

# Planning by Dynamic Programming:

Lecture 3

Daniel Silver

1. introduction,
2. Policy Evaluation.
3. Policy iteration
4. value iteration
5. extensions to dynamic programming
6. Contraction mapping
- Critical Thinking

Dynamic sequential or temporal component to the problem.

Programming optimizing a "program" i.e. a policy.

★ CF know programming

★ A method for solving complex problems

★ By breaking them down into subproblems.

- Solve the subproblems

- Combine solutions to subproblems

★ Dynamic programming is a very general solution method for problems which have two properties:

★ Optimal substructure

- Principle of optimality applies.

- Optimal solutions can be decomposed into subproblems

★ Overlapping subproblems

- Subproblems recur many times

- Solutions can be cached and reused

★ Markov decision processes satisfy both properties.

- Bellman equations give recursive decomposition

- value function stores and reuses solutions.

- Dynamic programming assumes full knowledge of the MDP.

- It is used for planning in an MDP

- For prediction

Input MDP:  $\langle S, A, P, R, V \rangle$  and policy  $\pi$

s.a.m



- Or, MDP  $\langle S, P^{\pi}, R^{\pi}, V \rangle$
- Output, value function  $V_{\pi}$
- Or, for control
  - Input, MDP  $\langle S, P, R, \gamma \rangle$
  - Output, optimal value function  $V_{*}$
  - and, optimal policy  $\pi_{*}$

\* Dynamic programming is used to solve many other problems eg.

- \* Scheduling algorithms.
- \* String algorithms (sequence alignment)
- \* Graph " (shortest path)
- \* Graphical models (Viterbi algorithm)
- \* Bioinformatics

→ **Problem:** evaluate a given policy  $\pi$

→ **Solution:** iterative application of Bellman expectation backup.

$V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V_{\pi}$

→ arbitrary initial value

• easy synchronous backups

• At each iteration  $k+1$

• for all states  $s \in S$

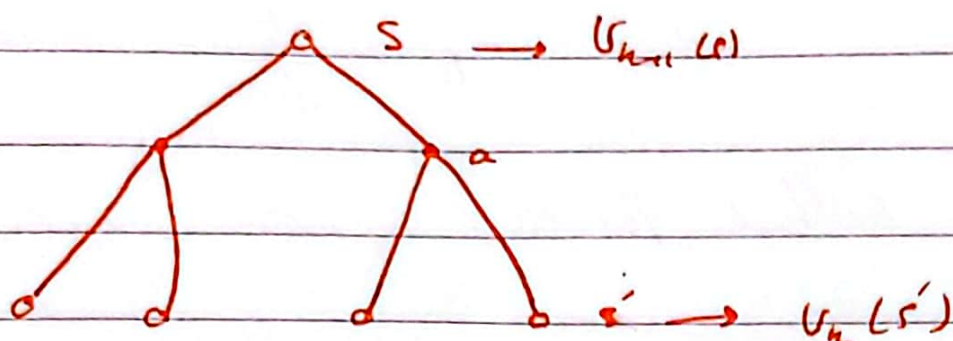
• update  $V_{k+1}(s)$  from  $V_k(s')$

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

• we will discuss asynchronous backups later

• Convergence to  $V_{\pi}$  will be proven at the end of the lecture.





$$U_{h+1}(s) = \sum_{a \in A} \pi(a|s) (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a U_h(s'))$$

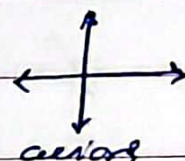
$$U_{h+1} = R^\pi + \gamma P^\pi U_h$$

this is a unique value function that satisfies this for any  $\pi$ .

Example: Small Grid world

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

$\gamma = 1$   
for all transitions



• undiscounted episodic MDP ( $\gamma = 1$ )

Normal states 1, ..., 14

One terminal state (shown twice as shaded squares)

Actions leading out of grid leave state unchanged.

Reward is 1 until the terminal state is reached

Agent follows uniform random policy

$$\pi(a|1) = \pi(a|2) = \pi(a|3) = \pi(a|4) = 0.25$$

Small Grid world

s.a.m



it is not optimal policy, we try to improve it

## Policy iteration

done by a feedback process.

- Given a policy  $\pi$

Evaluate the policy  $\pi$

$$V_{\pi}(s) = E[R_{t+1} + \gamma V_{\pi}(R_{t+2} + \gamma V_{\pi}(R_{t+3} + \dots | S_{t+1})]$$

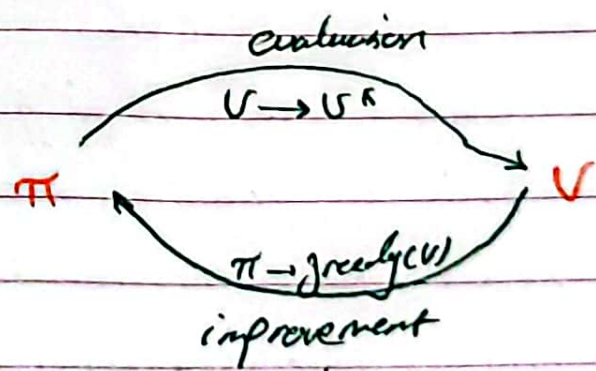
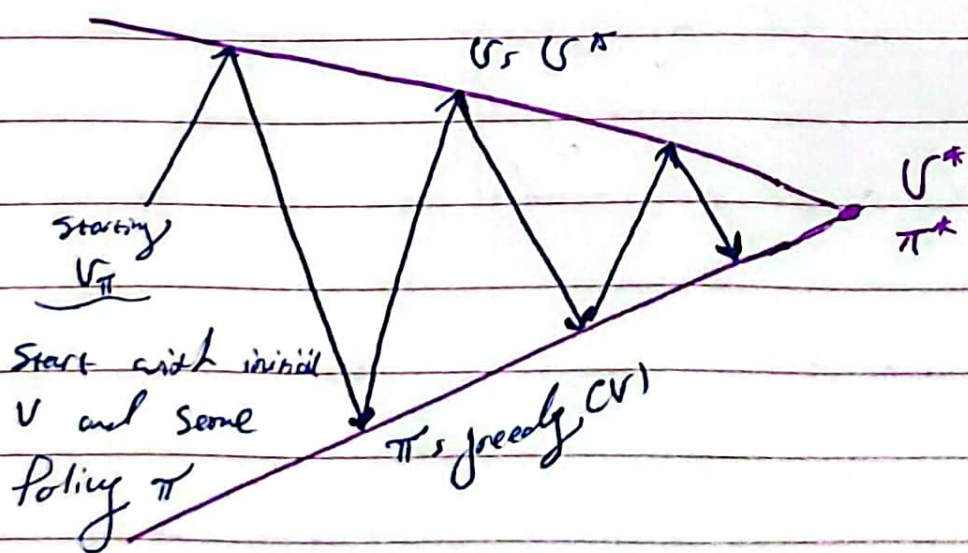
Improve the policy by acting greedily with respect to  $V_{\pi}$   
 $\pi' = \text{greedy}(V_{\pi})$

In small grid world improved policy was optimal,  $\pi' = \pi^*$ .

In general, need more iterations of improvement/evaluation.

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

★ Recall from last lecture: there is always at least one deterministic optimal policy for any MDP.



**Policy evaluation** Estimate  $V_{\pi}$

Iterative policy evaluation.

**Policy improvement** Generate  $\pi' \geq \pi$

Greedy policy improvement

sam



## Example:

- States: Two locations, maximum of 20 cars at each.
- Actions: Move up to 5 cars between locations overnight. (from location A to B)
- Rewards: \$10 for each car rented (must be available)
- Transitions: Transitions: Cars returns and requested randomly.
  - Poisson distribution,  $n$  returns/requests with prob  $\frac{\lambda^n}{n!} e^{-\lambda}$
  - 1st location: average requests = 3, average returns = 3
  - 2nd " : " " = 4, " " = 2

The Content Map: How Car Rentals' Policy, Discrete  
of Car Rental

## Important facts:

How sensitive is that the convergence rate to  $\pi_0$ ?

A: we'll see later at least for some of the algorithms the convergence rate is independent of  $\pi_0$ , but in practice that doesn't answer the question which is that you might just have your  $\pi_0$  your initial value may be way off and your first policy may be way off and so clearly if you start with the optimal policy you'll get to the optimal policy faster, so it does matter but the convergence rate is independent of the initial value and policy.

## a little bit formally.

- Consider a deterministic policy,  $a = \pi(s)$
- we can improve the policy by acting greedily,  
$$\pi'(s) = \arg \max_{a \in A} q_{\pi}(s, a)$$

acting greedily means that we basically are looking at the value of being in a state and taking a particular action and then following your policy after that.

sam



- This improves the value from any state s over one step.

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

→ if we are greedy about the greedy policy at least improves the value that we're getting over one step and this immediate one step, it's not wrong, get about what happens after that let's just consider one step and see whether we get more value over one step.

- It therefore improves the value function,  $V_{\pi'}(s) \geq V_{\pi}(s)$

$$V_{\pi}(s) \leq q_{\pi}(s, \pi(s)) = E_{\pi}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s]$$

$$\leq E_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, \pi'(S_{t+1})) | S_t = s]$$

$$\leq E_{\pi}[R_{t+1} + \gamma R_{t+2} + \gamma^2 q_{\pi}(S_{t+2}, \pi'(S_{t+2})) | S_t = s]$$

$$\leq E_{\pi}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s] = V_{\pi'}(s)$$

if improvements stop,

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

Then the Bellman optimality equation has been satisfied

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

Therefore  $V_{\pi}(s) = V_{\pi^*}(s)$  for all  $s \in S$

So  $\pi$  is an optimal policy.

Q: Does policy evaluation need to converge to  $V_{\pi}$ ?

Or should we introduce a stopping condition  $\overline{V_{\pi}}$

\* eg  $\epsilon$ -convergence of value function

Or simply stop after  $k$  iterations of iterative policy evaluation?

For example, in the small gridworld this was sufficient so

s.a.m



achieve optimal policy.

why not update policy every iteration? i.e. stop after  $k=1$ .  
This is equivalent to value iteration (next section)

## Value Iteration

Any optimal policy can be subdivided into two components:

- An optimal first action  $A_*$
- Followed by an optimal policy from successor state  $s'$

### Theorem (Principle of Optimality)

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

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

- if we know the solution to subproblems  $V_*(s')$
- Then solution  $V_*(s)$  can be found by one step lookahead

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

- The idea of value iteration is to apply these updates again and again and again iteratively.
- Intuition: Start with bad rewards and work backwards
- Still works with noisy, stochastic MDPs

- Problem: find optimal policy
- Solution: iterative application of Bellman optimality backup
- $V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V_*$
- Using synchronous backups
  - At each iteration  $k=1$

sam



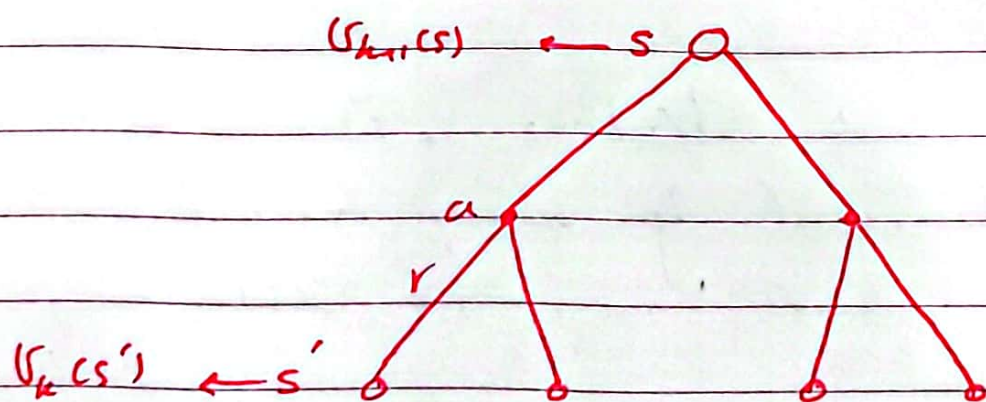
- For all States  $s \in S$
- Update  $U_{k+1}(s)$  from  $U_k(s')$
- Convergence so  $U^*$  will be proven later
- Unlike policy iteration, there is no explicit policy
- Intermediate value functions may not correspond to any policy

Difference between value iteration and policy iteration

one thing which is apparent is that we're not building an explicit policy at every step we're just working directly in value space. we kind of used our value function to build a policy used our policy to build the next value function.

value iteration goes directly from value function to value function so value function to value function

↳ may not correspond to particular  $U^*$



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

$$U_{k+1} = \max_{a \in A} (R^a + \gamma P^a U_k)$$

s.a.m



Problem	Bellman Equation	Algorithm
Prediction	Bellman Expectation Equation	Iterative Policy Evaluation
Control	" " " " + Greedy Policy Improvement	Policy Iteration
Control	Bellman Optimality Equation	Value Iteration

- ★ Algorithms are based on state-value function  $V_{\pi}(s)$  or  $V_{*}(s)$
- ★ Complexity  $O(mn^2)$  per iteration for  $m$  actions and  $n$  states
- ★ Could also apply to action-value functions  $Q_{\pi}(s, a)$  or  $Q_{*}(s, a)$
- ★ Complexity  $O(m^2n^2)$  per iteration  
high complexity

### Extensions:

- DP methods described so far used synchronous backups
- i.e. all states are backed up in parallel
- Asynchronous DP backs up states individually, in any order
- For each selected state apply the appropriate backup
- Can significantly reduce computation
- Guaranteed to converge if all states continue to be selected

### Three simple ideas for asynchronous dynamic programming:

- In place dynamic programming
- Prioritized Sweeping
- Real-time dynamic programming



## Replace Dynamic Programming

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

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

✓ Replace value iteration only stores a copy of value function for all  $s$  in  $S$

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

→ how important it is to update any state in your mdp, which state should be prioritised sweeping. you update first.

Use magnitude of Bellman error to guide state selection, eg

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

Backup the state with the largest remaining Bellman error

✓ Update Bellman error of affected states after each backup

✓ Requires knowledge of reverse dynamics (predecessor state)

✓ Can be implemented efficiently by maintaining a priority queue.

## Real-Time Dynamic Programming

Idea: only states that are relevant to agent → select the states that the agent actually visits.

Use agent's experience to guide the selection of states.

After each time-step  $s_t, A_t, R_{t+1}$  → collect real samples.

Backup the state  $s_t$

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

sam



\* DP uses full width backups

✓ For each backup (sync or async)

- Every successor state and action is considered

- using knowledge of the MDP transition and reward function

✓ DP is effective for medium-sized problems

✓ (millions of states)

✓ For large problems DP suffers Bellman's curse of dimensionality

- number of states  $n \propto |S|$  grows exponentially with number of state variables

✓ Even one backup can be too expensive