

Dynamic Programming





Rupam Mahmood







Estimates like *V* or *Q*

 $\pi(a|s)$





$$, R_t = r \mid S_{t-1} = s, A_{t-1} = a \}$$

Dynamic programming

- Dynamic programming (DP) is a way of knowing values and optimal policies when the model of the environment's dynamics is given
- ✓ If in a problem the model is given to the agent, it can use DP
- Often the agent is only given an estimate of the model if at all or ...
- The agent estimates the model through experience (model learning)
- Or a designer tests their agent by comparing agent's performance against true values / optimal policies in a toy MDP
- Generally, mechanisms that use a given model as opposed to experience to improve performance are known as planning methods

Two key ideas underlying DP methods: (1) Fixed-point iteration

Bellman equation

$$v_{\pi}(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a)[r + \gamma v_{\pi}]$$

$$v_{k+1}(s) \doteq \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a)[r + \gamma v_{k}]$$

Bellman update

(s')]

Fixed-point equation x = f(x)

(s')]

 $x_{k+1} \doteq f(x_k)$

Fixed-point iteration

Converges under some conditions

Iterative policy evaluation

Iterative Policy Evaluation, for estimating $V \approx v_{\pi}$

Input π , the policy to be evaluated Initialize V(s), for all $s \in S^+$, arbitrarily except that V(terminal) = 0Loop: $\Delta \leftarrow 0$ Loop for each $s \in S$: $v \leftarrow V(s)$ $V(s) \leftarrow \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,$ $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ until $\Delta < \theta$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation

$$a) \left[r + \gamma V(s') \right]$$

Two key ideas underlying DP methods: (2) Policy improvement

If for a policy π' we have $q_{\pi}(s, \pi'(s)) \ge v_{\pi}(s) \forall s$, then we also have $v_{\pi'}(s) \ge v_{\pi}(s), \forall s$

A greedy policy w.r.t. q_{π} is such a policy

 $\max_{a} q_{\pi}(s, a) \ge v_{\pi}(s) \forall s$

Policy iteration

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \cdots \xrightarrow{I} \pi_* \xrightarrow{E} v_*$$

- 1. Initialization $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in S$
- 2. Policy Evaluation Loop: $\Delta \leftarrow 0$ Loop for each $s \in S$: $v \leftarrow V(s)$ $V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r+\gamma V(s')]$ $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
- 3. Policy Improvement policy-stable $\leftarrow true$ For each $s \in S$: old-action $\leftarrow \pi(s)$ $\pi(s) \leftarrow \operatorname{arg\,max}_a \sum_{s',r} p(s',r|s,a) [r+\gamma V(s')]$ If $old\text{-}action \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$ If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

Value iteration

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation Initialize V(s), for all $s \in S^+$, arbitrarily except that V(terminal) = 0Loop: $\Delta \leftarrow 0$ Loop for each $s \in S$: $v \leftarrow V(s)$ $V(s) \leftarrow \max_{a} \sum_{s',r} p(s',r|s,a) [r]$ $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ until $\Delta < \theta$ Output a deterministic policy, $\pi \approx \pi_*$, such that $\pi(s) = \operatorname{argmax}_{a} \sum_{s',r} p(s', r | s, a) \left[r + \gamma V(s') \right]$

$$v_* = \max_{a} \sum_{\substack{s',r}} p(s)$$
$$v_{k+1} \doteq \max_{a} \sum_{\substack{s',r}} p(s)$$

$$+\gamma V(s')]$$

 $r', r \mid s, a) [r + \gamma v_*(s')]$

 $(s', r \mid s, a)[r + \gamma v_k(s')]$