

Chapter 4: Dynamic Programming

Programming

- DP problems solve finite MDP's of known envs.
- use value functions to organize and structure search of good policies
- turn Bellman equations into assignments
- + provides update rules for improving approximations of desired value functions

Policy Evaluation

- what is the state-value function v_π for an arbitrary π ?
- recall

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

- if env dynamics are known, then you have $|S|$ linear equations in $|S|$ unknowns
- + solution exists, but tedious.
- * start with incremental approach
- * $\{v_0, v_1, v_2, \dots\}$ mappings $S \rightarrow \mathbb{R}$ (take state to real number)
- * pick v_0 at random, and incrementally update

$$v_{k+1}(s) = \mathbb{E}_\pi [R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s]$$

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

- $v_k = v_\pi$ for a given point-in-time guess
- + $\{v_k\} \rightarrow v_\pi$ as $k \rightarrow \infty$
- + iterative policy evaluation

- + update v_{k+1} from v_k by replacing current state value with instantaneous of old state + immediate reward
- * called expected update

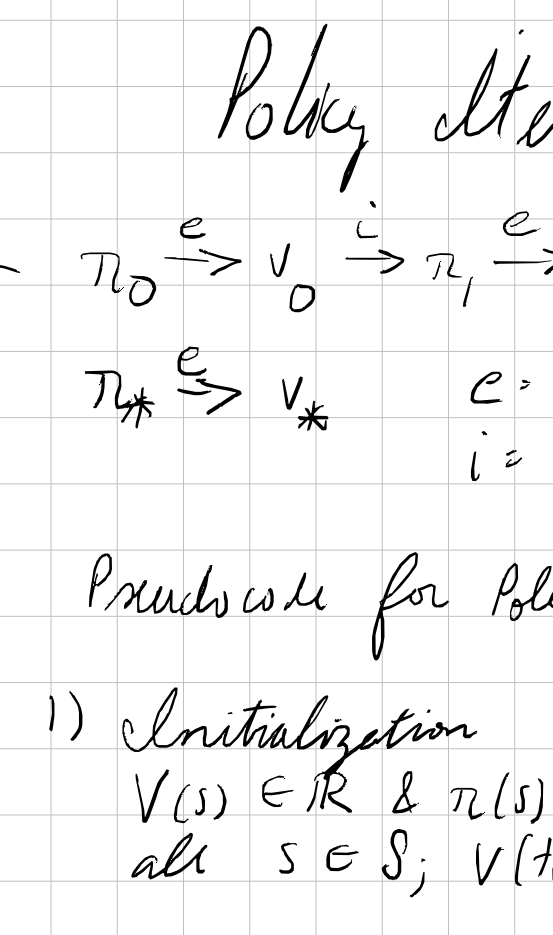
- * called expected because it uses expectation of other states, not a sample of the next one.
- * can use equation or backup diagram to illustrate expected update

- 2 ways to code this
- ① use 2 arrays, one for old values of $v_k(s)$, one for new $v_{k+1}(s)$
- ② use 1 array, update "in-place". (changes based on state-update order)

+ converges faster

+ Sweep through state space

EXAMPLE 4.1: GRIDWORLD 4×4



so $r(s,a,s') = -1 \forall s, s', a$

Refer to figure 4.1 for use of random policy for v_k and greedy policy based off v_k

EXERCISE 4.1 $q_\pi(11, d)$? $q_\pi(7, d)$

$q_\pi(11, d) = -1$

$q_\pi(7, d) = -3.4$ (assuming 3rd iteration)

EXERCISE 4.2: 15 under 13.

$v_\pi(15) = ?$

without modifying original transitions

$$v_\pi(15) = \frac{1}{4} [-1 + \gamma v_\pi(12) - 1 + \gamma v_\pi(13) - 1 + \gamma v_\pi(14) - 1 + \gamma v_\pi(15)]$$

if 13 dynamic change

$$v_\pi(15) = \frac{1}{4} [-1 + \gamma v_\pi(12) - 1 + \gamma v_\pi(9) - 1 + \gamma v_\pi(14) - 1 + \gamma v_\pi(15)]$$

2 equations in $v_\pi(15)$ and $v_\pi(13)$

EXERCISE 4.3: $q_\pi(s, a) = ?$ for 4, 3, 4, 4, 4, 5

$$v_\pi(s) = \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$$

$$q_\pi(s, a) = \mathbb{E} [R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

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

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

now iterative policy evaluator for action-values

$$q_{k+1}(s, a) = \sum_{s',r} p(s',r|s,a) \cdot [r + \gamma \sum_a \pi(a'|s') q_k(s', a')]$$

Policy Improvement

- how do we decide to try a different policy than the current one for a state s ?
- + pick new a , then continue normally
- * use $q_\pi(s, a)$ to see diff.

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

- * adjust one, then ask what is diff for $v_\pi(s)$
- * if better once, one might assume better for all
- policy improvement theorem

+ π, π' are deterministic policies $\forall s \in S \Rightarrow q_\pi(s, \pi'(s)) \geq v_\pi(s)$

+ then π' must be as good if not better than π : $v_{\pi'}(s) \geq v_\pi(s)$

proof: $v_\pi(s) \leq q_\pi(s, \pi'(s))$

$$= \mathbb{E} [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = \pi'(s)]$$

$$= \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$$

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) | S_t = s]$$

$$= \mathbb{E} [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = \pi'(S_{t+1})]$$

$$= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots$$

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

- policy improvement improves on original policy with new policy greedy w.r.t value function
- works just as well for stochastic policies

Policy Iteration

$\pi_0 \xrightarrow{e} v_0 \xrightarrow{i} \pi_1 \xrightarrow{e} v_1 \xrightarrow{i} \pi_2 \dots$

$\pi_k \xrightarrow{e} v_k$ $e = \text{eval}$
 $i = \text{improve}$

Pseudocode for Policy Iteration for $\pi \approx \pi^*$

1) Initialization
 $V(s) \in \mathbb{R}$ & $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in S$; $V(\text{terminal}) = 0$

2) Policy Evaluation
Loop while $\Delta < \theta$
 $\Delta = 0$
Loop for $s \in S$
 $v = V(s)$
 $V(s) = \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$
 $\Delta = \max(\Delta, |v - V(s)|)$

3) Policy Improvement
policy-stable = true
for each $s \in S$:
old-action = $\pi(s)$
 $\pi(s) = \text{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$
if old-action $\neq \pi(s)$
policy-stable = false
if policy-stable = true
return $v, v_k, \pi \approx \pi_k$
else
goto policy eval

Refer to figure 4.2 for Jack's Car Example

- policy iteration can occur in very few iterations

EXERCISE 4.4: Modify pseudocode of policy iteration to remove subtle bug

if policy of v_k & v_{k+1} is within some small delta, stop. π and π' may differ, but if they are both optimal, v_π and $v_{\pi'}$ will converge on v_*

EXERCISE 4.5: Write pseudocode for action-value policy iteration

$\Delta = \infty$ eval
while $\Delta > \epsilon$ (some threshold)
for $\forall s \in S$
for $\forall a \in \mathcal{A}(s)$
 $q = Q(s, a)$

mistake: need $\sum_{s',r} p(s',r|s,a) \left(Q(s, a) - r - \gamma \sum_{s'} \pi(s') Q(s', a') \right)$

just using $\Delta = \max |q - Q(s, a)|$

policy-stable = true improve
for each $s \in S$
 $a_{\text{old}} = \pi(s)$
 $\pi(s) = \text{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma \sum_{s'} \pi(s') Q(s', a')]$

if $a_{\text{old}} \neq \pi(s)$ mistake: just
policy-stable = false argmax $Q(s, a)$
if policy-stable = false you have given the def of $Q(s, a)$
return $Q^* q^*, \pi \approx \pi^*$ This is still right but overcomplicated. You already have the function you want

$v_*(s) = \text{argmax}_a q(s, a)$

EXERCISE 4.6 ϵ -soft $\epsilon/|A(s)|$ at least for action probability

1) initialization: $\pi(s) = \epsilon/|A(s)| \forall s$

2 & 3) $\pi(a|s)$ must be either $1 - \epsilon + \epsilon$ for greedy or $\epsilon/|A(s)|$

Termination: since v_* for Δ may not converge, must converge around $\pi(a|s)$

Value Iteration

- policy iteration requires protracted iterative computation of state space
- + convergence to v_* occurs in $\lim \rightarrow \infty$
- * too to stop sooner?
- special case: stop evaluation after single sweep of state space
- + Value Iteration combine policy improvement with single eval loop.

$$v_{k+1}(s) = \max_a \mathbb{E} [R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s, A_t = a]$$

$$= \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')]$$

- naturally occurring from Bellman if viewed as an update rule

- algo:

init: $\theta > 0$ threshold
 $v(s) \forall s \in S, v(\text{terminal}) = 0$

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

- faster convergence by interleaving multiple eval & improve sweeps.

Asynchronous Dynamic Programming

- single sweep, in order, LONGER TIME
- asynchronous, don't care for order
- + must still update all state if it hasn't, even if done out of order
- allows for smarter, more fine-tuned state update ordering
- can make real-time updates without needing all other state
- + experience state informs current interpretation of reality

Generalized Policy Iteration

- free and firm interleaving of eval and improve
- generalized policy iteration \Rightarrow fixed
- + fundamental for RL $\pi \xrightarrow{\text{improve}} v \xrightarrow{\text{eval}} \pi$
- + stabilize own time $\pi \xrightarrow{\text{improve}} v$
- + competing and cooperating at the same time $\pi \xleftrightarrow{\text{improve}} v$

Efficiency of DP

- worst case - polynomial in state + actions
- + guaranteed to find optimal policy in polynomial time even if total deterministic policies are k^n (n states k actions)
- curse of dimensionality affects the large state problem, not DP itself
- + DP handles quite well

Summary

- Final note: extracting based on estimate = bootstrapping
- + DP requires model of world + bootstraps
- * some RL use model & bootstrap. Some only model since only bootstrap.

EXERCISE 4.10: $g_{k+1}(s, a)$?

$$g_{k+1}(s, a) = \sum_{s',r} p(s',r|s,a) [r + \gamma \max_{a'} g_k(s')]$$

EXERCISE 4.8: There's no reward to oversteering. Rather than risk more for no reward, bet small parcel amounts till a full double of bet = 100