# Reinforcement Learning (Planning Algos)

**Planning Algorithms:** given knowledge of $P(s'|s,a)$ and $\mathbb{E}\, r(s,a)$  $\forall a \in A, s, s' \in S$

how do we compute the optimal **policy** $\pi^*$ ?

Sometimes called "Dynamic Programming" — used very broadly. Means different things to different people

## ① Value Iteration

In mid '90s, Puterman said this algo was most widely used, due to its simplicity, but he advises against using it
(due to slow convergence)

Bellman's Optimality Eq'n:

not a scalar, so "max" seems ill-defined, but $\pi$ decouples across $s$, so interpret component-wise

$(\forall s \in S)$   $V_{\pi^*}(s) = \max\limits_{a \in A} \left( \mathbb{E}[r(s,a)] + \gamma \sum\limits_{s' \in S} P(s'|s, a) \cdot V_{\pi}(s') \right)$

$= \max\limits_{\pi} \left( R_\pi + \gamma P_\pi \cdot V \right)$   Since choosing action $a$, based on $s$, is same as a **policy**

$:= \Phi(V)$  ← looks funny since $V$ depends on $\pi$, but OK if we require...

**Fixed Pt. Eq'n**

$V^* = \Phi(V^*)$   (Bellman's Eq'n, restated)

↗ non-linear operator

**Algorithm:** Value Iteration   aka Picard iteration

Initialize $V_0$ arbitrarily

Iterate $V_{t+1} = \Phi(V_t)$ for $t = 1, 2, \dots$

(stopping criteria: $\| V_{t+1} - \Phi(V_t) \| < \frac{1-\gamma}{\gamma} \cdot \varepsilon$

## Thm (17.11) Value iteration converges to an optimal $V^*$

**proof**  Use contraction-mapping (aka Banach or Banach-Picard fixed pt.) theorem

i.e., show $\exists\, c < 1$ st $\forall\, V, U$  $\| \Phi(V) - \Phi(U) \| \leq c \| V - U \|$

(any norm, as long as it a Banach space... in fact complete metric space works)

We'll show $\Phi$ is Lipschitz cts w/ constant $c = \gamma$, and use $\| \cdot \| = \| \cdot \|_\infty$

↳ I'm going to abuse notation to simplify it... you can make it rigorous by fixing a row (i.e., fix $s$)

$\Phi(V) - \Phi(U) = \max\limits_{\pi} (R_\pi + \gamma P_\pi V) - \max\limits_{\pi'} (R_{\pi'} + \gamma P_{\pi'} U)$

(argmax...)

$= R_\pi + \gamma P_\pi V - \max\limits_{\pi'} (R_{\pi'} + \gamma P_{\pi'} U)$

$\leq R_\pi + \gamma P_\pi V - (R_\pi + \gamma P_\pi U)$   since $\tilde{\pi}$ suboptimal

$= \gamma P_\pi (V - U)$

and similarly $\Phi(U) - \Phi(V) \leq \gamma P_{\pi'} (U - V)$

So $\|\Phi(v) - \Phi(u)\|_\infty \leq \gamma \|P_\pi (v-u)\|_\infty$

$\qquad\qquad\qquad \leq \gamma \|P_\pi\|_\infty \cdot \|v-u\|_\infty$   by def. operator norm

$\qquad\qquad\qquad\qquad\quad \underbrace{\qquad}_{\to \,=1 \text{ regardless of which } \pi}$

$\qquad\qquad\qquad = \gamma \cdot \|v-u\|_\infty$   $\square$

## How do we get $\pi^*$ back from $V^*$?

Bellman Eq'n,

$$V^*(s) = \max_{a \in A} \; \mathbb{E}\, r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) \cdot V^*(s)$$

$\qquad\qquad\quad \underset{\text{assuming } |A| < \infty}{} \quad \underbrace{\;\;}_{\text{known}} \qquad\qquad \underbrace{\;\;}_{\text{known}} \quad \underbrace{\;\;}_{\text{known}}$

$\Rightarrow$ we can solve the maximization problem to find $a$

and this $a$ is $\pi^*(s)$.

## Value iteration w/ Gauss-Seidel acceleration

$\qquad\qquad \underbrace{\qquad\qquad\qquad}_{\text{if you don't know this, just ignore this part}}$

$\simeq$ same cost per iteration, faster convergence

Always converges in our case due to properties of transition matrices,
$\qquad$ c.f. Puterman Thm 6.3.7

## ② Policy Iteration

Instead of solving for $V$ and $\pi$ is implicit, now work directly w/ $\pi$

### ALGORITHM: Policy Iteration

Initialize $\pi_0$ arbitrarily

For $t = 1, 2, \dots$ *

$\qquad$ Find $V_{\pi_t}$ by solving $(I - \gamma P_{\pi_t}) V = R_{\pi_t}$  // expensive! Solve a linear system of equations

$\qquad$ $\pi_{t+1} = \arg\max_\pi \left( R_\pi + \gamma P_\pi \cdot V_{\pi_t}\right)$  // "greedy" update

$\qquad\qquad$ ↳ same as "How do we get back $\pi^*$ from $V^*$"

$\qquad$ break if $\pi_{t+1} = \pi_t$

*Policy iter. terminates in a finite # of steps
$\qquad \nearrow$ if $|S| = n$, policy iter. is $O(n^3)$/iter.
$\qquad\qquad$ value iter. is $O(n^2 \cdot |A|)$/iter.

You can rewrite as $V_{t+1} = V_t - \left(\gamma P_{\pi_R} - I\right)^{-1} \left(\underset{\text{vague for now}}{\Phi(V_t)} - V_t\right)$   c.f. Puterman

so at a fixed pt, $0 = \underbrace{(\gamma P_\pi - I)^{-1}}_{\text{nonsingular}} \left(\Phi(v) - v\right)$

$\qquad\qquad \Rightarrow \Phi(v) = v$  (Bellman's Eq)

also like a *Newton method*  $f(v) = \Phi(v) - v$   ... and in practice, policy iter.
$\qquad\qquad\qquad\qquad\qquad f'(v) = (\gamma P_\pi - I)$   converges much faster than value iter.

A common practical version "modified policy iteration"

$\qquad$ solve $(I - \gamma P_{\pi_t}) v = R_{\pi_t}$ approximately, using the fact

$\qquad (I - \gamma P)^{-1} = \sum_{k=0}^\infty (\gamma P)^k$  (Neumann Series)

$\qquad$ so $(I - \gamma P)^{-1} R \approx \sum_{k=0}^A (\gamma P)^k R$ $\qquad$ compute $S = R$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $R_1 = \gamma P \cdot R$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $S \leftarrow S + R_1$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $R_2 = \gamma P \cdot R_1$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $S \leftarrow S + R_2$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\dots$

### Monotonicity* property

$\qquad$ If $V \geq U$ (meaning $V(s) \geq U(s) \; \forall s$)

$\qquad$ then $\Phi(v) \geq \Phi(u)$

For our proof, we'll need $v \geq u \Rightarrow (I-\gamma P_\pi)^{-1} v \geq (I-\gamma P_\pi)^{-1} u$    (for $\pi = \pi_t$)

**proof** since $\forall \pi, \|P_\pi\|_\infty = 1 \Rightarrow$ (Neumann Series converges) $\left( \text{since } \|\gamma P_\pi\|_\infty = \gamma < 1 \right)$

$$(I-\gamma P_\pi)^{-1} = \sum_{k=0}^{\infty} (\gamma P_\pi)^k$$

So $(I-\gamma P_\pi)^{-1}(v-u) = \sum (\gamma P_\pi)^k (v-u) \geq 0$   □

$\underbrace{\qquad}_{\geq 0 \text{ entries}}$

**Lemma 17.12** If $(V_t)$ constructed via policy iteration, then $V_t \leq V_{t+1} \leq V^*$ $(\forall t)$

**proof** because $\pi_{t+1}$ chosen greedily (to maximize $R_\pi + \gamma P_\pi \cdot V_t$)

$$R_{\pi_{t+1}} + \gamma P_{\pi_{t+1}} \cdot V_t \geq R_{\pi_t} + \gamma P_{\pi_t} \cdot V_t = V_t$$

so $R_{\pi_{t+1}} \geq (I - \gamma P_{\pi_{t+1}}) V_t$

hence $\underbrace{(I-\gamma P_{\pi_{t+1}})^{-1} R_{\pi_{t+1}}}_{V_{t+1}} \geq (...)^{-1} (...) V_t = V_t$    via monotonicity    □

**Thm** Policy Iteration converges to an optimal policy (for a finite MDP)

**Proof** By greedy update, if $V_{t+1} = V_t \Rightarrow$ satisfy Bellman's Eq'n so optimal

By lemma above, $V_t \leq V_{t+1}$, so

$V_t$ not optimal $\Rightarrow V_t < V_{t+1}$

Note for a finite MDP, # possible policies $= |A|^{|S|} < \infty$

and we can't repeat any policies since $V_t < V_{t+1}$

$\Rightarrow \exists$ maximal policy $V_{t^*}$, and $V_{t^*+1} = V_{t^*} \Rightarrow$ it's optimal.    □

**Corollary:** converge in $\leq |A|^{|S|}$ iterations

③ Linear Programming

Bellman Eq'n: $(\forall s)$ $V^*(s) = \max_{a \in A} \overbrace{R_a(s)}^{\mathbb{E} \, r(s,a)} + \gamma \sum_{s' \in S} \overbrace{P_a(s'|s)}^{\mathbb{P}(s'|s,a)} \cdot V^*(s')$

i.e., $\forall a,$
$$V^*(s) \geq R_a(s) + \gamma \sum_{s' \in S} P_a(s'|s) V^*(s')$$

or in matrix notation

$$\vec{V} \geq \vec{R}_a + \gamma P_a \cdot \vec{v} \quad (\forall a)$$

↗ a set of $|S|$ linear inequalities

↑ So $|A|$ of these $|S|$ sets of ineq.

So

Find $\vec{V} \in \mathbb{R}^{|S|}$ s.t. $\vec{V} \geq \vec{R}_a + \gamma P_a \vec{v}$   $\forall a \in A$

is a Linear Program (if a finite MDP).

$\exists$ weakly polynomial time algo to solve LP (1980's, ellipsoid, Karmarkar, etc.)

and in fact $\exists$ strongly polynomial time algo to solve the LP arising from MDP (Yinyu Ye 2000's)

In practice, due to extremely **large** state space S, LP formulation usually not efficient

## Summary

① For small MDP ( $|S|, |A|$ reasonably sized),
we can satisfactorily solve the planning problem

② Value Iteration : many iterations, cost $O(|S|^2 |A|)$/ iteration

Policy Iteration : few iterations, cost $O(|S|^3)$/ iteration

LP formulation : LP w/ $|S|$ variables, $|S| \cdot |A|$ constraints ⟩ I forgot exactly
(not relevant since
many special
tricks/algos)
complexity of LP something like variables$^3$· constraints$^2$

Backgammon : $|S| = 10^{20}$

Chess : $|S| = 10^{47}$

③ As $|S| \to \infty$, need special tricks.
Obviously still ongoing research

Go : $|S| = 10^{170}$ for $19 \times 19$ board

See wikipedia "Game complexity"