

Planning, Learning and Acting





Rupam Mahmood

March 9, 2020



First: Summary of Part 1



Learning and planning

Model of the environment: used by agent to predict environment's response to actions

Planning: a process that takes model as an input and produces/improves a policy Dynamic programming methods are planning methods which is used to produce/improve a policy



Model learning: the process of using real experience to improve the model

For example, if we are estimating the transition probabilities based on real experience

For example, TD or MC methods for prediction and control

Indirect reinforcement learning: the process of using real experience to improve the value function and policy using the model through planning

- For example, if the agent has the transition probabilities p(s',r|s,a) as in dynamic programming
- A planning method uses the model to simulate the environment to produce simulated experience

values _____ policy

- **Direct reinforcement learning:** the process of using real experience to directly improve the value function and policy

An architecture for planning, learning and acting at the same time

quiz: indicate indirect RL with an arrow (which of the two red ones?)



Tabular Dyna-Q

Initialize Q(s, a) and Model(s, a) for all $s \in S$ and $a \in A(s)$ Loop forever:

- (a) $S \leftarrow \text{current}$ (nonterminal) state
- (b) $A \leftarrow \varepsilon$ -greedy(S, Q)

(c) Take action A; observe resultant reward, R, and state, S'(d) $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment) (f) Loop repeat n times:

 $S \leftarrow$ random previously observed state $A \leftarrow$ random action previously taken in S $R, S' \leftarrow Model(S, A)$

Which part is direct RL? What is this method called?

Which part is model learning?

Which part is planning?

 $Q(S, A) \leftarrow Q(S, A) + \alpha \left[R + \gamma \max_{a} Q(S', a) - Q(S, A) \right]$





How many planning steps should we take?

Tabular Dyna-Q

Initialize Q(s, a) and Model(s, a) for all s Loop forever:

- (a) $S \leftarrow \text{current (nonterminal) state}$
- (b) $A \leftarrow \varepsilon$ -greedy(S, Q)
- (c) Take action A; observe resultant reward, R, and state, S'(d) $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment) (f) Loop repeat n times:

 $S \leftarrow$ random previously observed state $A \leftarrow$ random action previously taken in S $R, S' \leftarrow Model(S, A)$

 $Q(S,A) \leftarrow Q(S,A) + \alpha \left[R + \gamma \max_{a} Q(S',a) - Q(S,A) \right]$

$$\in S$$
 and $a \in \mathcal{A}(s)$

Dyna-Q+

 $\tau(s, a)$ denotes the number of time steps (s, a) has not been tried

2. Actions that have not been tried from a previously visited state are allowed to be considered in planning Where would you put these steps in Dyna-Q to get Dyna-Q+?

Tabular Dyna-Q

Initialize Q(s, a) and Model(s, a) for all $s \in S$ and $a \in A(s)$ Loop forever:

- (a) $S \leftarrow \text{current}$ (nonterminal) state (b) $A \leftarrow \varepsilon$ -greedy(S, Q)
- (c) Take action A; observe resultant reward, R, and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) Q(S, A)]$ (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)

(f) Loop repeat n times:

 $S \leftarrow$ random previously observed state $A \leftarrow$ random action previously taken in S $R, S' \leftarrow Model(S, A)$

1. Adds a bonus $\kappa \sqrt{\tau(s, a)}$ to reward in planning

 $Q(S, A) \leftarrow Q(S, A) + \alpha \left[R + \gamma \max_{a} Q(S', a) - Q(S, A) \right]$

Dyna-Q+: calculating visitation counts

Consider an MDP with one actions (L) and two states with the following episode



Calculate $\tau(s, a)$ for all state-action pairs at each step