

# Chapter 4: Dynamic programming

SB: Chap 4

4-1

## 4.1 Policy evaluation

- Goal: Compute  $V_\pi$  for given policy  $\pi$
- Direct method: Solve linear equation

$$V_\pi = \rho_\pi + \gamma P_\pi V_\pi$$

$$\begin{pmatrix} V_\pi \end{pmatrix} = \begin{pmatrix} \rho_\pi \end{pmatrix} + \gamma \begin{pmatrix} P_\pi \end{pmatrix} \begin{pmatrix} V_\pi \end{pmatrix}$$

$$\Rightarrow V_\pi = (I - \gamma P_\pi)^{-1} \rho_\pi$$

- Complexity:  $O(|S|^3)$

- Iterative method:  $V_0 \rightarrow V_1 \rightarrow \dots \rightarrow V_\pi$

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

$$V_{k+1} = \rho_\pi + \gamma P_\pi V_k$$

Bellman backups

matrix product

$$= T_\pi(V_k)$$

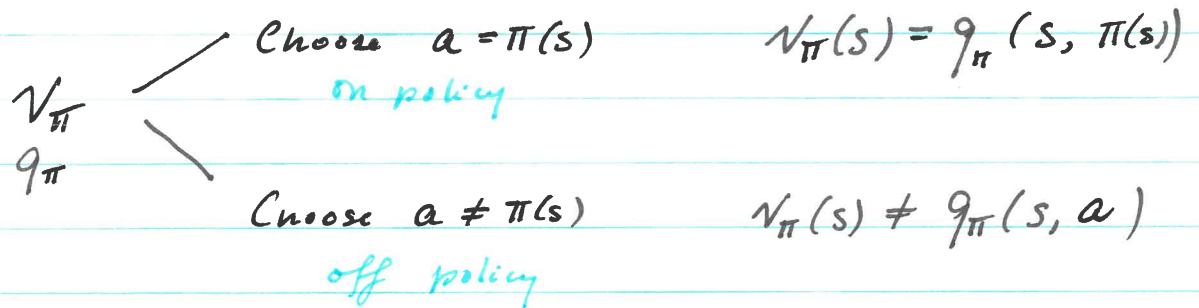
$$T_\pi(v) = \rho_\pi + \gamma P_\pi v$$

Bellman operator

- Initial value  $V_0$  arbitrary, except  $V_0(\text{terminal state}) = 0$
- $V_\pi$  fixed point of  $T_\pi$
- Convergence:  $V_k \rightarrow V_\pi$  as  $k \rightarrow \infty$
- $T_\pi$  is a contraction
- Complexity:  $|S|^2 \times M$   
matrix product # iterations
- Not exact  $V_\pi$  for finite # iterations

## 4.2. Policy iteration

- Consider deterministic policies  $\pi(s)$



- Policy improvement theorem: Let  $\pi$  and  $\pi'$  such that

$$q_{\pi'}(s, \pi'(s)) \geq q_{\pi}(s, \pi(s))$$

off policy
on policy

for all  $s \in \mathcal{S}$ . Then  $\pi'$  is as good or better than  $\pi$  in the sense that

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

better value under  $\pi'$

for all  $s \in \mathcal{S}$ .

- Greedy policy selection:

$$\pi'(s) = \arg \max_a q_{\pi}(s, a)$$

- Equation for updating  $\pi$
- Fixed point:  $\pi_*$  for  $q_*$  (optimal policy)
- Strict improvement unless optimal

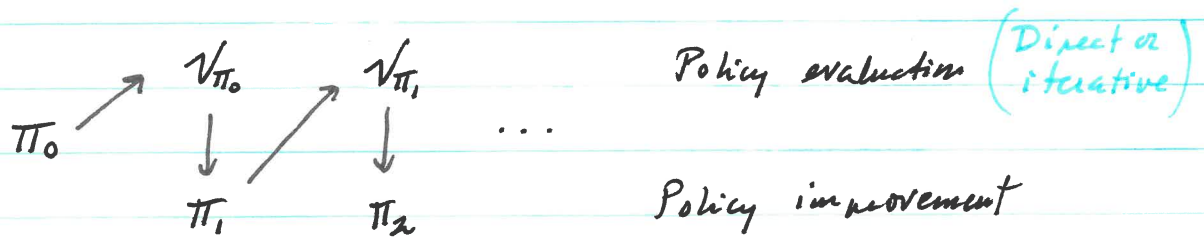
- $\epsilon$ -greedy policy selection:

$$\pi'(s) = \begin{cases} \arg \max_a q_{\pi}(s, a) & \text{Prob } 1-\epsilon \\ u \in \mathcal{S}' & \text{Prob } \epsilon \end{cases}$$

explanation

- No strict improvement

# Policy iteration algorithm



Convergence:

$$\begin{aligned} \pi_k &\rightarrow \pi_* \\ V_k &\rightarrow V_* \\ q_k &\rightarrow q_* \end{aligned} \quad \text{as } k \rightarrow \infty$$

Note: Policy evaluation for  $q_\pi$

$$\begin{aligned} q_\pi(s, a) &= r(s, a) + \gamma \sum_{a', s'} q_\pi(s', a') \pi(a'/s') p(s'/s, a) \\ &= r(s, a) + \gamma \sum_{s'} q_\pi(s', \pi(s')) p(s'/s, a) \\ &= r(s, a) + \gamma \sum_{s'} V_\pi(s') p(s'/s, a) \end{aligned}$$

on policy

$$\rightarrow q_\pi = r + \gamma P_a V_\pi$$

Note:

## Planning

- Model known
- Solve to get  $V_\pi$
- Solve to get  $V_*, \pi_*$

$$V_\pi, V_*$$

## Learning

- MDP model unknown
- Estimate  $V_\pi$  by exploration/sampling
- Converge to  $V_*, \pi_*$

$$q_\pi, q_*$$



## 4.4. Action value iteration

$$q_{k+1}(s, a) = E[ R_{t+1} + \underbrace{\delta \max_{a'} q_k(s_{t+1}, a')}_{\substack{\text{previous} \\ \text{estimate} \\ v_k(s_{t+1})}} \mid S_t = s, A_t = a]$$

time  $t+1$   
↓ ↓

update

$$q_{k+1}(s, a) = p(s, a) + \delta \max_{a'} \sum_{s'} q_k(s', a') p(s' | s, a)$$

- Max defines policy improvement at stage  $k+1$
- Fixed point:  $q_*$  Bellman optimality equation
- Convergence:  $q_k \rightarrow q_*$  as  $k \rightarrow \infty$

• Policy improvement:  $\pi_{k+1}(s) = \arg \max_a q_k(s, a)$

• Variant:

$$\underbrace{q_k(s, a)}_{q_k} = E[ R_{t+1} + \underbrace{\delta v_k(s_{t+1})}_{v_k} \mid S_t = s, A_t = a]$$

$q_k \leftarrow$

$v_k$

$$q_k(s, a) = p(s, a) + \sum_{s'} v_k(s') p(s' | s, a)$$

$$= p(s, a) + \delta (P_a v_k)(s)$$

action  $a$

- Max included in  $v_k$