

Reinforcement Learning.

- States s_t
- transition function $p(s_{t+1}, r_{t+1} | s_t, a_t)$
- Actions a_t
- Episodes $s_0, a_0, s_1, a_1, \dots$
 - Episodic tasks $s_0, \dots, s_T \rightarrow s'$
 - Continuing tasks s_0, \dots

Returns $G_t = R_{t+1} + R_{t+2} + \dots$

Infinite continuing episodes, introduce a discount factor γ

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$0 \leq \gamma < 1$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

- Policy $\pi(a|s)$ = Probability of executing action a when in state s .

- Value function under policy π , $V_{\pi}(s) \rightarrow \mathbb{R}$

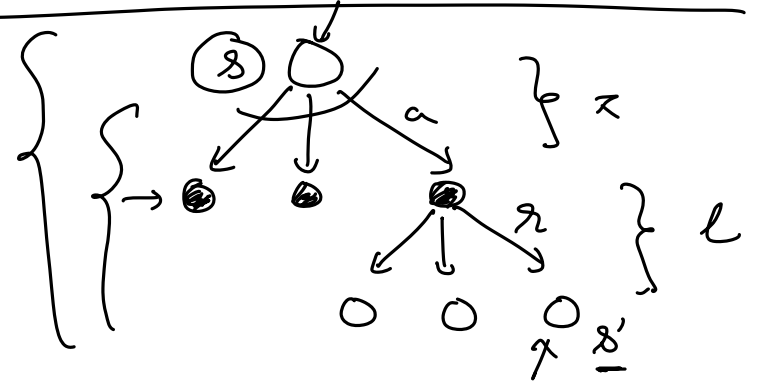
expected return when following policy π starting from state s .

$$V_{\pi}(s) = E_{\pi} [G_t | s_t = s]$$

$$= E_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma E_{\pi} [G_{t+1} | S_{t+1} = s']]$$

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



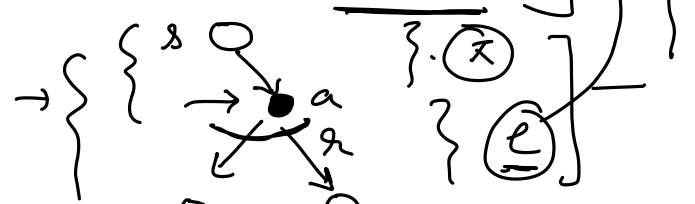
Action Value Function

$q_{\pi}(s, a)$ = Expected return if I were to execute action a from state s & then execute π .

$$q_{\pi}(s, a) = E_{\pi} [G_t | s_t = s, A_t = a]$$

$$= E_{\pi} [R_{t+1} + \gamma G_{t+1} | s_t = s, A_t = a]$$

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



$$V_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s, a)$$



Optimal value function

$$V_{*}(s) = \max_{\pi} V_{\pi}(s)$$

$$q_{\pi^*}(s,a) = \max_{\pi} q_{\pi}(s,a) \quad |$$

Optimal policies.

$\pi \geq \pi'$ if & only if $v_{\pi}(s) \geq v_{\pi'}(s) \quad \forall s \in \mathcal{S}$
 π^* is optimal if $\pi^* \geq \pi \quad \forall \pi$

$$q_{\pi^*}(s,a) = E [R_{t+1} + \gamma v_{\pi^*}(s_{t+1}) \mid S_t = s, A_t = a]$$

Bellman Optimality Equation

$$\begin{aligned} v_{\pi^*}(s) &= \max_a q_{\pi^*}(s,a) \\ &= \max_a E_{\pi^*} [G_t \mid S_t = s, A_t = a] \\ &= \max_a E_{\pi^*} [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\ &= \max_a E [R_{t+1} + \gamma v_{\pi^*}(s_{t+1})] \end{aligned}$$

$$v_{\pi^*}(s) = \max_a \sum_{s',r} p(s',r \mid s,a) [r + \gamma v_{\pi^*}(s')] \quad (1)$$

- if I give you $v_{\pi^*}(s)$, can you act in this MDP?

If instead, I gave you $q^*(s,a)$, can you act?

$$\leftarrow \arg \max_a \underline{q^*(s,a)}$$

Solving known MDPs.

- Policy evaluation

given a policy π , compute v_π

- Policy Improvement

given a policy π , obtain an improved policy π'

- Policy Iteration

- Value Iteration

Policy Evaluation

(15)

Given π . I want to evaluate $v_\pi(s)$

(a) Solving a system of equations

$$- v_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) \left[r + \gamma v_\pi(s') \right]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) r$$

$R^\pi(s)$

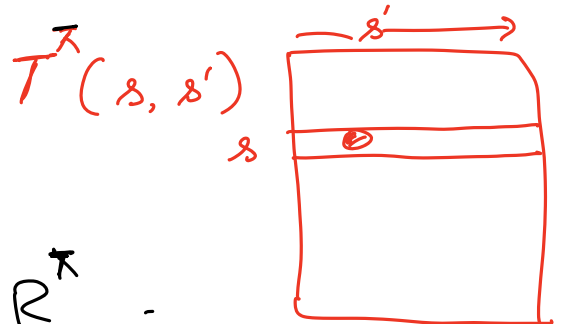
$$+ \gamma \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) v_\pi(s')$$

R^π

$$\sum_{s', r} \left[\sum_a p(s', r | s, a) \pi(a|s) \right] v_\pi(s')$$

$$V_\pi = R^\pi + \gamma T^\pi V_\pi$$

$$V_\pi = (I - \gamma T^\pi)^{-1} R^\pi$$

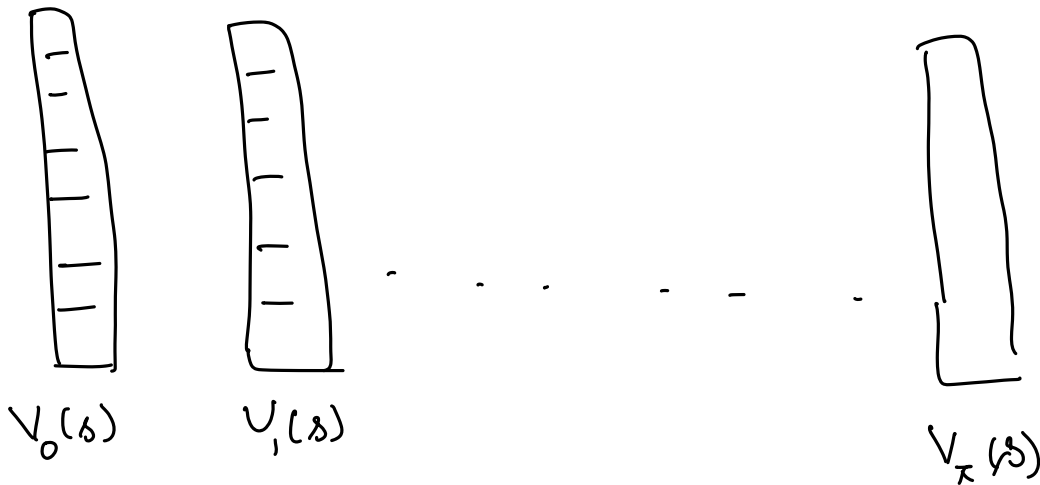


(b) Policy Iteration.

- Initialize $v_0(s) = 0$
- for $k = 1 \dots N$

for $s \in S$

$$v_{k+1}(s) \leftarrow \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_k(s')]$$



$$\Rightarrow \|V_{k+1} - U_{k+1}\|_{\infty} \leq \gamma \|U_k - U_k\|_{\infty}$$

γ -contraction

$$V_{k+1} \leftarrow R^{\pi} + \gamma T^{\pi} V_k$$

Policy Improvement

$$\pi, v_{\pi}(s) \rightarrow \pi'$$

Policy Improvement Theorem

Given a value function for a policy π , we can obtain a greedified policy π' :

$$\rightarrow \pi'(s) = \underset{a}{\operatorname{argmax}} \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

π better than π

$$\underline{v_{\pi'}(s)} \geq v_{\pi}(s) \quad \forall s$$

$$q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s) -$$

$$\underline{v_{\pi}(s)} \leq \underline{q_{\pi}(s, \pi'(s))} \dots \dots \dots v_{\pi'}(s)$$

$$= E_{\pi} [\underline{R_{t+1} + \gamma \underline{G_{t+1}}(s_t = s, A_t = \pi'(s))}] E_{\pi'}$$

$$= E_{\pi'} [R_{t+1} + \gamma \underline{v_{\pi'}(s)}_{t+1} \mid S_t = s]$$

$$\leq E_{\pi'} [R_{t+1} + \gamma \underline{q_{\pi}(s, \pi'(s_{t+1}))} \mid S_t = s]$$

$$= E_{\pi'} [R_{t+1} + \gamma \dots \dots \dots E_{\pi} [R_{t+2} + \gamma v_{\pi}(s_{t+2}) \mid S_{t+1}, A_{t+1} = \pi'(s_{t+1})]]$$

$$= E_{\pi'} [R_{t+1} + \gamma \dots \dots \dots E_{\pi'} [R_{t+2} + \gamma v_{\pi}(s_{t+2}) \mid \dots]]$$

$$\leq E_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 \dots]$$

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

for $\pi'(s) \leftarrow \operatorname{argmax}_a q_{\pi}(s, a)$

$$v_{\pi'}(s) \geq v_{\pi}(s) \quad \forall s$$

π' is an improvement over π .

Suppose new policy π' is only as good as

π , then

$$v_{\pi}(s) = \max_a E (R_{t+1} + \gamma v_{\pi}(s_{t+1}) \mid s_t = s, A_t = a)$$

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

Policy Iteration

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \dots \rightarrow v_{\pi^*}$$

Value Iteration

Turns out, one can be fairly lazy & only partially complete policy iteration steps, & things still work out.

Loop

$$\Delta \leftarrow 0$$

for each $s \in S$

$$v \leftarrow \underline{v(s)}$$

$$v(s) \leftarrow \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - \underline{v(s)}|)$$

until $\Delta \leq \epsilon$

- Model-free policy evaluation
- Model-free control

Model-free policy evaluation

Estimate $V_{\pi}(s)$ using a monte carlo estimate
 $E[G_t | S_t = s]$

First Visit Monte Carlo

Every Visit Monte Carlo

Sample episodes.

within an episode,

the first time you visit a state s

$$\begin{cases} N(s) \leftarrow N(s) + 1 \\ S(s) \leftarrow \underline{S(s)} + \underline{G_t} \end{cases}$$

G_t is return from this point onwards from the first time you visited S_t .

$$\underline{V(s)} \leftarrow S(s) / N(s)$$

Incremental Update

$$\begin{cases} N(s_t) \leftarrow N(s_t) + 1 \\ V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (G_t - V(s_t)) \end{cases}$$

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t - V(s_t)) -$$

↑
Nudge my $V(s_t)$ towards G_t .

Temporal Difference Learning TD TD(0)

$$V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma G_{t+1} - V(s_t))$$

$$\rightarrow V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

$$(s_t, a_t, s_{t+1}, r_{t+1})$$

TD target:

an estimate of my return based on current value function.

TD error.

Model free policy evaluation

TD

- Can start learning from incomplete episodes.
- low variance updates but they can be biased.

$$\underline{V(s_t)} \leftarrow \underline{V(s_t)} + \alpha (\underline{R_{t+1}} + \gamma \underline{V(s_{t+1})} - \underline{V(s_t)})$$

- u

Monte Carlo

high variance estimates.

$$\sum_{i=0}^{\infty} R_{t+i}$$

but it is an unbiased estimate.