

Reinforcement Learning

- Sarath Chandar

A. Dynamic Programming



4. Dynamic Programming

- Collection of algorithms that can be used to compute optimal policies given a perfect model of the environment as a Markov decision process (MDP).

Limitations :- 1. Assumes a perfect model.

2. Computationally expensive.

However, it provides essential foundation for the understanding of the future algorithms.

Assume that the environment is a finite MDP.

i.e. State S , action A and reward R are finite, the dynamics are given by

$$P(s', r | s, a) \quad \forall s \in S, a \in A(s), r \in R$$

and $s' \in S^+$ (where S^+ is S plus

key idea: use of value functions to

organize and structure the

search for good policies.

a terminal state for
episodic setting)

Optimal value fns v_* , q_* satisfy the Bellman optimality equations.

$$v_*(s) = \max_a \mathbb{E} [R_{t+1} + \gamma v_*(s_{t+1}) \mid s_t=s, a_t=a]$$

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

$$q_*(s, a) = \mathbb{E} [R_{t+1} + \gamma \max_{a'} q_*(s_{t+1}, a') \mid s_t=s, a_t=a]$$

$$= \sum_{s', r} p(s', r \mid s, a) [r + \gamma \max_{a'} q_*(s', a')]$$

DP algorithms are obtained by turning Bellman Equations into update rules for improving approximations of the desired value functions.

Policy Evaluation (Prediction) :

How to Compute the state value fn v_π for an arbitrary policy π ?

$$v_\pi(s) = \mathbb{E}[G_t \mid s_t=s]$$

$$= \mathbb{E}_{\pi} \left[R_{t+1} + \gamma G_{t+1} \mid S_t = s \right]$$

$$= \mathbb{E}_{\pi} \left[R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s \right]$$

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V_{\pi}(s')] \quad \textcircled{*}$$

where $\pi(a|s)$ = prob. of taking action 'a' in state 's' under policy π .

The existence and uniqueness of V_{π} are guaranteed as long as either $\gamma < 1$ or eventual termination is guaranteed from all states under policy π .

Note: If the env's dynamics are completely known, then $\textcircled{*}$ is a system of $|S|$ simultaneous linear equations in $|S|$ unknowns. The solution is straight forward but tedious.

Iterative soln:- Consider a sequence of approximate value fns $v_0, v_1, v_2 \dots$ each mapping S^+ to \mathbb{R} .

The initial approx. is chosen arbitrarily (except that the terminal state, if any, must be given value 0).

Each successive approx. is obtained by using the Bellman equation for v_{π} as an update rule.

$$v_{k+1}(s) = \mathbb{E}[R_{t+1} + \gamma v_k(s_{t+1}) | s_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_k(s')]$$

$\forall s \in S$.

- ① $v_k = v_{\pi}$ is a fixed point for this update rule.
- ② the sequence $\{v_k\}$ can be shown to converge to v_{π} as $k \rightarrow \infty$ under the same conditions that guarantee the existence of v_{π} .

This algorithm is known as iterative policy evaluation.

These updates are expected updates because they are based on an expectation over all possible next states rather than on a sample next state.

Note 1: One could use two arrays, one for old values $V_k(s)$ and one for new values $V_{k+1}(s)$. Or one can also do in-place updates which converges faster.

Note 2: One round of update of all the states is called as 'sweep' through the state space.

Note 3: Iterative Policy evaluation only converges in limit. However, in practice we stop when the updates are sufficiently small.

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

Input π , the policy to be evaluated

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in \mathcal{S}$:

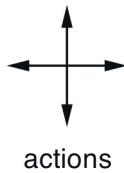
$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$

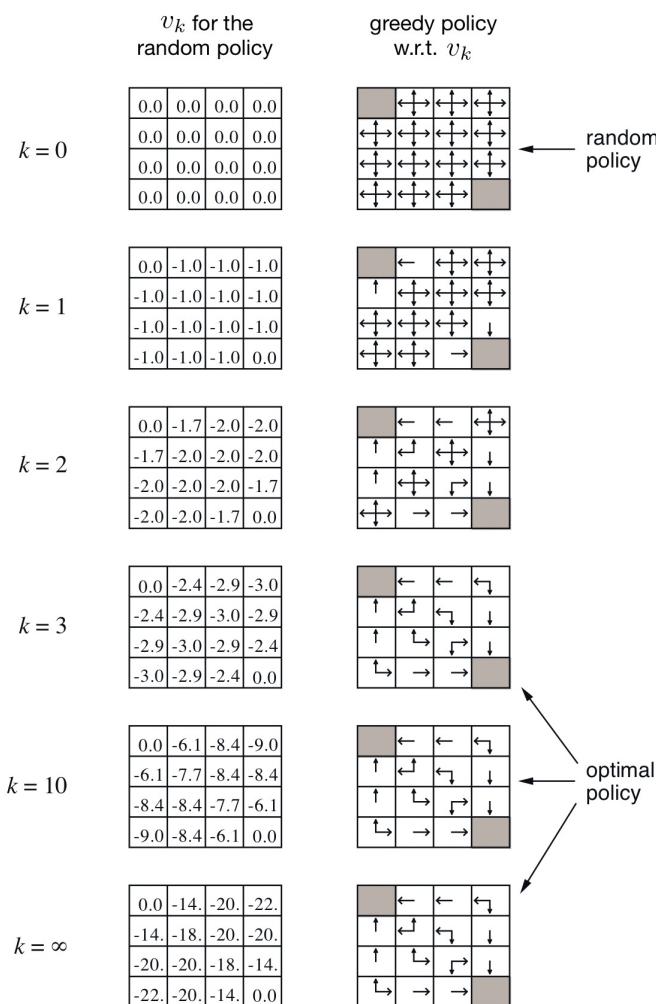
Example:



	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$R_t = -1$
on all transitions

undiscounted,
Episodic setting.



Policy improvement:-

One reason for computing the value function for a policy is to help find better policies.

Given v_{π} for some arbitrary policy π ,

For some state s , we would like to know

whether or not we should change the policy
to deterministically choose an action $a \in \Pi(s)$.

Consider selecting 'a' in 's' and thereafter

following the existing policy π . The value

of this way of behaving is

$$q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma v_{\pi}(s_{t+1}) \mid s_t = s, A_t = a]$$

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

If this is greater than $v_{\pi}(s)$, then the

new policy π' which takes action 'a' in state

's' and follow π for rest of the states would

be a better policy.

Policy improvement theorem:- (PIT)

Let π and π' be a pair of deterministic policies such that, for all $s \in S$

$$q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s)$$

Then the policy π' must be as good as, or better than π .

$$\text{i.e. } v_{\pi'}(s) \geq v_{\pi}(s) \quad \forall s \in S.$$

Proof:

$$v_{\pi}(s) \leq q_{\pi}(s, \pi'(s))$$

$$= \mathbb{E} [R_{t+1} + \gamma v_{\pi}(s_{t+1}) \mid S_t = s, A_t = \pi'(s)]$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma v_{\pi}(s_{t+1}) \mid S_t = s]$$

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma q_{\pi}(s_{t+1}, \pi'(s_{t+1})) \mid S_t = s]$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma \mathbb{E}_{\pi'} [R_{t+2} + \gamma v_{\pi}(s_{t+2}) \mid S_{t+1},$$

$$A_{t+1} = \pi'(s_{t+1})] \mid S_t = s]$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 v_{\pi}(s_{t+1}) \mid S_t = s]$$

$$\begin{aligned}
&\leq \mathbb{E}_{\pi}, \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v_{\pi}(s_{t+3}) \mid s_t = s \right] \\
&\vdots \\
&\leq \mathbb{E}_{\pi}, \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \mid s_t = s \right] \\
&= V_{\pi'}(s)
\end{aligned}$$

Given the value fn. for some Policy π ,
 Consider the following greedy policy π' wrt
 the value fn.

$$\begin{aligned}
\pi'(s) &= \underset{a}{\operatorname{argmax}} q_{\pi}(s, a) \\
&= \underset{a}{\operatorname{argmax}} \mathbb{E} \left[R_{t+1} + \gamma v_{\pi}(s_{t+1}) \mid s_t = s, A_t = a \right] \\
&= \underset{a}{\operatorname{argmax}} \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]
\end{aligned}$$

By PIT, this greedy policy should be as
 good as or better than the original policy.

The process of making a new policy that improves on an original policy, by making it greedy w.r.t. the value fn of the original policy is called Policy improvement.

Suppose the new greedy Policy π' is as good as, but not better than π , then

$$V_{\pi} = V_{\pi'}$$

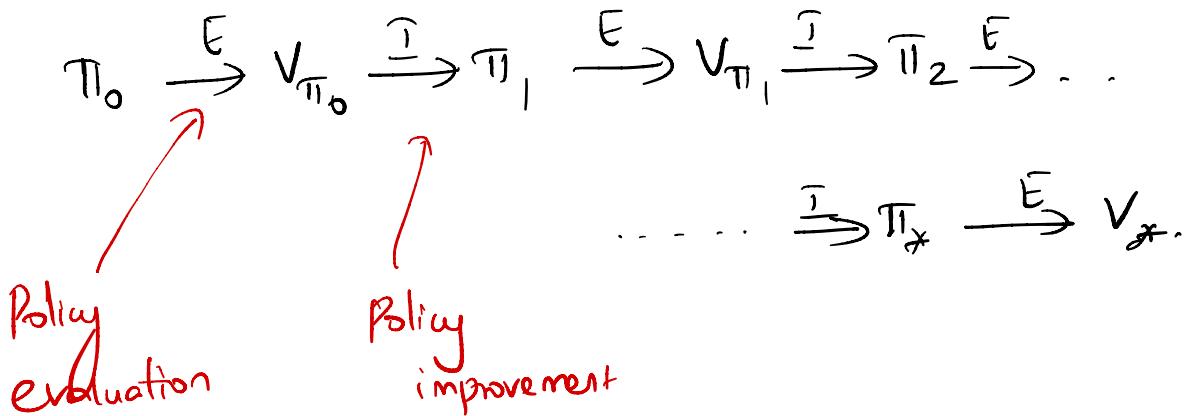
$$V_{\pi'}(s) = \max_a \mathbb{E}[R_{t+1} + r V_{\pi'}(s_{t+1}) \mid s_t = s, A_t = a]$$

This is same as Bellman Optimality eqn.

Thus $V_{\pi'} = V_x$. and both π, π' should be optimal policies.

Policy Iteration:-

Sequence of monotonically improving policies and value fns.



Each policy is guaranteed to be a strict improvement over the previous one (unless it is already optimal).

Because a finite MDP has only a finite number of policies, this process must converge to an optimal policy and optimal value function in a finite number of iterations.

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

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

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

3. Policy Improvement

$policy-stable \leftarrow true$

For each $s \in \mathcal{S}$:

$$old-action \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

If $old-action \neq \pi(s)$, then $policy-stable \leftarrow false$

If $policy-stable$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Value Iteration :- (VI)

Limitation of Policy iteration:- Each of its iteration involves Policy evaluation, which itself is an iterative computation. However, Policy evaluation only converges in limit. When to truncate policy evaluation step?

Value iteration stops policy evaluation after just one sweep.

The update rule for Value iteration combines the policy improvement and truncated policy evaluation steps.

$$\begin{aligned} V_{k+1}(s) &= \max_a E[R_{t+1} + \gamma v_k(s_{t+1}) \mid s_t = s, A_t = a] \\ &= \max_a \sum_{s',r} P(s',r|s,a) [r + \gamma v_k(s')] \end{aligned}$$

$\forall s \in S$

For arbitrary v_0 , the sequence of v_k 's can be shown to converge to v_* under some conditions

that guarantee the existence of v_x .

Note 1: v^I is obtained by simply turning the Bellman optimality eqn to an update rule.

Note 2: v^I is identical to policy evaluation except that it requires \max be taken over all actions.

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

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that

$$\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

Many variations of interleaving Policy evaluation
with policy improvement is possible!

Asynchronous DP:-

DP requires sweep over entire state space. This is not possible in many problems.

Ex: Backgammon has over 10^{20} states!

Asynchronous DP \rightarrow in-place DP algorithms with no systematic sweeping.

Just update whatever state is next available.

This makes it easier to intermix computation with real-time interaction. To solve a given MDP, we can run an iterative DP algorithm at the same time that an agent is actually experiencing the MDP.

Agent's experience \rightarrow used to decide what state to update.

Latest value & Policy info from DP $\left\{ \begin{array}{l} \\ \end{array} \right\} \rightarrow$ can guide the agent's decision making.

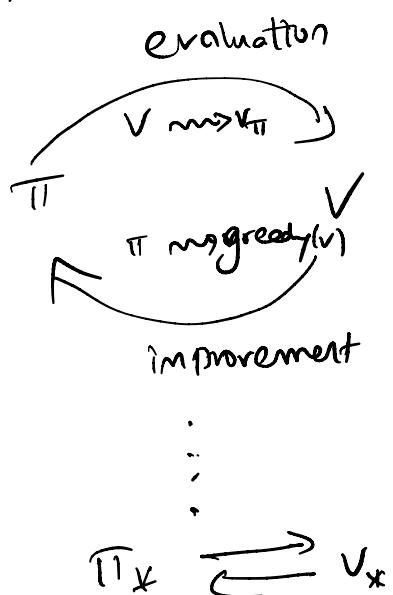
Generalized Policy Iteration:- (GPI)

Policy iteration consists of two simultaneous, interacting processes

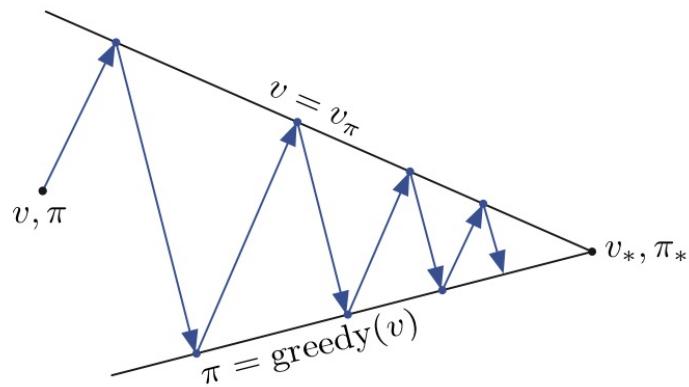
→ one making the value fn consistent with the current policy (Policy evaluation)

→ the other making the policy greedy w.r.t. the current value fn (Policy improvement)

GPI - idea of letting Policy evaluation and Policy improvement interact, independent of the granularity and other details of two processes.



These two processes are both competing and cooperating.



Efficiency of DP:-

DP finds an optimal policy in time that is polynomial in the number of states and actions.

n = # of states

k = # of actions.

Total # of deterministic policies = k^n .

DP is exponentially faster than direct search!

Linear programming can be used to solve MDP but they are not as scalable as DP.