Dynamic Programming (DP)
- a collection of algorithms to compute
optimal policies, given a model An
the dynamics of the environment: p(s,rls,a)
-> assume a finite MDP (episodic task)
- hey idea: use value functions as a basis
to search for policies
- recall Bellman eq's:
vx(s) = max E[R++++y V*(S+++)|S+=s, A+=a]
= max
$$\sum_{s',r} p(s',rls,a)[r+y V*(s')]$$

q*(s,a) = E[R++++ y max q*(S++,a')|S+=s, t=a]
= $\sum_{s',r} p(s',rls,a)[r+y max q*(s'a')]$
DP algorithms: turn Bellman eqs into
update rules tor improved
approximation to V*, q*

- /-

- Policy Evaluation:
compute value direction
$$V_{\pi}(s)$$
 for a given
fixed policy $\pi(a/s)$.
vecall: $V_{\pi}(s) = I_{\pi}[G_{t}]S_{t}=s]$
 $= \sum \pi(a/s) \sum p(s',r/s,a)[r+y U_{\pi}(s')]$
 \Rightarrow system of $|S|$ simultaneons linear exclusions
for $|S|$ unknowns $V_{\pi}(s)$
 \Rightarrow alternative approach: solve iteratively, i.e.
find a sequence $V_{0}, V_{1}, V_{2}, \dots$ of value
functions $, s.t = V_{\mu}(s) \xrightarrow{k \to \infty} V_{\pi}(s) + s$.
 $V_{\mu+1}(s) := I_{\pi}[R_{k+1} + y V_{\mu}(S_{k+1}) | S_{k} = s]$
 $= \sum \pi(a/s) \sum p(s',r | s,a)[r+y V_{\mu}(s')]$
has a fixed point at $V_{\mu} = V_{\pi}$
 \Rightarrow iterative policy evalutation:
. each iteration steps updates $V_{\mu}(s) + s \in S$
in-place

Iterative Policy Evaluation, for estimating $V \approx v_{\pi}$

Input π , the policy to be evaluated Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation Initialize V(s), for all $s \in S^+$, arbitrarily except that V(terminal) = 0

Loop: (iferations)

$$\Delta \leftarrow 0$$
Loop for each $s \in S$: (states)
 $v \leftarrow V(s)$
 $V(s) \leftarrow \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

$$\frac{HW}{HW}: do Exercise 4.3$$

$$\frac{Policy}{Policy} \frac{Policy}{Policy} \frac{Policy}{Policy}$$

-3-

$$\pi'(a/s) = \begin{cases} 1, \text{ for } a = \overline{a} \\ \overline{p} \\ 20, \text{ otherwise} \end{cases} (=) \overline{a} = \pi'(s)$$

deterministic

$$\frac{P_{root}I}{V_{\pi}(s)} = q_{\pi}(s, \pi'(s))$$

$$\frac{(4et)}{F} E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_{t}=s, A_{t}=\pi'(s)]$$

$$\pi' deter.$$

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

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

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

$$= E_{\pi'}[R_{t+1} + \gamma E[R_{t+2} + \gamma V_{\pi}(S_{t+2}) | S_{t+1}, A_{t+1}=\pi'(S_{t+1})] | S_{t}=s]$$

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

$$= E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^{2}V_{\pi}(S_{t+2}) | S_{t+3}] | S_{t}=s]$$

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

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

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

-> given policy
$$\pi$$
, consider very greedy policy:
 $\pi'(s) := \arg\max [\pi(s, \alpha)]$
 $= \arg\max \mathcal{F}[R_{t+1} + \int V_{\pi}(S_{t+1})/S_{t} = s, A_{t} = \alpha]$
 $= \arg\max \sum_{s,r} p(s'_{r}|s_{r}\alpha)[r + \int V_{\pi}(s')]$
• by construction, π' meets the requirement
of the policy improvement theorem
 $=> \pi' = \pi$
Here $V_{\pi'} = V_{\pi}$, and from def. of π' we have
 $V_{\pi'}(s) = \max \mathcal{F}[R_{t+1} + \int V_{\pi'}(S_{t+1})|S_{t} = s, A_{t} - \alpha]$
 $= \max \sum_{s,r} p(s'_{r}r|s_{r}\alpha)[r + \int V_{\pi'}(s')]$
note: policy improvement them. holds also
for stochastic policies:
if $\sum \pi'(\alpha|s) q_{\pi'}(s, \alpha) \geq V_{\pi'}(s)$

-S-

Policy Iteration

$$T_{o} \xrightarrow{E} V_{T_{o}} \xrightarrow{I} T_{I} \xrightarrow{E} V_{T_{o}} \xrightarrow{T} \sum V_{T_{o}} \xrightarrow{T} V_{T_{o}} \xrightarrow{E} V_{T_{o}}$$

iterate between policy evaluation (E) &
policy improvement (I)
. each operation E, I itself is iterative

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$ 1. Initialization $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in S$ 2. Policy Evaluation Loop: $\Delta \leftarrow 0$ Loop for each $s \in S$: $v \leftarrow V(s)$ $V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) \left[r + \gamma V(s')\right]$ $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ until $\Delta < \theta$ (a small positive number determining the accuracy of estimation) 3. Policy Improvement policy-stable $\leftarrow true$ For each $s \in S$: $\textit{old-action} \leftarrow \pi(s)$ $\pi(s) \leftarrow \arg\max_{a} \sum_{s',r} p(s',r \,|\, s,a) \big[r + \gamma V(s') \big]$ If old-action $\neq \pi(s)$, then policy-stable \leftarrow false If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Value Iteration
->do we need to run each E step intil
convergence? -> NO!
· combine E & I steps:
Vk+: (s) = max / E (R++, + y Vu (S++,) | S+= s, A+=a]
= max
$$\sum_{s':v} p(s':r/s,a) [v + y Vu(s')]$$

· turn Bellman eq. into update rule ittelf

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation Initialize V(s), for all $s \in S^+$, arbitrarily except that V(terminal) = 0

Loop:

$$\begin{array}{l} \Delta \leftarrow 0 \\ | \text{ Loop for each } s \in \mathbb{S}: \\ | v \leftarrow V(s) \\ | V(s) \leftarrow \max_a \sum_{s',r} p(s',r \,|\, s,a) \left[r + \gamma V(s')\right] \\ | \Delta \leftarrow \max(\Delta, |v - V(s)|) \\ \text{until } \Delta < \theta \end{array}$$
Output a deterministic policy, $\pi \approx \pi_*$, such that $\pi(s) = \operatorname{arg\,max}_a \sum_{s',r} p(s',r \,|\, s,a) \left[r + \gamma V(s')\right] \end{array}$

-8-

-> drawbacks;