ;  
16 foreach  $s \in S$  do  
17    $\pi(s) \leftarrow \arg \max_a \{r(s, a) + \gamma \cdot V(s')\}$ ;
```

1.3 Method2: Prioritized Sweeping

Prioritized Sweeping improves In-place DP by choosing the state with the largest residual to update. Right after updating a state, we push its predecessors back into an array(priority queue) In my implementtation, I use a simple list instead of a priority queue for:

- Simplicity: Python’s built-in list is easy to use and understand.
- Duplication: The same state would be pushed multiple times, and the checking is too annoying.
- Efficiency: The performance in the test case is still too small.

It is worth noting that I have to initialize the value function $V(s)$ to $-\infty$ instead of 0. I think it is because the first residual can thus be quite large, and the update can propagate faster.

Algorithm 2: Prioritized Sweeping

Input : GridWorld \mathcal{G} , discount γ ; threshold θ

Output : Updated $V(\cdot)$ and greedy $\pi(\cdot)$

```

1 Initialization;
2  $V(s) \leftarrow -\infty$  for all  $s \in S$ ;
3  $\text{Pred}[s] \leftarrow \emptyset$  for all  $s \in S$ ;
4 Compute initial residual $[s]$  for all  $s$ ;
5 while some residual $[s] \geq \theta$  do
6    $s^* \leftarrow \arg \max_s \text{residual}[s]$ ;
7   Recompute  $V(s^*)$  by one-step backup;;
8    $V(s^*) \leftarrow \max_a \{r(s, a) + \gamma \cdot V(s')\}$ ;
9   residual $[s^*] \leftarrow 0$ ;
10  foreach  $p \in \text{Pred}[s^*]$  do
11    Recompute backup value for  $p$ ;
12    residual $[p] \leftarrow |\text{backup}(p) - V(p)|$ ;

13 foreach  $s \in S$  do
14    $\pi(s) \leftarrow \arg \max_a \{r(s, a) + \gamma \cdot V(s')\}$ ;
```

1.4 Method3: Real-Time DP

Real-Time DP, different from the previous two methods, updates value function only along the trajectory generated by the current policy.

Therefore, to ensure all state policy are updated, I choose all states as the start state in a random order in each episode.

Despite these modifications, the performance is still bad compared to Prioritized Sweeping.

Algorithm 3: Real-Time Dynamic Programming

Input : GridWorld \mathcal{G} , discount γ ; threshold θ

Output : Updated $V(\cdot)$ and greedy $\pi(\cdot)$

```
1 Initialization;
2  $V(s) \leftarrow 0$  for all  $s \in S$ ; // initial optimistic values
3 LegalActions[ $s$ ]  $\leftarrow$  all actions for each  $s$ ;
4 Unsolved  $\leftarrow S$ ; // all states initially unsolved
5 while not converged do
6   Choose a random permutation of unsolved states;
7   foreach  $s \in \text{permutation}$  do
8     if  $s$  is terminal then
9       continue
10     $steps \leftarrow 0$ ,  $max\_steps \leftarrow 4|S|$ ;
11    while  $steps < max\_steps$  do
12       $v_{new} \leftarrow \max_a \{r(s, a) + \gamma \cdot V(s')\}$ ;
13       $diff \leftarrow |v_{new} - V(s)|$ ;
14       $V(s) \leftarrow v_{new}$ ;
15      if  $diff < \theta$  then
16        remove  $s$  from Unsolved
17       $a^* \leftarrow \arg \max_a \{r(s, a) + \gamma \cdot V(s')\}$ ;
18       $\pi(s) \leftarrow a^*$ ;
19       $(s', r, done) \leftarrow \text{step}(s, a^*)$ ;
20      if  $done$  then
21        break
22      if  $s' = s$  then
23        remove  $a^*$  from LegalActions[ $s$ ]; break
24       $s \leftarrow s'$ ,  $steps \leftarrow steps + 1$ ;
```

2 What is your final method? How is it better than other methods you've tried?

2.1 Concept

In my asynchronous DP implementation, I find out that Prioritized Sweeping is the best method.

So, I implement an optimization based on Prioritized Sweeping.

In the current maze, if an action leads to the same state(a.k.a hit the wall), then this action is useless.

Therefore, I modify the Bellman equation to ignore these actions.

Algorithm 4: Optimized Prioritized Sweeping (Async DP)

Input : GridWorld \mathcal{G} , discount γ ; threshold θ

Output : Updated $V(\cdot)$ and greedy $\pi(\cdot)$

```
1 Initialization;
2  $V(s) \leftarrow -\infty$  for all  $s \in S$ ;
3  $\text{Pred}[s] \leftarrow \emptyset$  for all  $s \in S$ ;
4  $\text{LegalActions}[s] \leftarrow$  all actions for each  $s$ ;
5 foreach  $s \in S$  do
6    $q_{\max} \leftarrow -\infty$ ;
7   foreach  $a \in A$  do
8      $(s', r, \text{done}) \leftarrow \text{step}(s, a)$ ;
9     if  $s' = s$  then remove  $a$  from  $\text{LegalActions}[s]$  ;
10    if  $\text{done}$  then  $q \leftarrow r$  ;
11    else  $q \leftarrow r + \gamma \cdot V(s')$  ;
12     $q_{\max} \leftarrow \max(q_{\max}, q)$ ;
13    if  $s' \neq s$  then
14       $\text{add } s \text{ to } \text{Pred}[s']$ 
15   $\text{residual}[s] \leftarrow |q_{\max} - V(s)|$ ;
16   $\text{new\_values}[s] \leftarrow q_{\max}$ ;
17 while true do
18    $s^* \leftarrow \arg \max_s \text{residual}[s]$ ;
19   if  $\text{residual}[s^*] < \theta$  then
20     break
21    $v_{\text{new}} \leftarrow \text{get\_state\_value}(s^*)$ ;
22    $\text{err} \leftarrow |v_{\text{new}} - V(s^*)|$ ;
23   if  $\text{err} < \theta$  then
24      $\text{residual}[s^*] \leftarrow 0$ ; continue
25    $V(s^*) \leftarrow v_{\text{new}}$ ;
26    $\text{residual}[s^*] \leftarrow 0$ ;
27   foreach  $p \in \text{Pred}[s^*]$  do
28      $v_{\text{pred}} \leftarrow \text{get\_state\_value}(p)$ ;
29      $e_p \leftarrow |v_{\text{pred}} - V(p)|$ ;
30     if  $e_p > \theta$  then
31        $\text{new\_values}[p] \leftarrow v_{\text{pred}}$ ;
32        $\text{residual}[p] \leftarrow e_p$ ;
33 foreach  $s \in S$  do
34   if  $\text{LegalActions}[s] = \emptyset$  then
35     continue
36    $\pi(s) \leftarrow \arg \max_{a \in \text{LegalActions}[s]} \{r(s, a) + \gamma \cdot V(s')\}$ ;
```

3 Experiment Results & Discussion

3.1 Experiment Setup

I choose the sample map, and three gridworld map(10x10, 15x15, 20x20) generated by chat-GPT with different sizes. the maps are shown below.

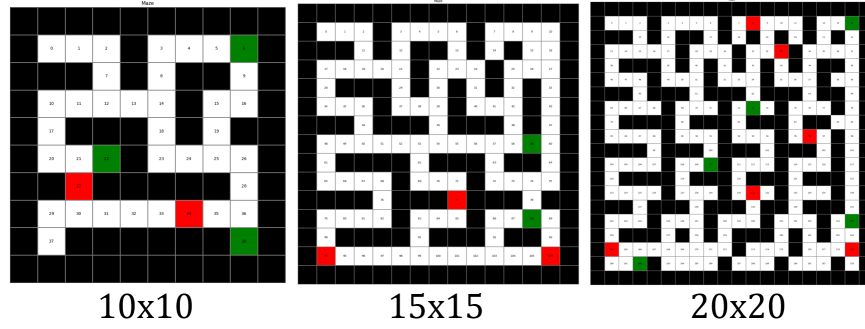


Figure 1: from left to right: 10x10 map, 15x15 map, 20x20 map

3.2 Experiment Results

To begin with, for the sample map, All the methods can converge to the optimal policy. like the figure below.

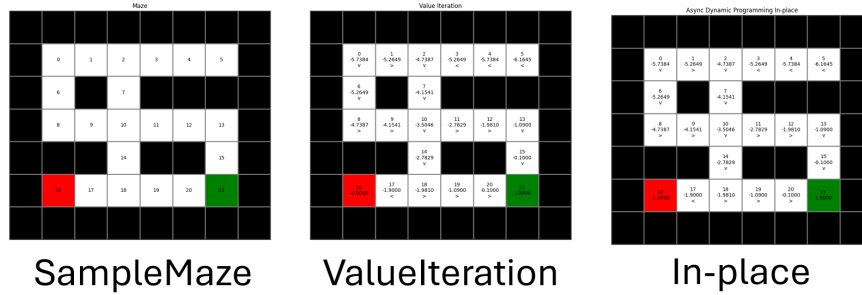


Figure 2: Optimal policy for sample map_1

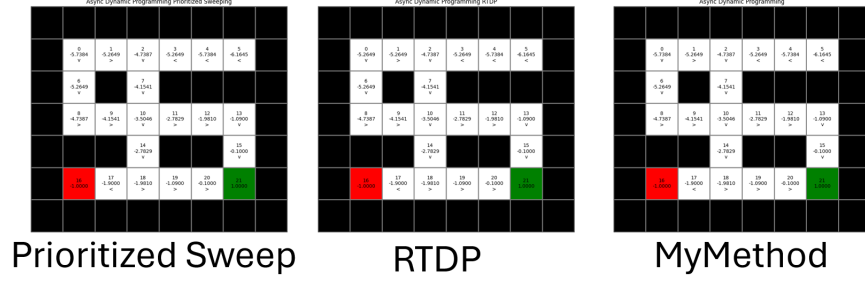


Figure 3: Optimal policy for sample map_2

Method	Sample Map	10x10 Map	15x15 Map	20x20 Map
Value Iteration(base line)	1144	2028	8560	20952
In-place DP	1056	1716	8560	16296
Prioritized Sweeping	692	988	3044	5604
Real-Time DP	2243	2904	13386	26677
Optimized PS	469	690	2078	3850

Table 1: Required steps for convergence across different gridworld maps

3.3 Discussion

The Optimized method I propose consistently outperforms the other methods as shown in Table 1.

The reasons are as follows:

- The prioritized sweeping mechanism ensures that states with the highest residuals are updated first, accelerating the propagation of value updates through the state space.
- By ignoring useless actions, the Bellman update focuses on more promising actions, leading to faster convergence.
- The optimization reduces the number of states by 30-40% in all test cases, compared with the original Prioritized Sweeping method.