

RL HW1 Report

Yu-Xiang, Luo

April 1, 2024

Q1

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

Unfortunately, I did not come out a way to implement prioritized sweeping and real-time DP. I can only perform the in-place DP.

- in-place DP: **Iteration steps: 968**, the value iteration sync DP is **1056** steps
- Prioritized sweeping: No implementation
- Real-time DP: No implementation

The reason why in-place DP is effective: Let's consider a simple scenario: a $n \times 1$ grid world, with S_{n-1} be goal state.

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

step_reward = -1, goal_reward = 10

Value Iteration

Solved in 288 steps

Start state: 0, End state: 7

Async Dynamic Programming

Solved in 64 steps

Start state: 0, End state: 7

The async DP use `[i for i in range(7, -1, -1)]` as Iteration order

This is a simple case, but I think is enough to explain my thought. Initially, $\forall v \in V, v = 0$, which means in the first iteration, only $V(7)$ is changed since `step(7, action).reward = 1`. While, at the second iteration, sync DP will only change the value of $V(6)$ since only $V(7)$ change. Then, $V(5)$ and $V(6)$, during next iterations. In other words, The reward of $V(7)$ needs at least 7 iteratios to pass through $V(0)$, yet in in-place DP, the best traverse order(`for i in range(7, -1, -1)`) can pass $V(7)$ reward to $V(0)$ in one iteration. Therefore, if we randomize the traversal(to avoid always meet the worst case, like start state is 0 and goal state is $n - 1$), then the expected number of update would be greater than sync DP's, **or at least as good as sync DP's**.(Sometimes it doesn't improve because the reward isn't good enough for the agent. Like the default reward for our HW1, `step_reward = -1, goal_reward = 1` is not enough)

Value Iteration

Solved in 1056 steps

Start state: 0, End state: 21

Async Dynamic Programming

Solved in 616 steps

Start state: 0, End state: 21

Randomize the order of traversal in every iteration, `goal_reward = 100`

I do try this method on the maze TA offer, and the result is as above.

Q2

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

Since we're allowed to store all step we've travelled. That is, we can store all `grid_world.step(state, action)`. Then, the worst iteration steps would be $4 * state_space$.

3 replies

**TA 黃柏睿 Bo-Ruei Huang B09901171** 3 days ago
Model的理解沒有錯，Model-based RL也是一個以提升sample efficiency為目標的RL主流分支，所以不會禁止同學建Model

**洪國瀚 Johnson Hung B09902120** 1 day ago
助教你好，想請問所以把走過的(s, a) -> (s', r) 存起來並在下一次遇到(s, a)時使用是合法的嗎？

**TA 黃柏睿 Bo-Ruei Huang B09901171** 1 day ago
是

 1 

By little improvement, step on goal state and trap state can be reduced to 1.