Markov Decision Process

"Model"  $(S, A, r, p)$          Policy $\pi(s) \in A$

states   actions   reward   transition

Planning

Reinforcement Learning

Given model, output
a policy.
(a) Finite horizon
(b) Infinite horizon

(b) Planning — Infinite Horizon          $\gamma \in [0, 1)$
                                          "Discount factor"

MDP value function

$$V^{\pi}(s) = \mathbb{E}_{s \sim p}\left[ \sum_{t=0}^{\infty} \gamma^t \, r(s_t, \pi(s_t)) \,\Big|\, s_0 = s \right]$$

Policy $\pi$ is better than or equal to
policy $\pi'$ if $V^{\pi}(s) \geq V^{\pi'}(s)$ in all states $s$.

Optimal policy $\pi^*$ better than or equal to
all policies $\pi$.

Optimal value function
$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

[4] Bellman equation / Principle of Optimality

$$(\square) \quad V^*(s) = \max_a \left[ r(s,a) + \gamma \sum_{s'} p(s'|s,a) V^*(s') \right]$$

$$\begin{array}{ccc} \text{Value} \\ \text{state} \end{array} = \begin{array}{c} \text{Value of} \\ \text{the optimal action} \end{array} + \begin{array}{c} \text{Value of} \\ \text{continuing} \\ \text{optimally} \end{array}$$

② Algorithm 1 : Value Iteration

Solve $(\square)$ iteratively

Initialize $V(s) = 0$ , all states $s$

Repeat $\quad V'(s) \leftarrow \max_a \left[ r(s,a) + \gamma \sum_{s'} p(s'|s,a) V(s') \right]$

all $s$

$V(s) \leftarrow V'(s)$ , all $s$

Theorem VI converges to the optimal value function $U^*$

Note Policy extraction

$$\pi^*(s) \in \arg\max_a \left[ r(s,a) + \gamma \sum_{s'} p(s'|s,e) V^*(s') \right]$$

Idea     Show that $\underline{V}' \leftarrow B(\underline{V})$, for
"Bellman operator" $B$, is a
contra'on.

---

(Aside)

Consider $f: \mathbb{R}^D \to \mathbb{R}^D$, update $x' \leftarrow f(x)$,
fixpoint $f(x^*) = x^*$.

Function $f$ is a <u>contraction</u> when

$$\| f(x) - f(y) \| < \| x - y \| \quad , \quad x \neq y$$

eg., $f(x) = \dfrac{x}{2}$. $\quad (8, 2) \quad (4, 1) \quad (2, \frac{1}{2}) \ldots$

     fixpoint $x^* = 0$

<u>theorem</u>   Given contraction property,
then $f(x)$ has a unique fixpoint and
$x' \leftarrow f(x)$ converges.

<u>Proof</u>   ⓐ $f$ has unique fixpoint, else

$$\| f(x^*) - f(y^*) \| = \| x^* - y^* \| \text{ , violate}_{\text{contraction}}$$

ⓑ $f$ ~~has~~ converge to fixpoint; consider $x \neq x^*$

$$\| f(x) - x^* \| = \| f(x) - f(x^*) \| < \| x - x^* \|$$

Fact    Bellman operator $B$ is a contraction for norm $\|\underline{V}\| = \max_s |V(s)|$

Comments

① VI converges to $V^*$ asymptotically

② Policy (Optimal) extracted from $V$ in a finite # of steps

③ Extraction

$$\pi(s) \in \arg\max_a \left[ r(s,a) + \gamma \sum_{s'} p(s'|s,a) V(s) \right]$$

③ Algorithm 2 : Policy Iteration

$$\pi^{(0)} \xrightarrow{E} V^{(0)} \xrightarrow{I} \pi^{(1)} \xrightarrow{E} V^{(1)} \xrightarrow{I} \pi^{(2)} \longrightarrow$$

Initialize $\pi^{(0)}$ (arbitrary)

Repeat  ① Evaluate $V^{\pi}$ ( $\pi$ is current policy)

      ② Improve:

$$\pi'(s) \leftarrow \arg\max_a \left[ r(s,a) + \gamma \sum_{s'} p(s'|s,a) V^{\pi}(s') \right]$$

for all s

$$\pi \leftarrow \pi'$$

Theorem  PI converges to optimal policy in a finite # of steps ( $\Rightarrow$ also get optimal value function)

(Note)  Suppose for state s, exists an a

$$r(s,a) + \gamma \sum_{s'} p(s'|s,a) V^{\pi}(s') > V^{\pi}(s)$$

( taking action once, then following $\pi$ is better than following $\pi$ )

Can show that $V^{(k+1)} > V^{(k)}$ if the policy changes.

# Notes

① Policy evaluation

Given $\pi$, can solve a system of equations
to get $V^\pi$.

$$V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) V^\pi(s') ; \text{all } s$$

$|S|$ unknowns , $|S|$ equations

In vector form:

$$\underline{V}^\pi = \underline{R}^\pi + \gamma \underline{P}^\pi \underline{V}^\pi \quad , \quad P^\pi \text{ is } |S| \times |S|$$
$$\text{transition matrix}$$

$$\Leftrightarrow (I - \gamma \underline{P}^\pi) \underline{V}^\pi = \underline{R}^\pi$$

$$\Leftrightarrow \underline{V}^\pi = \underbrace{(I - \gamma \underline{P}^\pi)^{-1}}_{\text{full rank}} \underline{R}^\pi$$

② PI requires fewer iterations than
   VI to solve for optimal $\pi^*$

Comparison:

VI  $V'(s) \leftarrow \max_a \left[ r(s,a) + \gamma \sum_{s'} p(s'|s,a) V(s') \right]$

   $O(|S||A| L)$ per iteration
   $L$ is the max # reachable
        next states

PI  $\pi'(s) \leftarrow \arg\max_a \left[ r(s,a) + \gamma \sum_{s'} p(s'|s,a) V^{\pi}(s') \right]$

   $O(|S||A| L + |S|^3)$ per iteration
                 $\underbrace{\qquad}_{\substack{\text{policy} \\ \text{evaluation}}}$

# Reinforcement Learning

Learn from environment, no knowledge of
$\nu$ or $p$ ("the model").

New challenge: explore (learning new things)

vs. exploit (leverage what you know to do well)

## Two main approaches

| Model-Based | Model-free |
|---|---|
| → Learn model (predict next state, reward) | → Don't learn model |
| → Use planning, to decide how to act | → Directly learn a policy, or an action-value function (Q-function) |
| ✓ Can accomodate changes in rewards, transitions | ✓ Simple, inexpensive computation |
| ⊘ Costly (lots of computation) | ⊘ If things change, have to "act a lot" to learn a new behavior |

⑤ Model-free RL: Value-based methods

💡 learn Q-function

$$Q^\pi(s,a) = r(s,a) + \gamma \sum_{s'} p(s'|s,a) V^\pi(s')$$

$$\underbrace{\text{value of taking } a, \text{ followed by policy } \pi}$$

$$Q^*(s,a) = r(s,a) + \gamma \sum_{s'} p(s'|s,a) V^*(s')$$

$$\pi^*(s) = \arg\max_a Q^*(s,a)$$

$$\underbrace{V^*(s')}$$

Bellman Equation

$$Q^*(s,a) = r(s,a) + \gamma \sum_{s'} p(s'|s,a) \max_{a'} Q^*(s',a')$$

| Reinforcement Learner |

Initialize $Q(s,a)$ values  (tabular)

Repeat:
  ① Act based on $Q$ values
  ② Use $s, a, r, s', a'$ to update
     $Q$ values
     [ Better approximate $Q^*$ ]

① Act : ε-greedy agent

$$\pi(s) \leftarrow \begin{cases} \arg\max_a Q(s,a) & w.\ prob\ 1-\varepsilon \\ random & w.\ prob\ \varepsilon \end{cases}$$

② Learn $Q^*$ through "Temporal difference" updates

Ⓐ SARSA   Each time get new experience
(on policy)   $s, a, r, s', a'$

$$Q(s,a) \leftarrow Q(s,a) + \alpha_t \left[ r + \gamma Q(s', a') - Q(s,a) \right]$$

learning rate
(period $t$)

(1-step estimate of $Q^\pi(s,a)$)

TD-error

Ⓑ Q-LEARNING   Each time get new
(off-policy)   experience $(s, a, r, s')$

$$Q(s,a) \leftarrow Q(s,a) + \alpha_t \left[ r + \gamma \max_{a'} Q(s', a') - Q(s,a) \right]$$

1 step estimate of $Q^*(s,a)$,

TD-error

(b) SARSA is "on policy" because it estimates $Q^\pi$ for policy followed (including $\epsilon$-exploration)

(c) Q-LEARNING is "off policy" because it estimates $Q^*$ while following $\pi$ (converge to $Q^*$ as long as $(s,a)$ visited often enough)