

How to solve MDPs?

One way is Dynamic Programming

↳ policy evaluation

↳ policy iteration

↳ value iteration

DP → break problem into subproblems

solve them & concatenate

Bellman equation gives

recursive decomposition

Dynamic Programming assumes full knowledge of MDP, used for planning

for prediction MDP $\langle S, A, P, R, \gamma \rangle + \pi$ → IPE → value function v_π

for control MDP $\langle S, A, P, R, \gamma \rangle$ → → optimal value function v_* & optimal policy π_*

Policy Evaluation

method: iterative application of Bellman expectation

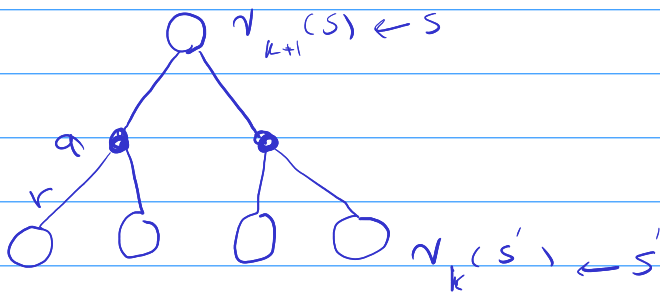
using synchronous backups:

at each iteration $k+1$:

for all states $s \in \mathcal{S}$:

update $v_{k+1}(s)$ from $v_k(s')$

where s' is the successor state



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

$$v^{k+1} = R^\pi + \gamma R^\pi v^k \quad [\text{gridworld example green}]$$

policy evaluation: estimate v_π , iterative policy evaluation

Policy Iteration

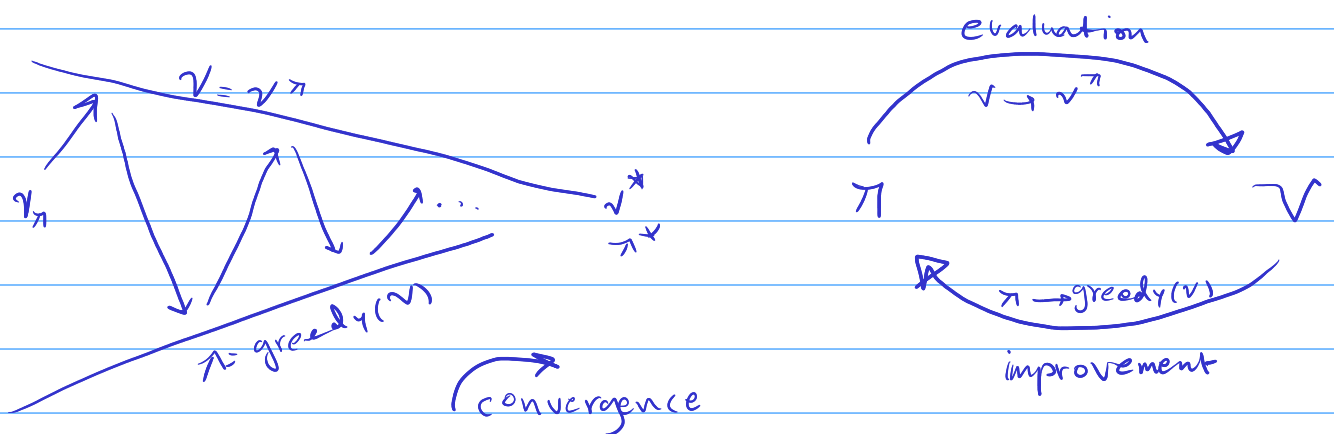
given a policy

→ evaluate the policy $\pi \left[v_{\pi}(s) = E[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s] \right]$

→ Improve the policy by acting greedily with v_{π}
 $\pi' = \text{greedy}(v_{\pi})$

Policy improvement: Generate $\pi' \succ \pi$

greedy policy improvement



Convergence?

considering a deterministic policy $a = \pi(s)$
we can improve the policy by acting greedily $\pi'(s) = \arg\max_{a \in A} q_{\pi}(s, a)$
which improves the value from any step s over one step
 $q_{\pi}(s, \pi'(s)) = \max_{a \in A} q_{\pi}(s, a) \geq q_{\pi}(s, \pi(s)) = v_{\pi}(s)$
and therefore improves the value function $v_{\pi'}(s) \geq v_{\pi}(s)$

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

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 Bellman optimality equation becomes satisfied

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

Principle of Optimality

a policy $\pi(a|s)$ achieves optimal value from state s $V_\pi(s) = V_*(s)$ iff
for any s' reachable from s ,
 π achieves the optimal value from s' , $V_\pi(s') = V_*(s')$

→ at each state, the optimal policy includes
taking the best action +
following the optimal policy from the next state

Value Iteration (1)

if we know the solution to subproblem $V_*(s')$

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

(one-step lookahead)

Basically start with final rewards and work backwards

① ... \rightsquigarrow ... ⑥

in practice goal state may not be obvious, but
it works by state sweep over entire state space

Value Iteration (2)

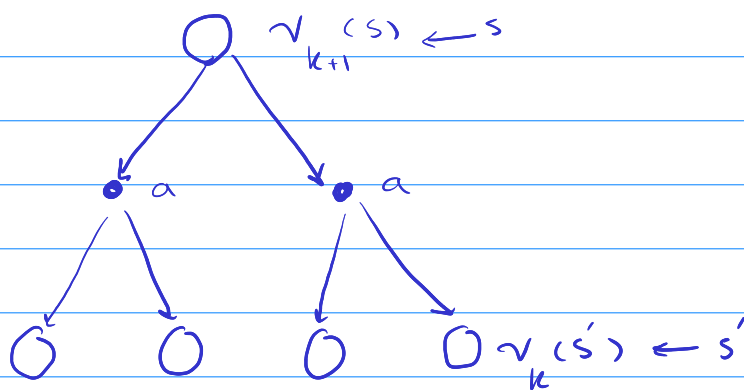
$$v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_*$$

at each iteration $k+1$

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

update $v_{k+1}(s)$ from $v_k(s')$

* intermediate value function may not correspond to any policy



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

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

Problem	equation	algorithm
prediction	Belman Expectation	iter policy eval
control	BE + Greedy policy impr	policy iteration
control	Belman Optimality	value iteration