

David Silver Lecture on Dynamic Programming

- Agenda:
 - ① Policy Evaluation
 - ② Policy Iteration
 - ③ Value Iteration
 - ④ Extension to DP

- Introduction: DP = Optimization for sequential problems
 - Breaks problems into subproblems and reassembles subproblems.
 - DP requires two properties:
 - ① Optimal Substructure: Optimal solution to problem constructed from optimal solution to subproblems
 - ∴ You can divide problem into subproblems
 - ② Overlapping Subproblems
 - ∴ You can cache solutions to subproblems to save time.

- Bellman Equation gives decomposition into subproblems
 - Recursive decomposition
- Caching done by value function: Maximum reward from any state onward.

- Planning: We are given structure of MDP and we want to solve it.
 - Prediction: Input: MDP $\langle S, A, P, R, \gamma \rangle$ and policy π
Output: Value Function

→ POLICY EVALUATION

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

Output : What's the best policy π^* and optimal value function v^*

- Policy Evaluation: Given MDP and policy, evaluate value function

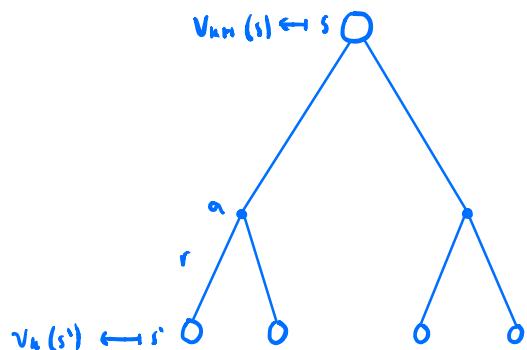
→ Use Bellman Expectation Equation to do policy evaluation

$$v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_n$$

↑
Value function
(start off with
arbitrary value,
typically zero)

→ Use synchronous backup

→ Sweep over all states in every iteration



$$V_{k+1}(s) := \sum_{a \in A} \pi(a|s) \left(R_a + \gamma \sum_{s' \in S} P_{s,s'} V_k(s') \right)$$

→ Turn this into iterative update

- Value function helps us determine better policies (more to come)

- Iterative Policy Evaluation:

Given a policy:

① Initialize all scores to some value

② Iteratively estimate the value function at each state s as a probability sum of value functions from all states s' which can transition to s , plus the reward associated with the transition to s .

→ Iterate Bellman Equation

③ Repeat step 2 to convergence

• Policy Iteration:

• Given a policy π , how can you guarantee improvement to π ?

→ Solve this in 2 steps:

① Evaluate Policy (as in previous section)

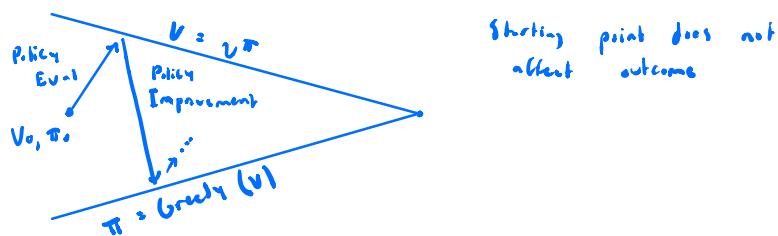
② Improve policy by acting greedily w.r.t. V_π

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

→ Iterates Evaluate and Improvement steps

→ Returns a deterministic optimal policy

• After initial policy, all other policies are deterministic



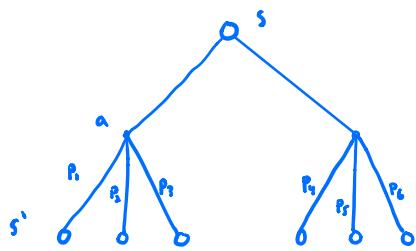
• Example Problem: Car rentals and returns at two locations

→ Redistribute cars each night

- ↳ What policy maximizes income
- ↳ State = # of cars in each lot
- ↳ Action = # of cars to move from 1 location to the other

- Process:
- Begin with some policy: Trivial policy is move 0 cars
 - Evaluate policy to build up value function
 - Take the action that yields the highest value function and use this as your policy in the next iteration
 - Reiterate

• Take a Step Back: Policy Evaluation



$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) (R_{ss} + V_{\pi}(s'))$$

Take the probability weighted sum of rewards associated with all possible transitions.

- Precise Example: Consider a deterministic policy $a = \pi(s)$

- How do we improve policy? Act Greedily!

$$\pi'(s) = \underset{a \in A}{\text{argmax}} q_{\pi}(s, a)$$

Use the value function based on current policy
Iterate over all actions

- ↳ Pick actions in a way that maximize action value q

q is immediate reward with action a + reward at next state.

→ Prove that this improves action value in one step

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

Follow policy π' for 1 step, then follow π afterwards

Definition of max

Definition of argmax that defined π'

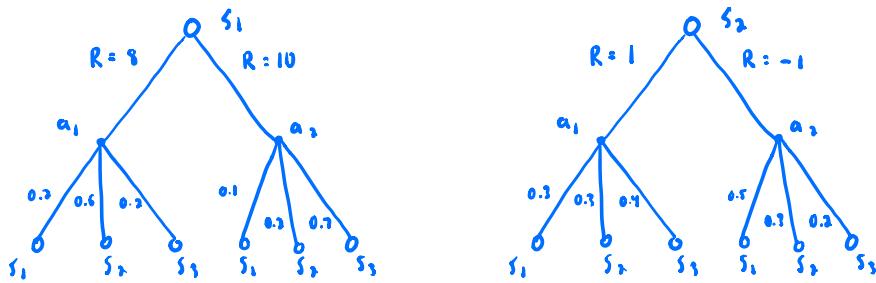
∴ Value function improves over 1 step.

∴ $v_{\pi'}(s) \geq v_{\pi}(s)$

Show greedy policy selection improves value for one step

Continue at 50:00

• Value Iteration in Homework Problem : Assignment 4, Problem 1



- Initialize Scores at maximum over actions

$$v_0(s) = \frac{s_1}{10} = \frac{s_2}{1} = \frac{s_3}{0}$$

- At each step $i+1$, take the Bellman operator of the value function of each state s at i . Estimate value function of each state s as a probability weighted sum of transition rewards and rewards at s'

$$q(s_i, a_i) = R_{ss'} + \underbrace{\sum_{s' \in A} p_{ss'} v_{i+1}(s')}_{\text{Bellman Operator}} =$$

$$q(s_1, a_1) = 8 + (0.2 \cdot 10 + 0.6 \cdot 1 + 0.2 \cdot 0) = 10.6$$

$$q(s_1, a_2) = 10 + (0.1 \cdot 10 + 0.2 \cdot 1 + 0.7 \cdot 0) = 11.2 \leftarrow v_i(s_1)$$

$$q(s_2, a_1) = 1 + (0.3 \cdot 10 + 0.7 \cdot 1 + 0.4 \cdot 0) = 4.3 \quad \left. \right\} v_i(s_2)$$

$$q(s_2, a_2) = -1 + (0.5 \cdot 10 + 0.7 \cdot 1 + 0.2 \cdot 0) = 4.3 \quad \left. \right\} v_i(s_2)$$

$$q(s_1, a_1) = 8 + (0.2 \cdot 11.2 + 0.6 \cdot 4.3 + 0.2 \cdot 0) = 12.82$$

$$q(s_1, a_2) = 10 + (0.1 \cdot 11.2 + 0.2 \cdot 4.3 + 0.7 \cdot 0) = 11.98$$

$$q(s_2, a_1) = 1 + (0.3 \cdot 11.2 + 0.3 \cdot 4.3 + 0.4 \cdot 0) = 5.65$$

$$q(s_2, a_2) = -1 + (0.5 \cdot 11.2 + 0.7 \cdot 4.3 + 0.2 \cdot 0) = 5.89$$

$$|(V_{k+1} - V^*)(s)| \leq |V_k(s) - V^*(s)|$$

Policy Iteration

$\Pi \vdash MDP \rightarrow MRP$

$\forall MRP \exists V(s)$ s.t. $V = R + \gamma PV$

until convergence

$$V(s) = R + \gamma \sum_{s'} V(s') p(s, s')$$

Greedily Improve Policy

$$\Pi_{k+1}(s) = \arg\max_a \sum_i (R_{ss'i} + V(s')) p_{ssi}^a$$

$$\Pi^*(s) = \arg\max_a \sum_i (\hat{R}_{ss'i} + V^*(s')) p_{ssi}$$

MRP = Expected Discount Reward from a .

$\dots \in \mathbb{R}$

given own

Policy Evaluation

$$\Pi_k + MDP \rightarrow MRP$$

$$User \quad V = (I - \gamma P)^{-1} R$$

as inner loop in
pricing contraction

Concern BL Contraction

Mapping)