

## Planning by Dynamic Programming

A) Dynamic : sequential or temporal component to the problem.

We try to solve it step by step, something changes in the intermediary & we change our step or action accordingly

Programming : Optimising a "program" i.e a policy

We use DP to solve large & complex problems by breaking them into subproblems & then solve them & combine solutions to subproblems

Dynamic Programming is a very general solution method for problems which have two properties

→ Optimal substructure :

principle of optimality applies where solving the subproblems actually leads to an optimal overall problem.

→ Overlapping Subproblems :

the subproblems keep appearing again & again, solving or learning to solve subproblems helps us learn something.

solution thus can be cached & reused

MDP's satisfy both these properties

- \* The Bellman equation gives the recursive decomposition
- \* The value function stores & learns solutions

③ Dynamic Programming assumes full knowledge of the MDP

It is used for Planning in an MDP

→ For prediction : (if we move around randomly how much reward do we get)

Input : MDP  $\langle S, A, P, R, \gamma \rangle$  & policy  $\pi$

or : MRP  $\langle S, P^\pi, R^\pi, \gamma \rangle$

Output : Value function  $V_\pi$

→ For control :

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

Output : Optimal value function  $V^*$

$\Sigma$  : Optimal policy  $\pi^*$

(Without Section I formulate more of  $(-\kappa)$ )

④ Policy Evaluation

Problem : Evaluate a given policy  $\pi$

Solution : iterative application of Bellman Expectation backup

We start off with an arbitrary value function

i.e.  $v_1 = 0$  & then we do a look ahead

Bellman eq & iterate till

$v_1 \rightarrow v_2 \dots \rightarrow v_\pi$

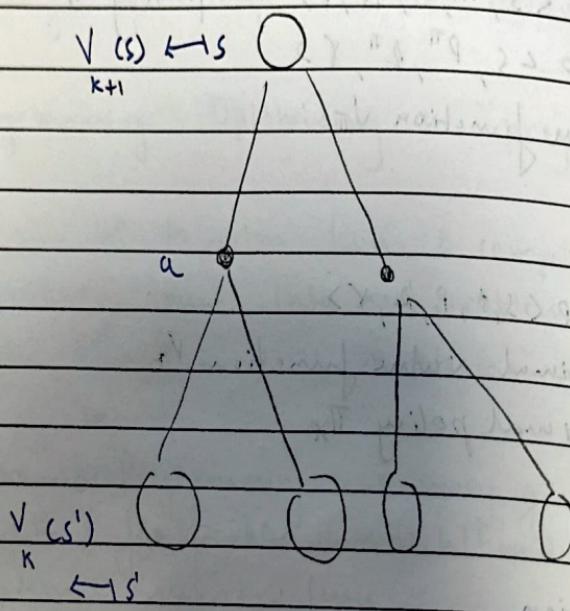
We'll do this through synchronous backup i.e.

we sweep through all states of our bellman equation

## Synchronous Backup

- At each iteration  $k+1$
- For all states  $s \in S$
- Update  $v_{k+1}(s)$  from  $v_k(s')$
- where  $s'$  is the successor state of  $s$

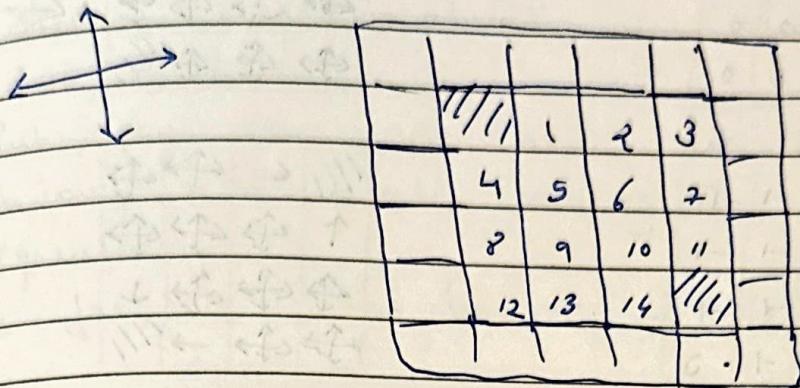
### D) The Bellman Expectation Equation



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

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

## Evaluating "Random Policy in the small grid World"



- Undiscounted MDP (i.e  $\gamma = 1$ )
- Non-terminal states = 1 - - - 14
- One terminal state
- Actions leading out of grid leave state unchanged
- Reward is -1 until the terminal state is reached
- Agent follows a random policy

$$\pi(a_t | s_t) = \pi(a_t | e_t) = \pi(a_t | s_t) = \pi(a_t | w_s) = \frac{1}{4}$$

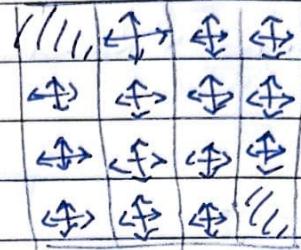
Here we're basically trying to answer what if we randomly keep moving, when would we reach the greyed box or how long would it take in expectation?

VK for the  
Random policy

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

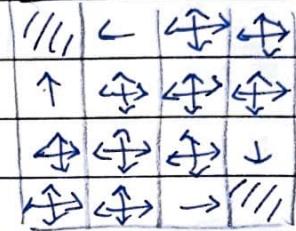
 $K=0$ 

Greek Policy

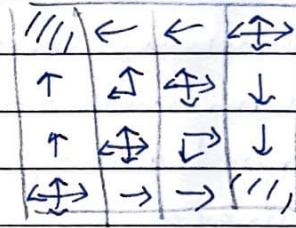
wrt  $\sqrt{K}$ 

← Random  
Policy

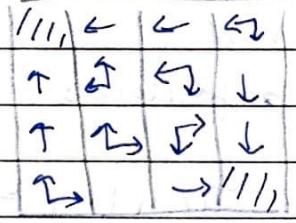
0	1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

 $K=1$ 

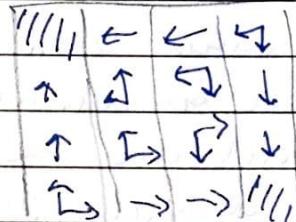
0	-1.7	-2	-2
-1.7	-2	-2	-2
-2	-2	-2	-1.7
-2	-2	-1.7	0

 $K=2$ 

0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.4
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0

 $K=3$ 

0	-6.1	-8.4	-9
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9	-8.4	-6.1	0

 $K=10$ 

Optimal  
Policy

0.0	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	0.0

 $K=50$ 

In the above diagrams the left column boxes symbolize the reward or value of each state, at the start we initialize  $v = 0$   $v_k$  with 0.

But you can see how we iterate through all states using the bellman expectation equation

e.g. for  $k=1$ , the box \* is calculated as

$$v_{k=1}^* = (-1 + 0) \times \frac{1}{4} + (-1 + 0) \times \frac{1}{4} + (-1 + 0) \times \frac{1}{4} + (-1 + 0) \times \frac{1}{4}$$

$\downarrow$   
-1 the

reward

for taking

a step  $\rightarrow 0$  is the

previous value of

this state we go into after

the action

2 we do it randomly for each dir

2 avg it

$$\therefore v_{k=1}^* = -1$$

You slowly start to notice how values update as each iteration progresses.

With  $k=0$ , the values seem to have converged

2 you notice on avg it takes about 18-20 steps if you follow a random policy in this grid world

Now, based on the  $V_K$  we calculate, the eight column boxes tell us that instead we try to pick actions greedily depending on our reward we start to notice that the directions on the converged  $V_{K=3}$  are same as for  $V_{K=3}^{k=\infty}$  which implies that we find an optimal policy well before  $K = \infty$ .

### f) Policy Iteration

We wish to make our policy better - How do we do it?

→ Given a policy  $\pi_t$

We evaluate the policy  $\pi_t$

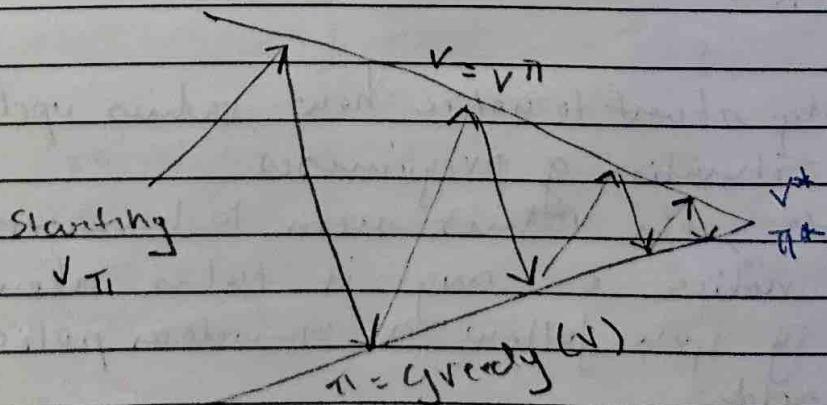
$$V_{\pi}(s) = E[R_{t+1} + \gamma R_{t+2} + \dots | s_t = s]$$

Improve the policy by acting greedily wrt to  $V_{\pi}$   
 $\pi' = \text{greedy}(V_{\pi})$

→ In a small gridworld improved policy was optimal,  
 $\pi' = \pi^*$

→ In general, need more iterations of improvement / evaluation

→ But this process of policy iteration always converges to  $\pi^*$



- We start off with some arbitrary input to our value function like  $v$  & some policy.  
 So we shall do policy evaluation on our up arrows & policy improvement on the down arrows.
- So first we calculate our value function wrt to our original policy then we act greedily or update policy depending on our value function & we repeat this till we converge to optimal value function & policy.

evaluation

$$v \rightarrow v^*$$

$$\pi \rightarrow \text{greedy}(\pi)$$

improvement

$$\pi^* \rightarrow v^*$$

So it does not matter where we start with what value function or policy we will always end up converging to the optimal policy or value function.

#### 4) Policy Improvement

- Consider a deterministic policy,  $a = \pi(s)$
  - We can then improve the policy by acting greedily.
- $$\pi'(s) = \operatorname{argmax}_{a \in A} q_{\pi}(s, a)$$
- Taking the argmax improves the value from any one state  $s$  over atleast one step.

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

This above equation simply means to determine if our greedy policy is atleast better for one step than just our original policy i.e.  $q_{\pi}(s, \pi'(s))$  says if we take a greedy step  $(\pi'(s))$  from state  $s$  i.e. a single step & continue with the original policy  $\pi$  there on, it is better than if we had continued with the original policy  $v_{\pi}(s)$  from the given state  $s$ .

i.e.  $q_{\pi}(s, \pi'(s))$  tell us simply how much value we get if we start in  $s$ , follow action 'a' for one step (the greedy policy  $\pi'$ ) & then follow  $\pi$  afterwards.

$v_{\pi}(s)$  tells us the value we get following  $\pi$  from  $s$ .

→ Since  $\max_{a \in A} q_{\pi}(s, a)$  gives us the max value

of action from state  $s$  so that does tell us  $q_{\pi'}(s, \pi'(s))$

so max over all actions has to be atleast as good as our particular action we choose based on  $\pi$ .

It thus improves the value function;  $V_{\pi'}(s) > V_{\pi}(s)$

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

This tells us that if we pick this greedy policy, the total amount of reward we get by being greedy is atleast as good as the policy we had before greedyfying it.

→ If improvement stops,

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

Then the Bellman optimality eq. has been satisfied

$$V_{\pi}(s) = \max_{a \in A} q_{\pi}(s, a)$$

thus,  $\pi(s) = \star$  for all  $s \in S$

so thus if policy  $\pi(s) = \max_{a \in A} q_\pi(s, a)$

then  $\pi(s)$  is optimal &  $\pi$  is optimal.

#### 4) Modified Policy Iteration

In the earlier grid world example we noticed that we already found the optimal policy at  $k=3$ , then why  $k=4, \dots$  as value iterations needed?

∴ Does policy evaluation need to converge to  $V_\pi$ ?

- We could introduce a stopping condition e.g.  $\epsilon$ -convergence of value function, if we are close in the changing values of the bellman equation by a tiny amount  $\epsilon$  we stop.
- Or simply we stop after  $K$  iterations
- Why not update policy every iteration? i.e. stop after  $K$  iterations  
This is equivalent to value iteration

## I) Principal of Optimality

Any optimal policy can be subdivided into two components

- An optimal first action  $A_*$ , so if your first action is optimal & you follow an optimal policy from where you end up, the overall behaviour is optimal

### - Theorem

A policy  $\pi^*(a|s)$  achieves the optimal value from state  $s$ ,  $v_{\pi^*}(s) = v^*(s)$ , if & only if

- For any state  $s'$  reachable from  $s$
- $\pi^*$  achieves the optimal value from state  $s'$ ,  $v_{\pi^*}(s') = v^*(s')$

## I) Deterministic Value Iteration

- If we know the solution to subproblems  $v^*(s')$
- Then solution  $v^*(s)$  can be found by a one step look ahead

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

The idea of value iteration is to apply these updates iteratively

Intuition : Start with final rewards & work backwards

### k) Value Iteration

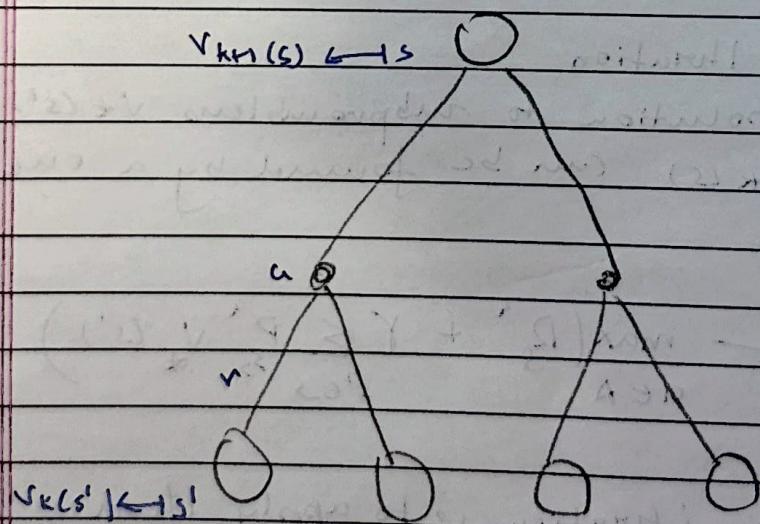
Problem: find optimal policy  $\pi^*$

Solution: iterative application of bellman optimality equation back-up

$$V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V^*$$

Using Synchronous backups

- At each  $K+1$  iteration
- For all states  $s \in S$
- Update:  $V_{K+1}(s)$  from  $V_K(s')$
- Unlike policy iteration, there is no explicit policy
- Intermediate value function may not correspond to any policy



$$V_{K+1}(s) = \max_{u \in A} \left( R_s^u + \gamma \sum_{s'} P_{ss'}^u V_K(s') \right)$$

$$V_{K+1} = \max_{u \in A} R^u + \gamma P^u V_K$$

## 1) Understanding it All:

Policy Evaluation:

Here we are given a fixed policy  $\pi$  & we compute  $V\pi$

$$V_{\pi+1}(s) = \sum_a \pi(a|s) \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V_\pi(s') \right)$$

The policy does not change here, we are given a policy that tells us how & with what probability we take action (eg. Random in grid world - eg. or rather in each dirn N,S,W,E) based on this fixed policy we keep calculating  $V_k$ 's for  $k=1, 2, \dots, \infty$  till the values converge & thus  $V_\infty = V\pi$  or the evaluation for that particular policy.  
 Basically answering the question if we commit to this policy forever, how good is it?

## Policy Iteration & Improvement

Now the next step is to ask: knowing how good your policy is, how do we improve it?

For this first we look at the policy improvement theorem where we define

$$q_n(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V_n(s')$$

$$\text{& a new policy } \pi'(s) = \arg \max_a q_n(s, a)$$

which guarantees  $V_{\pi}(s) \geq V_{\pi'}(s)$   $\forall s$

Now why,  $q_{\pi}(s, a)$  tells us the value of taking an action  $a$  now & then following  $\pi$ . However  $q_{\pi'}(s, \pi'(s))$  tells us instead of taking action  $a$  based on the policy we take, a bound on  $\pi'(s)$  which is  $\max q_{\pi'}(s, a)$  & then following  $\pi$ .

Choosing the action with the maximum  $q_{\pi}$  cannot logically make things worse, since the  $V_{\pi}(s')$  of the later part of the equation is characterizing all future rewards. Thus, the max  $q_{\pi}(s, a)$  cannot necessarily be  $\downarrow$  the avg.  $q_{\pi}(s, a)$   
 $\pi'(s)$

using the policy improvement theorem tells us if we act greedily we get a policy that's worse than before.

Now comes the Improvement.

First we take a policy, any initialized policy, then valuation on that policy  $\pi$  thus  $V_{\pi, k=1}, V_{\pi, k=2}, \dots, V_{\pi, k=\infty}$  or it converges & we yet  $V_{\pi, k=\infty}$  then we update our policy  $\pi$ , to the new  $\pi'$ , based on greedily selecting the action, & we repeat process till  $\pi$  stops changing & we get some policy back.

This leads to many evaluations & policy updates  
 & can be expensive so one possibility is  
 modified Policy iteration.

Here we either evaluate till  $K=1, 5, 10$  steps  
 & then update policy greedily & repeat.

However there is one extreme case  $K=1$  this  
 is value iteration.

Value Iteration means we evaluate the  
 policy for just one sweep  $K=1$  & then  
 the evaluation is produced by acting  
 greedily immediately & we never store  
 policy

$$v_{k+1} = \max_a (R_s^a + \gamma \sum_{s'} P_{ss'}^a v_k(s'))$$

Once we have  $v^*$  or  $v$  converges,  
 finding optimal policy is trivial because

$$\pi^*(s) = \arg \max_a (R_s^a + \gamma \sum_{s'} P_{ss'}^a v^*(s'))$$

M) Problem.

Bellman Equation

Algorithm

Prediction

Bellman Expectation Eq<sup>n</sup>

Iterative

control

Bellman Expectation Eq<sup>n+1</sup>

Policy Evaluation

Greedy policy improvement

Policy Iteration

control

Bellman Optimality eq<sup>n</sup>

Value Iteration

- Each of these algorithms are based on state-value function  $V_{\pi}(s)$  or  $V_{\pi^*}(s)$
- So time complexity is  $O(mn^2)$  per iteration for  $m$  actions &  $n$  states because there are  $n$  states & for each  $n$  states there are  $n$  possible successor states.
- Could also be applied to action-value function  $q_{\pi}(s, a)$  or  $q_{\pi^*}(s, a)$
- Complexity  $O(m^2n^2)$  per iteration.

## N) Asynchronous Dynamic Programming

In this case all states are backed up in parallel, Asynchronous DP backs up states individually & in any order

We apply the appropriate backup for each selected state

This reduces computation significantly

This is guaranteed to converge provided all states are contained to be selected

Ideas for Asynchronous DP

→ Inplace DP

Synchronous value iteration stores two copies of the value function for all  $s \in S$

$$V_{\text{new}}(s) \leftarrow \max_{a \in A} \left( R^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\text{old}}(s') \right)$$

$$V_{\text{old}} \leftarrow V_{\text{new}}$$

Inplace only stores one copy of value function for all  $s \in S$

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

So here once we compute the  $V(s)$  for a state if that is part of the computation for another state, we use its updated value & not its previous or original value

So what matters here is the order of updation, so if you have an order of whom you update the state closer to the goal & keep going backwards, this helps you be efficient. & This order of updation actually matters.

### → Prioritized Sweeping

So given the fact ordering is important we need to come up with a measure to identify which states are really important & which one should you update first.

So we can use the magnitude of the Bellman error to guide our state selection eg

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

so we tend to propagate those states first which have the highest emerging bellman error

### Realtime - DP

Have a real agent to get relevant state information use the agents experience to guide state selection