

# 2 DYNAMIC PROGRAMMING (DP)

1

## 2.1 Introduction

### - DP and its Requirements

- DP is a method for solving complex problems by:
  - ↳ breaking them into sub problems.
  - ↳ solve sub prob  $\rightarrow$  combine solution to sub prob.
- Requirements for DP are:
  - $\rightarrow$  Optimal sub-structure: an opt solution can be decomposed into opt solution of sub-problems.
  - $\rightarrow$  Overlapping subproblems: subproblems occur many times; solutions can be cached and reused.
  - $\rightarrow$  MDPs satisfy both properties
    - ↳ Bellman equation gives recursive decomposition.
    - ↳ value function stores  $\gamma$  return solution.

### - Planning by DP

In planning someone tells us the structure of the MDP (the dynamics and rewards) and given perfect knowledge of how the environment works we want to solve the MDP.

Two special cases of planning can be solved in MDP:

#### • Plan to solve a prediction problem

$\rightarrow$  Input: MDP  $(S, A, P, R, \gamma)$  and policy  $\pi$

$\rightarrow$  Output: value-function  $V_{\pi}$

} Someone tells us the MDP and also gives a policy and asks you, how much reward would one get walking around in this world.

#### • Plan is used for control

$\rightarrow$  Input = MDP  $(S, A, P, R, \gamma)$

$\rightarrow$  Output: optimal value function  $V_{*}$  amongst all policies what's the most reward that can be achieved from this MDP

} Someone tells us the MDP, but now instead of giving a policy they want to know what's the best policy;

- Iterative Policy Evaluation

- Problem : evaluate a given policy  $\pi$
- Solution : take bellman eq and turn into iterative update

↳  $V_1, V_2, \dots, V_{\pi}$

We begin with initializing the value space randomly. At  $V_1$  we would have plugged in our 1 step look ahead bellman eq. , we iterate this many time and we eventually end up w/ the val function.

↳ The way this is done is using synchronous backups i.e.

↳ @ each iteration we go sweep across all states in MDP

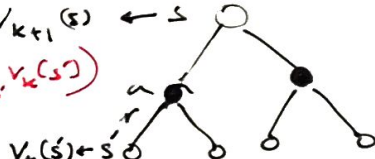
@ each iteration  $k+1$

$\forall$  states  $s \in S$

Update  $V_{k+1}(s)$  from  $V_k(s')$

where  $s'$  is a successor state of  $s$

$$\rightarrow V_{k+1}(s) = \sum_{a \in A} \pi(a|s) \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right)$$

$$V^{k+1} = R^{\pi} + \gamma P^{\pi} V^k$$


- Grid-world example

→ Undiscounted episode MDP ( $\gamma=1$ )

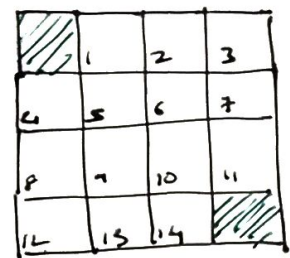
→ Non-terminal states 1... 14

→ One terminal state

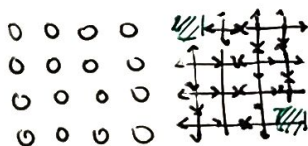
→ Actions out of grid leave state unchanged

→ Reward = -1 until terminal state reached

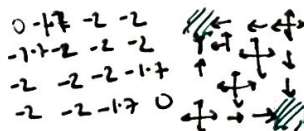
→ Agent follows uniform random policy  $\pi(a|\cdot) = \pi(d|\cdot) = \pi(s|\cdot) = \pi(w|\cdot) = 0.25$



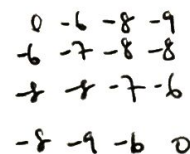
$k=0$



$k=2$



$k=10$



We notice with the gradual inc in iterations the policy gets better until optimality is reached.

## Policy Iteration

For the previous setting we evaluated a given policy and here we'd want to find the best possible policy.

### - Making Policy better

- Given policy  $\pi$

→ evaluate policy  $\pi$ :  $V_{\pi}(s) = E[R_{t+1} + \gamma V_{\pi}(s_{t+2}) + \dots | S_t = s]$

→ Improve policy by acting greedily w.r.t  $V_{\pi}$

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

To explain further, we begin at state  $s$ , cal val function of state around it <sup>on to</sup> which an action can be effectuated and the greedily move in the direction w/ highest value. Again, evaluate the state you land in and move again greedily.

→ In the small grid world problem this was enough but not in real world prob where you'll need more iterations of improvement/evaluation.

→ This <sup>↑</sup> process always converges to an optimal policy  $\pi^*$

### - Illustration

- We start w/

some i/p's  
(random vals. (0s))

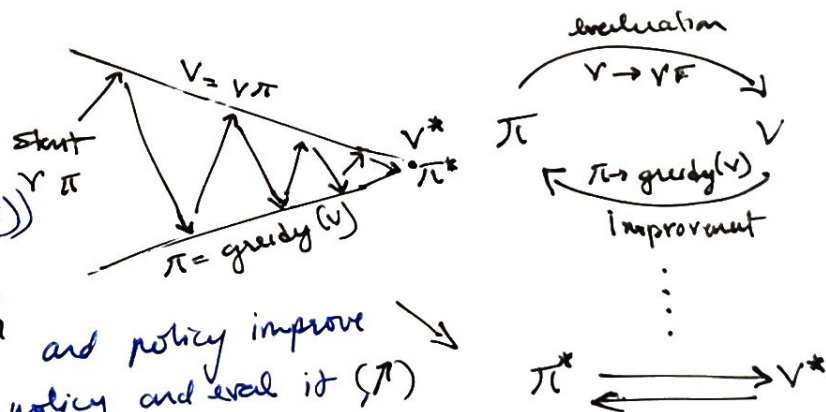
and some policy

↳ policy eval  $\nearrow$  and policy improve  $\searrow$

↳ start w/ policy and eval it ( $\nearrow$ )

then ( $\searrow$ ) act greedily w.r.t

to val function to give us the new policy ----



- Check Rental car policy example in slides.



## - Policy Improvement

- Consider a deterministic policy,  $a = \pi(s)$  : for example acting greedily always makes my policy deterministic.
- We can improve policy by acting greedily  $\pi'(s) = \underset{a \in A}{\operatorname{argmax}} q_{\pi}(s, a)$
- If we act greedily, the greedy policy at least improves the values we get over 1 step (immediate 1 step)

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

→ This is indirectly saying that if we pick this greedy policy the value of greedy policy is at least as much as the policy before we "greedified" it. So we don't make things worse ever, at least make it better. optimal

- If improvements stop, we'd want to know if we achieved

$$v_{\pi}(s, \pi'(s)) = \underset{a \in A}{\operatorname{max}} q_{\pi}(s, a) = q_{\pi}(s, \pi'(s)) = v_{\pi'}(s)$$

→ This is basically the Bellman equation, i.e. if this is satisfied we have satisfied the Bellman optimality equation

$$v_{\pi}(s) = \underset{a \in A}{\operatorname{max}} q_{\pi}(s, a)$$

→ And if Bellman opt eq satisfied then  $\pi$  must be the optimal policy  $\therefore v_{\pi}(s) = v_*(s) \quad \forall s \in S$

## - Extensions to policy iterations

### • Modified Policy Iteration

→ How the policy eval need to converge to  $v_{\pi}$ ?

→ Or should we introduce a stopping condition : when the state-values are only changed by a tiny amt after each iteration then one can just stop. (small  $\epsilon$  would be sufficient for  $\pi^*$ )

... → Or simply stop after  $k$  iterations of iterative policy evaluation

→ Why not update policy every iteration? ie stop after  $k=1$ , calc values and proceed w.r.t to that value function and then proceed: → this is equivalent to value iteration

2.4

## Value Iteration

### - Principle of Optimality (Recap)

- Any opt policy can be subdivided into 2 compse
  - An optimal 1<sup>st</sup> action  $A_*$
  - followed by following an opt policy from successor state
- A policy  $\pi(a|s)$  achieves the opt value from state  $s$ ,  
 $V_{\pi}(s) = V_*(s)$  iff → for any state  $s'$  reachable from  $s$   
 $\rightarrow \pi$  achieves the opt val from state  $s'$ ,  
 $V_{\pi}(s') = V_*(s')$

### - Deterministic Value Iteration

- If we know the solution to sub problems  $V_*(s')$  (The opt values from  $s'$ ). Use this info to build opt value func for prev step
- All we need to do is; one step look ahead using bellman optimality eqn  $V_*(s) \leftarrow \max_{a \in A} R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_*(s')$
- The idea of value iteration is to apply updates iteratively; ie instead of someone giving  $V_*(s)$ , we first initialize it randomly and then iteratively arrive at opt sol.  
 $\rightarrow$  start w/ final reward & work backwards.

### - Value Iteration

- Problem: find optimal  $\pi$
- Solution: Iterative application of Bellman Optimality Backup  
 $(V_1, V_2, \dots, V_n)$
- Using synchronous backups  
 $\rightarrow$  @ each iteration  $k+1 \rightarrow$  for all states  $s \in S \rightarrow$  update  $V_{k+1}(s)$  from  $V_k(s')$

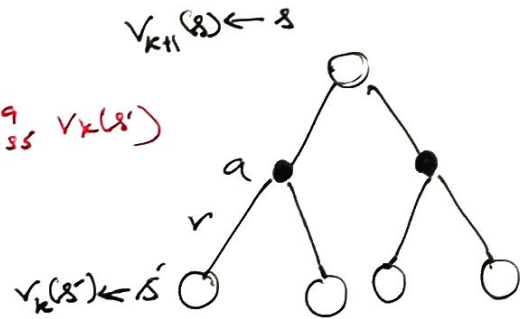


- Unlike policy iteration, there's no explicit policy.

↳ w/ policy iteration @ every step we constructed a val-function that was the val-function for a particular policy whereas w/ value iteration a given  $V_i$  may not correspond to any real policy.

$$V_{k+1}(s) = \max_{a \in A} \left( R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V_k(s') \right)$$

$$V_{k+1} = \max_{a \in A} R^a + \gamma P^a V_k$$



## - Synchronous Dynamic Programming Algos

Problem	Bellman Equation	Algo
Prediction	Bellman Expectation Eq	Iterative Policy Evaluation
Control	BEB + greedy policy improvement	Policy Iteration
Control	Bellman optimality Eq	Value Iteration

- All problems are planning problems, in all cases the MDP is given and we're just trying to solve for it.

↳ 2 types of planning problem: Prediction and Control  
 (Prediction):

- In case 1, we're trying to solve how much reward you can get from a particular policy; To figure that out we use the Bellman expectation equation, which is taken and turned into an iterative update to arrive @ an iterative policy eval algo.
- Talking about control, which shows how to get the most out of an MDP, how to max total reward, how

find  $v^*$  and hence the optimal policy.

7

- In case 2, the 1<sup>st</sup> control approach was to again use the Bellman exp eq and iterate over it to evaluate the policy, but then to alternate that process of evaluation w/ the process of policy improvement which gave us the policy iteration algo.
- In case 3, we turn the B.O. eq to an iterative update to get the value iteration algo.
- All these algos are based on state-val func  $V_\pi(s)$  or  $V_*(s)$ .
- Complexity if we have  $n$  actions &  $m$  states:  $O(mn^2)$  per iteration. Not too great for bigger problems.
- Could also apply to action-val func  $q_\pi(s, a)$  or  $q_*(s, a)$   
 $\hookrightarrow$  complexity  $O(m^2 n^2)$

Backups just mean getting  $V(s)$  by one-step look ahead.



## - Synchronous DP

- DP methods so far use synchronous backups i.e. all states are backed up in parallel.
- Async DP backs up states individually in any order. For each selected state apply appropriate backup. Reduces computation and guarantee to converge if all states continue to be selected.
- 3 ideas for Async DP

$\rightarrow$  In place DP: stores only 1 copy of val function

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

$\rightarrow$  Prioritized sweeping: use Bellman error to guide state selection

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

@ iteration  $k+1$

$\hookrightarrow$  Backup state w/ largest Bellman error  $\hookrightarrow$  update Bellman error of affected states after each backup. Can be implemented using priority queue

$\rightarrow$  Real time DP  $V(s_t) \leftarrow \max_{a \in A} (R_{s_t}^a + \gamma \sum_{s' \in S} P_{s_t s'}^a V(s'))$