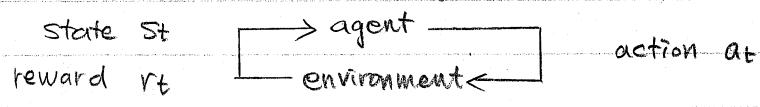


11/16

Review

* Reinforcement learning

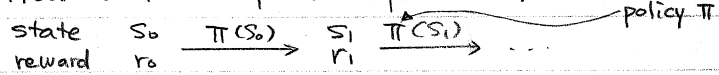


* Markov decision process

$$MDP = \{ \mathcal{S}, \mathcal{A}, P(s'|s, a), R(s) \}$$

states, actions, transitions, rewards.

* How to learn from experience?



* State value function

$$V^\pi(s) = E^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s \right]$$

expected long-term discounted return.

discount factor $0 \leq \gamma < 1$

Recursion relation:

$$\begin{aligned} V^\pi(s) &= E^\pi [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \mid s_0 = s] \\ &= R(s) + \gamma E^\pi [R(s_1) + \gamma R(s_2) + \gamma^2 R(s_3) + \dots \mid s_0 = s] \\ &= R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) E^\pi [R(s_1) + \gamma R(s_2) + \dots \mid s_1 = s'] \end{aligned}$$

$$V^\pi(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s') \quad \text{Bellman equation}$$

* Action-value function

$$\begin{aligned} Q^\pi(s, a) &= \text{expected return from initial state } s, \text{ taking action } a, \\ &\quad \text{then following } \pi \\ &= E^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s, a_0 = a \right] \\ &= R(s) + \gamma \sum_{s'} P(s'|s, a) V^\pi(s') \end{aligned}$$

* Optimality

- optimal policy π^*

Theorem: there is always at least one policy π^* for which $V^{\pi^*}(s) \geq V^\pi(s)$ for all states s and policies π .

Proof: by construction.

- optimal state-value function

$$V^*(s) \triangleq V^{\pi^*}(s)$$

- optimal action-value function.

$$Q^*(s, a) \triangleq Q^{\pi^*}(s, a)$$

- relations

$$V^*(s) = \max_a Q^*(s, a) \quad ; \quad Q^*(s, a) = R(s) + \gamma \sum_{s'} P(s'|s, a) V^*(s')$$

$$\begin{aligned}\pi^*(s) &= \operatorname{argmax}_a Q^*(s, a) \\ &= \operatorname{argmax}_a \sum_{s'} P(s'|s, a) V^*(s')\end{aligned}$$

Planning

Assume complete model of environment $MDP = \{S, A, P(s'|s, a), R(s)\}$ and γ .
How to compute $\pi^*(s)$ or equivalently $V^*(s)$ or $Q^*(s, a)$?

Q1. Policy evaluation

How to compute $V^\pi(s)$ for fixed policy π ?

From Bellman equation:

$$V^\pi(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s') \quad \text{for } s=1, 2, \dots, N$$

This is a system of N linear equations.

$$\begin{aligned}V^\pi(s) - \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s') &= R(s) \\ \sum_{s'} [I(s, s') - \gamma P(s'|s, \pi(s))] V^\pi(s') &= R(s) \\ [I - \gamma P^\pi] V^\pi &= R\end{aligned}$$

known nxn matrix
in terms of MDP
vector of unknowns
vector of rewards

$$V^\pi = (I - \gamma P^\pi)^{-1} R \quad \text{always invertible.}$$

Matrix inversion $O(N^3)$

HW: iterative solution of Bellman equations.

Q2. Policy Improvement

How to compute π' such that $V^{\pi'}(s) \geq V^\pi(s)$ for all s ?

* Define greedy policy $\pi'(s)$ by:

$$\begin{aligned}\pi'(s) &= \operatorname{argmax}_a [Q^\pi(s, a)] \\ &= \operatorname{argmax}_a [R(s) + \gamma \sum_{s'} P(s'|s, a) V^\pi(s')] \\ &= \operatorname{argmax}_a [\sum_{s'} P(s'|s, a) V^\pi(s')]\end{aligned}$$

* Thm: greedy policy π' everywhere performs equally or better than original policy π : $V^{\pi'}(s) \geq V^\pi(s)$ for all s .

Intuition: if better to choose action a in state s , then follow π ,
it's always better to follow action a in state s .

$$\begin{aligned}\text{* Proof: } V^\pi(s) &= Q^\pi(s, \pi(s)) \\ &\leq \max_a Q^\pi(s, a) \\ &= Q^\pi(s, \pi'(s))\end{aligned}$$

$$\text{Expand: } V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')$$

So far, better to take one step under π' and revert to π , than to follow π .

Apply "one-step" inequality to $V^\pi(s')$ on RHS:

$$V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) [R(s') + \gamma \sum_{s''} P(s''|s', \pi'(s')) V^\pi(s'')]$$

Better to take 2 steps under π' and revert to π than to always follow π .

Apply "t" times: better to take t steps under π' and revert to π than to always follow π .

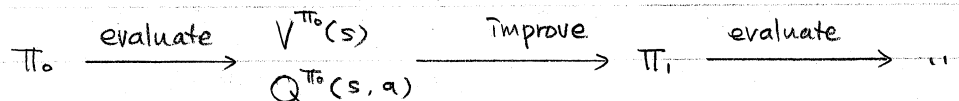
Let $t \rightarrow \infty$: it's always better to follow $\pi'(s)$ than $\pi(s)$.

$$\Rightarrow V^\pi(s) \leq V^{\pi'}(s). \text{ since RHS converges to } V^{\pi'}(s) \text{ assuming } \gamma < 1.$$

Q3. Policy iteration

How to compute π^* ?

- Algorithm:
- 1) initialize policy at random
 - 2) repeat until convergence
 - compute value function
 - derive greedy policy.



* Is this guaranteed to converge? Yes.

- There are finite # policies $|A|^{|S|}$
- Policies cannot be indefinitely improved.
- Typically converge in far less steps than $|A|^{|S|}$.

* Does it always converge to optimal π^* ? Yes!

Thm: suppose $\pi'(s) = \arg\max_a Q^\pi(s, a)$ and $V^{\pi'}(s) = V^\pi(s)$ for all s.

Then, $V^\pi(s) = V^*(s)$.

Note optimal value function $V^*(s)$ is unique even if there are many optimal policies.

* Proof strategy

- 1) Derive "Bellman optimality equation" satisfied by $V^\pi(s)$ when $V^\pi(s) = V^{\pi'}(s)$.
 - 2) Show that $V^\pi(s) \geq V^{\tilde{\pi}}(s)$ for all states and policies $\tilde{\pi}$.
- Hence $V^\pi(s) = V^*(s)$.

* Proof

- 1) From Bellman equation for $\pi'(s)$:

$$V^{\pi'}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi'}(s')$$

By assumption : $V^\pi(s) = V^{\pi'}(s)$ b/c we're converged.

$$\text{Hence : } V^\pi(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')$$

By assumption, π' is greedy with respect to $V^\pi(s)$:

$$\boxed{V^\pi(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^\pi(s')}$$

\Rightarrow Bellman optimality equation. set of nonlinear equation for $s=1, 2, \dots, n$

Different than linear Bellman equation for arbitrary policies.

* Why nonlinear?

$$\max_a f(a) = \lim_{\lambda \rightarrow \infty} \frac{1}{\lambda} \log \sum_a e^{\lambda f(a)}.$$