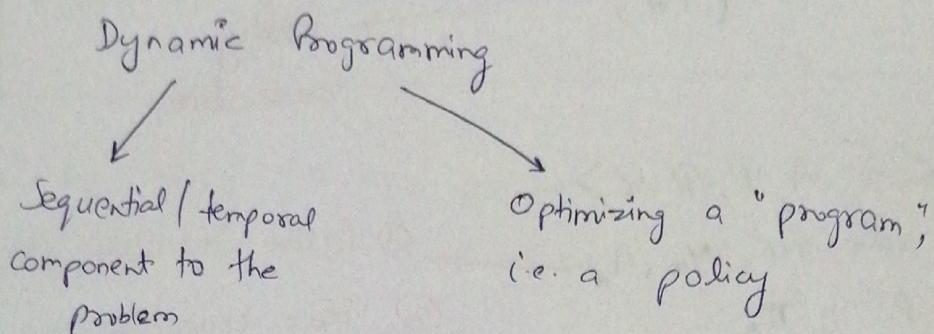


## Lecture - 3

### Dynamic Programming

In RL, DP refers to a collection of algorithms used to compute optimal policies given a perfect model of the environment as an MDP.



→ DP solves a problem by breaking them into subproblems

A DP problem has two components

① Optimal substructure

- Principle of optimality applies
- Optimal sol<sup>n</sup> can be decomposed into subproblems

② Overlapping subproblems

- Subproblems recur many times
- Their sol<sup>n</sup> can be cached & reused

→ MDP satisfies both these properties

① Bellman eqn gives recursive decomposition

② Value fun<sup>n</sup> stores & reuses sol<sup>n</sup>'s

## Planning by DP

- DP assumes full knowledge of the MDP
- It is used for planning in an MDP

### for prediction

Input: MDP  $\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$  and policy  $\pi$  (and action  $a$ )

Output: Value fun  $v_\pi(q_\pi)$   $\rightarrow$  Policy Evaluation

### for control

Input: MDP  $\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$

Output: optimal value fun  $v_*, q_*$   $\rightarrow$  ① Policy iteration  
optimal policy  $\pi_*$   $\rightarrow$  ② Value iteration

## Policy Evaluation

Given an arbitrary policy  $\pi$ , compute the (state) value fun of  $\pi$   $v_\pi$ .

from Bellman expectation eqn,

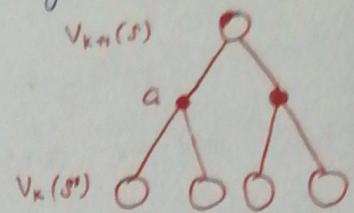
$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}} P_{ss'}^a (\delta(s, a, s') + \gamma v_\pi(s'))$$

where  $\delta(s, a, s') = \mathbb{E}[R_{t+1} | S_t = s, A_t = a, S_{t+1} = s']$

→ If the environment's dynamics are known, above eqn is a system of  $|\mathcal{S}|$  simultaneous eqns in  $|\mathcal{S}|$  unknowns  $(v_\pi(s), s \in \mathcal{S})$

→ Iterative application of above Bellman eqn gives an iterative method to compute  $V_\pi(s)$ ,

$$V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V_\pi(s)$$



$$V_{k+1}^{(s)} = \mathbb{E}_\pi [R_{t+1} + \gamma V_k(s_{t+1}) \mid s_t = s]$$

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s'} P_{ss'}^a (r(s, a, s') + \gamma V_k(s'))$$

→ The above sequence can be shown, in general, to converge to  $V_\pi$  as  $k \rightarrow \infty$ . This algo. is called iterative policy evaluation.

Iterative policy evaluation, for estimating  $V \approx V_\pi$

- ① Input  $\pi$  (Policy) and  $\theta$  (algorithm threshold)
- ② Initialize  $V(s)$ ,  $\forall s \in \mathcal{S}$  arbitrarily except that  $V(\text{terminal}) = 0$
- ③ Loop
  - ④  $\Delta \leftarrow 0$
  - ⑤ Loop for each  $s \in \mathcal{S}$ 
    - (i)  $v \leftarrow V(s)$
    - (ii)  $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} P_{ss'}^a (r(s, a, s') + \gamma V(s'))$
    - (iii)  $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
  - until  $\Delta < \theta$

## Policy Improvement

Once a policy has been evaluated, its value function can be used to construct another better policy.

A new policy  $\pi'$  can be constructed by acting greedily w.r.t. value function of previous policy.

$$\pi' = \text{greedy } (\pi)$$

This results in a deterministic policy  $\pi'$ , which improves over policy  $\pi$ .

i.e.  $\pi'(s) = \underset{a \in A}{\operatorname{argmax}} q_{\pi}(s, a) = \sum_{s'} P_{ss'}^a \left( r(s, a, s') + \gamma v_{\pi}(s') \right)$

\* This improves the value from any state  $s$  over one step

$$q_{\pi}(s, \pi'(s)) = \max_{a \in A} q_{\pi}(s, a) \geq q_{\pi}(s, \pi(s)) = v_{\pi}(s)$$

\* We can also show that new policy  $\pi'$  is better than policy  $\pi$  overall. i.e.  $v_{\pi}(s) \geq v_{\pi'}(s)$ .

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

$$V_{\pi}(s) \leq V_{\pi'}(s)$$

This is true for all states  $s \in \mathcal{S}$ . So, the new policy  $\pi'$  improves over original policy  $\pi$ .

When the improvement stops, i.e. when for two successive iterations  $V(\pi') = V(\pi)$ , then

$$\begin{aligned} V_{\pi'}(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma V_{\pi'}(s_{t+1}) \mid s_t = s, A_t = a] \\ &= \max_a \sum_{s'} P_{s,a}^s (r(s,a,s') + \gamma V_{\pi'}(s')) \end{aligned}$$

The Bellman optimality eqn has been satisfied, and we can conclude that

$$\boxed{\pi' = \pi_*}$$

i.e.  $\pi'$  must be an optimal policy, as well as  $\pi$ .

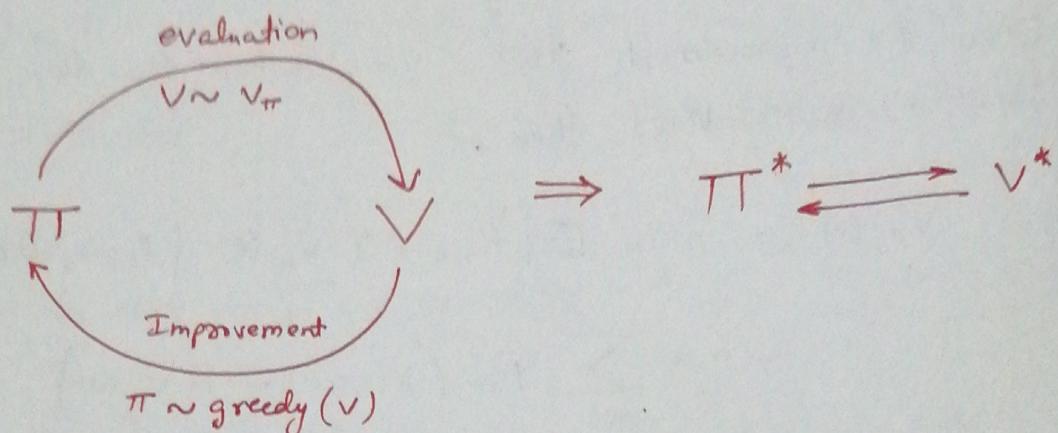
### Policy Iteration

Once a policy  $\pi$  has been improved using  $V_{\pi}$  to yield a better policy  $\pi'$ , we can compute  $V_{\pi'}$  again & improve it again to yield an even better policy  $\pi''$ .

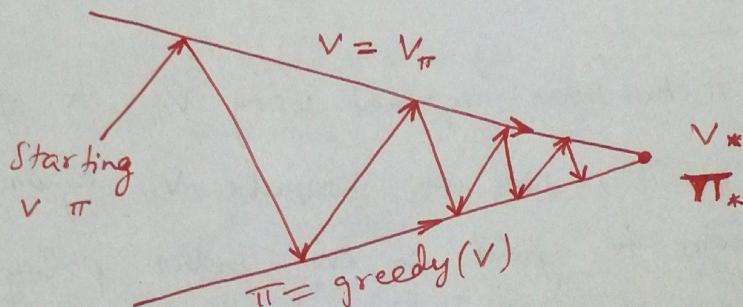
This way of finding optimal policy is called Policy Iteration.

→ Policy iteration consists of two simultaneous, interacting processes, one making the value fun<sup>n</sup> consistent with the current policy & other making policy greedy w.r.t. current value fun<sup>n</sup>. This is in general not necessary to alternate b/w PEval & PImp. In asynchronous DP, in some cases a single

State is updated in one process before returning to other states. As long as both processes continue to update all states (infinitely as  $t \rightarrow \infty$ ), convergence to optimal policy is guaranteed.



→ Generalized policy Iteration (GPI) is the general idea of letting policy-evaluation & policy-improvement processes interact, independent of granularity and other details of the two processes.



Policy iteration (using iterative policy eval.) for estimating  $\pi \approx \pi^*$  using  $V^*$

- ① Initialize  $V(s) \in \mathbb{R}$  &  $\pi(s) \in A(s)$  arbitrarily
- ② Perform policy evaluation for policy  $\pi$  to get  $V$
- ③ Policy improvement  
for each state  $s \in S$   

$$\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} \sum_{s'} P_{ss'}^a (r(s,a,s') + \gamma V(s'))$$
- ④ Go to step ② until convergence

## Exercise 4.5 (Sutton)

Policy iteration to estimate  $\pi_*$  using action-value function  $q_*$ .

### ① Initialization

$Q(s, a) \in \mathbb{R}$  and  $\pi(s) \in A(s)$  arbitrarily  $\forall s \in S$

### ② Policy evaluation

Loop

$$\Delta \leftarrow 0$$

Loop for each  $s \in S$  and  $a \in A(s)$

$$q \leftarrow Q(s, a)$$

$$Q(s, a) \leftarrow \sum_{s'} P_{ss'}^a (r(s, a, s') + \gamma \sum_{a'} \pi(a'|s') Q(s', a'))$$

$$\Delta \leftarrow \max(\Delta, |q - Q(s, a)|)$$

until  $\Delta < \theta$  (threshold determining accuracy)

### ③ Policy improvement

policy-stable  $\leftarrow$  True

for each  $s \in S$  and  $a \in A(s)$

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

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s'} P_{ss'}^a (r(s, a, s') + \gamma \sum_{a'} \pi(a'|s') Q(s', a'))$$

if old-action  $\notin \{a_i\}$  (set of equiprobable policies from  $\pi(s)$ )  
(actions)

then policy-stable  $\leftarrow$  false

if Policy-stable then stop and return  $Q \approx q_*$  and  $\pi \approx \pi_*$

else goto step ②

## Value iteration

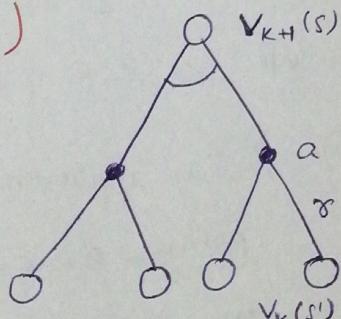
Policy iteration has a drawback that each of its iterations involves policy evaluation, which may itself be an iterative process requiring multiple sweeps through the state set  $\mathcal{S}$ .

However, without losing convergence guaranteed policy evaluation can be truncated in several ways. One special case is when policy evaluation is stopped after just one sweep (one update of each state). This is called value iteration algorithm.

$$V_{k+1}(s) = \max_a \mathbb{E} [R_{t+1} + \gamma V_k(s_{t+1}) \mid s_t = s, A_t = a]$$

$$= \max_a \sum_{s'} P_{ss'}^a (\gamma r(s, a, s') + \gamma V_k(s'))$$

$$\forall s \in \mathcal{S}$$



- \* This can again be shown to converge to  $V_*$  under conditions that guarantee the existence ~~convergence~~ of  $V_*$

## Principle of Optimality

A policy  $\pi(a|s)$  achieves the optimal value from state  $s$ ,  $V_\pi(s) = V_*(s)$  iff

- \* for any state  $s'$  reachable from  $s$ ,  $\pi$  achieves the optimal value from  $s'$ .  $V_\pi(s') = V_*(s')$

Value iteration can be thought of as applying principle of optimality iteratively.

$$V_*(s) = \max_{a \in A} R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_*(s')$$

$$= \max_{a \in A} \sum_{s'} P_{ss'}^a (\delta(s, a, s') + \gamma V_*(s'))$$

Value iteration for estimating  $\pi \approx \pi_*$

- ① Initialize  $V(s)$  &  $s \in S^+$  arbitrarily except  $V(\text{terminal}) = 0$
- ② Loop

$$\Delta \leftarrow 0$$

Loop for each  $s \in S$

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s'} P_{ss'}^a (\delta(s, a, s') + \gamma V(s'))$$

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

until  $\Delta < \theta$

- ③ Output deterministic policy  $\pi \approx \pi_*$  s.t

$$\pi(s) = \arg \max_a \sum_{s'} P_{ss'}^a (\delta(s, a, s') + \gamma V(s'))$$

Synchronous DP algorithms

Problem	Bellman Egn	Algorithm
Prediction	Bellman expectation egn	Iterative policy evaluation
Control	Bellman exp. egn + Greedy policy improvement	Policy iteration
Control	Bellman optimality egn	Value iteration

## Asynchronous DP

One major drawback of DP methods discussed so far is that they used synchronous backups; they require sweeps of the state set. This is very expensive if state set is very large.

One soln to this is asynchronous DP, which backs up states individually in any order. It, for each selected state, applies the appropriate backup.

- \* This is also guaranteed to converge if all states continue to be selected.

### ① In-place DP

Only store one copy of value func

$$V(s) \leftarrow \max_{a \in A} \sum_{s'} P_{ss'}^a (r(s,a,s') + \gamma V(s'))$$

### ② Real-time DP

Only update/backup states that are ~~backup~~ relevant to the agent. Use agent's exp. to guide selection of states.

### ③ Prioritised sweeping

Use magnitude of Bellman error to guide state selection.

Eg.

$$\left| \max_{a \in A} \left( \sum_{s'} P_{ss'}^a (r(s,a,s') + \gamma V(s')) - V(s) \right) \right|$$