

## Ch-4, Dynamic Programming

- Dynamic Programming (DP) - collection of algo that can be used to compute optimal policies given a perfect model of the environment as a MDP
- Assume - env is a finite MDP. - assume that its state, action and reward sets  $S, A$  and  $R$  are finite and that its dynamics are given by set of probs  $p(s', r | s, a)$  for all  $s \in S, a \in A(s), r \in R$  and  $s' \in S^+$  ( $S^+$  is  $S$  plus a terminal state if prob is episodic)
- Although DP ideas can be applied to prob with continuous state and action spaces, exact sol<sup>n</sup> are only possible in special cases.  
A common way - obtain approx sol<sup>n</sup> for tasks with continuous states and actions is to quantize the state and actions and then apply finite-state DP methods
- Key idea of DP - Use of value functions to organize and structure the search for good policies

We can easily obtain optimal policy once we have found the optimal value functions  $V_*$  or  $q_*$  which satisfy Bellman Optimality eq<sup>n</sup>:

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

$$= \max_a \sum_{s', r} p(s', r | s, a) \left[ r + \gamma V_*(s') \right]$$

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

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

Algo are derived from these Bellman eq<sup>s</sup> by turning them into update rules for iteratively improving approximation of the value function

## # Policy Evaluation (Prediction)

→ Policy evaluation is the process of computing the state value function  $V^\pi$  for a given policy  $\pi$ . This is also known as prediction problem.

$$\rightarrow V^\pi(s) \doteq E_\pi [G_t \mid S_t = s]$$

$$= E_\pi [R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

$$= E_\pi [R_{t+1} + \gamma V^\pi(S_{t+1}) \mid S_t = s]$$

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

where  $\pi(a|s)$  - prob to take action  $a$  in state  $s$  under policy  $\pi$   
 $p(s', r|s, a)$  - prob of transitioning to state  $s'$  and receiving reward  $r$  after taking action  $a$  in state  $s$ .

$\gamma$  - discount factor

• The existence and uniqueness of  $V^\pi$  are guaranteed as long as either  $\gamma < 1$  or eventual termination is guaranteed from all states under policy  $\pi$ .

→ To compute  $V^\pi$  we can use an iterative approach known as iterative policy evaluation. Starting with an initial approximation  $V_0$  (which can be arbitrary except for terminal state, if any, which must be given a value of 0), we



repeatedly apply the following update rule to each state  $s$ :

$$V_{k+1}(s) \doteq E_{\pi} [R_{t+1} + \gamma V_k(s_{t+1}) \mid s_t = s] \\ = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$$

This process is repeated until the value function converges to a fixed pt, which is guaranteed under certain conditions (e.g. if  $\gamma < 1$  or if eventual termination is guaranteed for all states under policy  $\pi$ )

→ To produce each successive approx,  $V_{k+1}$  from  $V_k$ , iterative policy evaluation applies the same operator to each state  $s$ : it replaces old value of successor states of  $s$  with a new value obtained from the old values of the successor states of  $s$ , and expected immediate rewards along all the one-step transitions possible under the policy being evaluated → Expected update.

In this eq<sup>n</sup>, the inner summation computes the expected return for each action  $a$ , considering all next possible states  $s'$  and rewards  $r$ . The outer summation then averages the expected return according to policy  $\pi$  which ~~represents~~ specifies the probs of taking each action  $a$  in state  $s$ .

This process is applied to each state  $s$  in state space, and the result is a new approximate  $V_{k+1}$  of the value function. By repeating this process repeatedly, value function converges to true value function  $V^{\pi}$  for given policy  $\pi$ .

The term "expected update" emphasizes that this process is

Based on expectation over all possible next state with the  
one single sampled next state. - key characteristic of  
DP algo

→ Iterative Policy Evaluation: for estimating  $V \approx V_\pi$

Input  $\pi$ , the policy to be evaluated

Algo parameter: a small threshold  $\theta > 0$  determining accuracy of  
estimate

Initialize  $V(s)$  for all  $s \in S'$ , arbitrarily except  $V(\text{terminal}) = 0$

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in S$ :

$v \leftarrow V(s)$

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

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

until  $\Delta < \theta$

## # Policy Improvement:-

→ Reason for computing the value function for a policy -  
help find better policies.

→ value function of a policy ( $V_\pi$ ): This is the expected  
return when starting in state  $s$  and following policy  $\pi$   
thereafter

→  $q_\pi(s, a) \doteq E [R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t = s, A_t = a]$   
 $= \sum_{s', r} p(s', r | s, a) [r + \gamma V_\pi(s')]$



# Policy Improvement - Computation of a improved policy given the value function for that policy

(3)

PAGE EDGE

Date: / /

This is the expected return when ~~starting~~ <sup>starting</sup> in state  $s$ , and following policy  $\pi$  thereafter.

→ Policy Improvement: Suppose we have a deterministic policy  $\pi$  and we want to know if we should ~~change~~ change the policy to always choose an action in state  $s$ . We compare  $q_{\pi}(s, a)$  with  $v_{\pi}(s)$ :

• If  $q_{\pi}(s, a) > v_{\pi}(s)$ , it means that taking action  $a$  in state  $s$  and then following  $\pi$  from  $s$ . Hence we should change the policy to always choose action  $a$  in state  $s$ .

→ Policy Improvement Theorem:

If for all states  $s$ ,  $q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s)$  where  $\pi'$  - new policy then  $\pi'$  is as good or better than  $\pi$ . That is,  $v_{\pi'}(s) \geq v_{\pi}(s)$  for all  $s$ .

If there is a strict inequality for any state, then  $\pi'$  is strictly better than  $\pi$ .

If  $q_{\pi}(s, a) > v_{\pi}(s)$  then the changed policy is indeed better than  $\pi$ .

→ Greedy Policy Improvement:

To improve a policy, we can construct a new policy  $\pi'$  that is greedy w.r.t value function of the original policy  $\pi$ . For each state  $s$ ,  $\pi'$  chooses the action that maximizes  $q_{\pi}(s, a)$ . This new policy  $\pi'$  is better or as good as the original policy  $\pi$ .

$$\pi'(s) \doteq \underset{a}{\operatorname{argmax}} q_{\pi}(s, a)$$

$$\begin{aligned}\pi'(s) &= \underset{a}{\operatorname{argmax}} E \left[ R_{t+1} + \gamma V_{\pi'}(S_{t+1}) \mid S_t = s, A_t = a \right] \\ &= \underset{a}{\operatorname{argmax}} \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma V_{\pi'}(s') \right]\end{aligned}$$

where  $\operatorname{argmax}$  - value of  $a$  at which exp that follows is maximized

Greedy policy takes the action that looks best in short term. By constructing, the greedy policy meets the conditions of the policy improvement theorem, so we know that it is 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 ~~with~~ greedy w.r.t value function of the original policy is called policy improvement.

Note Policy improvement thus must give us a strictly better policy except when the original policy is already optimal.

Optimal

→ Convergence to Policy: If the greedy policy  $\pi'$  is as good as, but not better than, the old policy  $\pi$ , then both  $\pi$  and  $\pi'$  must be optimal policies. If not, the process of policy improvement can be repeated <sup>with  $\pi'$</sup>  to get even better policy. At this process will eventually converge to policy which is optimal.

Summary: Policy improvement is a process of using the value function of a policy to construct a better policy. By repeatedly applying policy improvement, we can iteratively improve policies until we reach an optimal policy.



## # Policy Iteration:-

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

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} V_*$$

Where  $\xrightarrow{E}$  denotes policy evaluation  
 $\xrightarrow{I}$  denotes policy improvement

- 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 no. of policies, this process must converge to an optimal policy and optimal value function in a finite number of iterations.

This way of finding optimal policy - Policy Iteration

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

1) Initialization:

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

2) Policy Evaluation:

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in S$ :

$$V \leftarrow V(i)$$

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

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

until  $\Delta < \theta$  (a small +ve no. determining the accuracy of estimation)

### 3) Policy Improvement:

policy-stable  $\leftarrow$  true

For each  $s \in S$ :

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

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) (r + \gamma V(s'))$$

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

If policy-stable, then stop and return  $V \approx V + \alpha(\pi \approx \pi)$ ,  
else go to 2.

Note

Starting each policy evaluation with the value function for the previous policy typically results in faster convergence, as the value function often changes little from one policy to the next.

- Policy Improvement theorem assures us that these policies are better than the original random policy.

### Value Iteration:-

→ Drawback of policy iteration: Each of its iteration involves policy evaluation, which may itself be a protracted iterative computation requiring multiple sweeps through



the state set.

→ Value Iteration - variant of policy iteration

→ Value iteration simplifies this by truncating the policy evaluation step to just one sweep through the state set, effectively combining policy evaluation and policy improvement into a single step.

Update rule for value iteration is:

$$V_{k+1}(s) = \max_a E \left[ R_{t+1} + \gamma V_k(S_{t+1}) \mid S_t = s, A_t = a \right]$$

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

This rule is similar to the policy evaluation update, but it includes a maximization over all actions  $a$ , which incorporates policy improvement directly into update.

→ In practice, value iteration is terminated when the change in the value function is small enough, indicating that it has approximately converged to  $V^*$ . The algo then derives an optimal policy from the final value function by selecting, for each state, the action that maximizes the expected return according to the value function.

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

Algorithm parameters: a small threshold  $\delta > 0$  determine accuracy of estimation

Initializing  $V(1)$  for all  $s \in S^+$ , arbitrarily except

that  $V(\text{terminal}) = 0$

loop:

$\Delta \leftarrow 0$

loop for each  $s \in S$ :

$V \leftarrow V(s)$

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

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

until  $\Delta < \theta$

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

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

Note Value iteration is obtained simply by turning the Bellman optimality eq<sup>n</sup> into an update rule.

→ value iteration is often faster than policy iteration bcoz it avoids the separate, potentially lengthy policy evaluation step. It effectively combines one sweep of policy evaluation and one sweep of policy improvement in each iteration.

Moreover, variations of value iteration can be created by adjusting the number of policy evaluation sweeps b/w each policy improvement sweep, allowing ~~best~~ balance b/w computational efficiency and convergence speed.

Imagine you are playing a game, finding a best way from start → find in a maze with lot of paths

1) Policy Evaluation (Checking Your Path): You pick a path and walk - get some pts. Track of pts you got



for the whole path.

- 2) Policy Improvement (Finding a better path): If you find a step where changing a direction ~~make~~ gets you more pts, you change your path.
- 3) Policy Iteration (Keep Improving): Keep doing the above 2 steps over and over. You keep doing this till you don't find any more ways to get more pts.
- 4) Value Iteration (Shortcut): Instead of checking the whole path, you pick the best move at every step. You make the best move and we end up with that path faster.

## # Asynchronous Dynamic Programming:-

- Drawback of traditional DP method - involve ~~many~~ operations over the entire set of MDP which can be computationally infeasible for large-scale problem.
- Asynchronous DP addresses this issue by <sup>updating</sup> ~~updating~~ the values of states in any order, without the need for systematic sweeps.
- Key Characteristics of Asynchronous DP algorithm:
  - 1) In-place updates: values of state are updated in place using most recent value of other states that are available.
  - 2) Flexibility in state selection: Some states might be updated multiple times before others are updated once.

3) Convergence requirements: To ensure convergence all states must continue to be updated; no state can be permanently ignored. The sequence of state updates can be deterministic or stochastic.

4) Asynchronous value Iteration: An example of an asynchronous DP algorithm is a version of value iteration that updates the value of only one state at each step. Asymptotic convergence to the optimal value function  $V^*$  is guaranteed under certain conditions.

5) Efficiency: Improve the rate of progress by selecting states for updates in a way that allows value info to propagate efficiently. This can include updating states more frequently if they are less more relevant to optimal policies or skipping updates for states that are less relevant.

6) Real time iteration: Asynchronous algo can be interleaved with real time iteration, allowing a agent to use the latest value and policy info for decision-making while experiencing the environment. Updates can be applied to states as the agent visits them, focusing computational resources on the most relevant parts of state space.

→ Asynchronous DP - offer more flexible and potentially more efficient approach to solve large-scale MDPs.

## # Generalized Policy Iteration:-

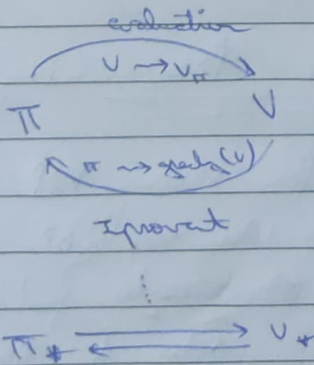
→ Policy Iteration consists of 2 simultaneous, iterative processes:

- Policy evaluation - making value function consistent with current policy.
- Policy improvement - making policy greedy w.r.t current value function.



- Policy iteration - these 2 processes alternate.
- Value iteration - for eg: only a single iteration of policy evaluation is performed in b/w each policy improvement
- Asynchronous DP - the evaluation and improvement are interleaved at each or even finer grain
- Generalized Policy Iteration (GPI) - encompasses a wide range of RL methods.
- It describes the interaction b/w 2 processes: policy evaluation and policy improvement. These processes work together to find the optimal value function and policy even though they might operate at different granularities and in various ways across diff algo.
- Almost all RL methods are GPI.

All have identifiable policies and value functions, with ~~the~~ policy always being improved w.r.t to value function and value function being driven towards value function for policy.



If both the evaluation process and improvement process stabilize then <sup>the</sup> value function and policy optimal.

GP - highlights the fundamental iteration for evaluating policy, which is at heart of learning optimal RL

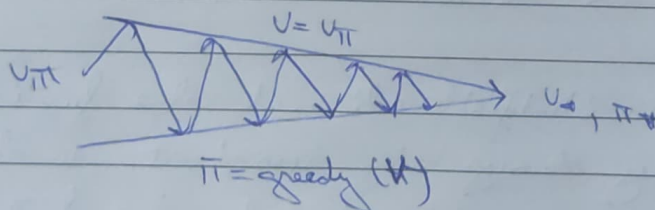
The value function stabilizes only when it is consistent with the current policy and the policy stabilizes only when it is greedy w.r.t. current value function.

Thus both processes stabilize only when a policy has been found that is greedy with respect to its own ~~evaluation~~ value function.

→ Evaluation and Improvement process in WPI can be viewed as both competing and cooperating.

They compete - pull in opposite directions. Making the policy greedy w.r.t. value function typically makes value function inconsistent for changed policy, and making value function consistent ~~inconsistent~~ with policy typically causes that policy no longer to be greedy.

I don't know, however the process iterates to find optimal value function and a optimal policy.



## Efficiency of Dynamic Programming:

→ DP is powerful method for solving MDP.

time taken ~~by~~ by DP

→ Efficiency: In worst case, ~~it takes~~ to find an optimal policy is polynomial in the number of <sup>states</sup> ~~states~~ and actions. This is exponentially faster than a direct search in policy space, which would require examining each of  $k^n$  possible deterministic policies ( $n$  - no of states,  $k$  - no of actions)



- Handling large State Spaces: Although the curse of dimensionality (exponential growth of state space with the number of state variables) poses challenge, DP is comparatively better suited to handle large states than methods like direct search or linear programming.
- Faster Convergence with Good Initialization: Policy iteration value functions converge much faster than their theoretical worst-case run times. Starting with good initial value functions or policies can further speed up convergence.
- Asynchronous Methods: For large state spaces, asynchronous DP preferred. These methods update the value of state in any order and do not require computation of memory for every state in each sweep.
- Bootstrapping:
- special prop of DP methods
  - all of them update estimate of the value of states based on the estimate of the values of successor states.
  - That is they update estimates based on ~~other~~ <sup>using the</sup> estimate.
  - May RL perform bootstrapping even though it does not require as DP requires it (DP requires complete and accurate model of env).