# Dynamic Programming (DP)

Compute optimal policies given a perfect model of the environment as a MDP.

Environment: finite MDP

$$p(s', r \mid s, a)$$

key idea: Use of value functions to organize and structure the search for good policies.

Bellman Optimality Equations:

$$V_*(s) = \max_a \mathbb{E}\left[ R_{t+1} + \gamma V_*(S_{t+1}) \mid S_t = s, A_t = a \right]$$

$$= \max_a \sum_{s', r} p(s', r \mid s, a)\left[ r + \gamma V_*(s') \right]$$

$$q_*(s, a) = \mathbb{E}\left[ R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

$$= \sum_{s', r} p(s', r \mid s, a)\left[ r + \gamma \max_{a'} q_*(s', a') \right]$$

Bellman Equations $\xrightarrow{\text{DP}}$ Assignments

(update rules for
~~improving approximations~~
~~of~~ the desired value
~~function~~)

<mark>Policy Evaluation (Prediction)</mark>

<u>iterative computation of the value functions for a given</u>
$$V_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s]$$ <u>policy</u>

$$= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

$$= \mathbb{E}_\pi[R_{t+1} + \gamma V_\pi(S_{t+1}) \mid S_t = s]$$

$$= \sum_a \pi(a \mid s) \sum_{s',r} p(s', r \mid s, a)[r + \gamma V_\pi(s')]$$

iterative policy evaluation:

$$V_{k+1}(s) \doteq \mathbb{E}[R_{t+1} + \gamma V_k(S_{t+1}) \mid S_t = s]$$

$$= \sum_a \pi(a \mid s) \sum_{s',r} p(s', r \mid s, a)[r + \gamma V_k(s')]$$

( expected update ) on each state

"based on an expectation over all possible next states
rather than on a sample next state"

$$V \leftarrow V(s)$$

$$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

"sweep" through the state space

## Policy Improvements (Theorem)

"We know how good it is to follow the current
policy from s, that is $V_\pi(s)$, but would it
be better or worse to change to the new
policy?"

Computation of an improved policy given the value
function for that policy

$$q_\pi(s,a) \doteq \mathbb{E}[R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t = s, A_t = a]$$

$$= \sum_{s',r} p(s',r|s,a) [r + \gamma V_\pi(s')]$$

$$\boxed{\begin{array}{l} q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s) \\[2mm] V_{\pi'}(s) \geq V_{\pi}(s) \end{array}}$$

$$V_{\pi}(s) \leq q_{\pi}(s, \pi'(s))$$

$$= \mathbb{E}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s, A_t = \pi'(s)]$$

$$= \mathbb{E}_{\pi'}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s]$$

$$\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s]$$

$$= \mathbb{E}_{\pi'}[R_{t+1} + \gamma \mathbb{E}_{\pi'}[R_{t+2} + \gamma V_{\pi}(S_{t+2}) \mid S_{t+1}, A_{t+1} =$$

$$\pi'(S_{t+1})] \mid S_t = s]$$

$$= \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{\pi}(S_{t+2}) \mid S_t = s]$$

$$\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 V_{\pi}(S_{t+3})$$

$$\mid S_t = s]$$

$$\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \cdots \mid S_t = s]$$

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

Greedy policy $\pi'$:

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

$$= \arg\max_a \mathbb{E}\left[ R_{t+1} + \gamma V_\pi(S_{t+1}) \mid S_t = s, A_t = a \right]$$

$$= \arg\max_a \sum_{s', r} p(s', r \mid s, a)\left[ r + \gamma V_\pi(s') \right]$$

$$V_\pi = V_{\pi'}$$

$$V_{\pi'}(s) = \max_a \mathbb{E}\left[ R_{t+1} + \gamma V_{\pi'}(S_{t+1}) \mid S_t = s, A_t = a \right]$$

$$= \max_a \sum_{s', r} p(s', r \mid s, a)\left[ r + \gamma V_{\pi'}(s') \right]$$

The process of making a new policy that __improve__ on an original policy, by making it greedy w.r.t. the value function of the original policy.

## Policy Iteration

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \cdots$$

$$\xrightarrow{I} \pi_* \xrightarrow{E} V_*$$

E: policy evaluation

I: policy improvement

## Value Iteration

Policy iteration: involve policy evaluation

$$\left( \begin{array}{c} \text{"require multiple sweeps through} \\ \text{the state set"} \end{array} \right)$$

"policy improvement + truncated policy evaluation"

"special case" of policy evaluation

$$V_{k+1}(s) \doteq \max_a \mathbb{E}\left[R_{t+1} + \gamma V_k(S_{t+1}) \mid S_t = s, A_t = a\right]$$

$$= \max_a \sum_{s',r} p(s', r \mid s, a)\left[r + \gamma V_k(s')\right]$$

$$\begin{cases} v \leftarrow V(s) \\ V(s) \leftarrow \max_a \sum_{s',r} p(s', r \mid s, a)\left[r + \gamma V(s')\right] \\ \Delta \leftarrow \max(\Delta, |v - V(s)|) \end{cases}$$

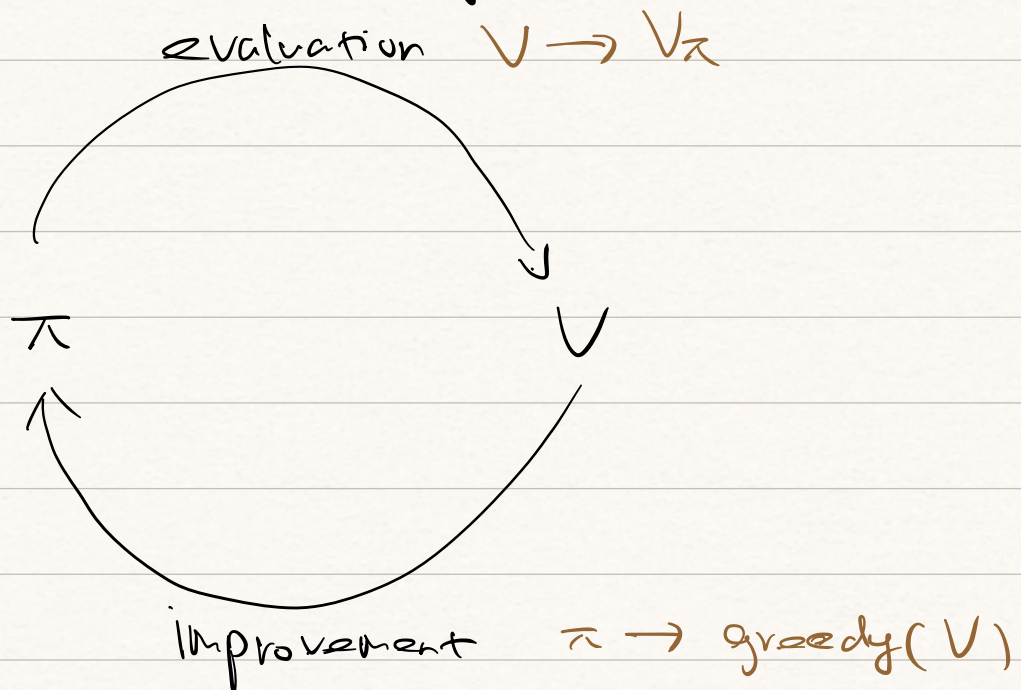$$\pi(s) = \arg\max_a \sum_{s',r} p(s', r \mid s, a)\left[r + \gamma V(s')\right]$$
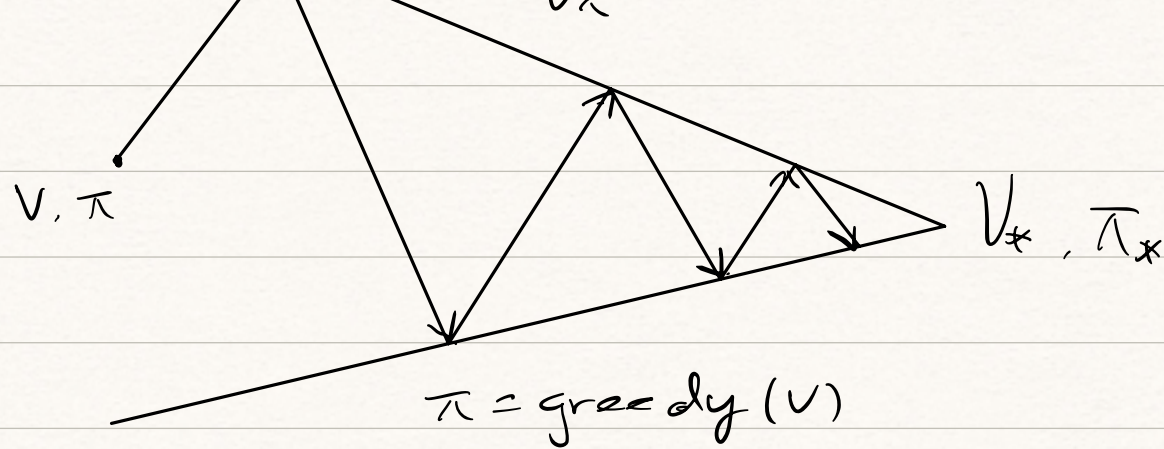
# Asynchronous DP

in-place iterative DP algorithm

Problem of policy iteration:  $\underline{\text{require sweeps}}$
$\text{of the state set}$

update the value of only one state, $S_k$, on each step $k$, using the value iteration update.

## Generalized Policy Iteration (GPI)

evaluation $V \rightarrow V_\pi$

$\pi$

$V$

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

$V = V_*$

$V, \pi$

$V_*, \pi_*$

$\pi = greedy(V)$

Bootstrapping !

All of them update estimates of the values of
states based on estimates of the values of
successor states. That is, they update estimates
on the basis of other estimates.