

Chapter 5: Monte Carlo Methods

06/08/19

Monte-Carlo methods :

- require experience only
- no assumption of complete knowledge of env.
- based on averaging sample returns
- * updates occur only on episode completion to have well defined returns
- random sampling

5.1 Monte Carlo Prediction

prediction of state value function, $v_{\pi}(s)$

- ↳ first visit MC (average of all returns following first)
- ↳ every visit MC (average of all returns)

Both of these converge to $v_{\pi}(s)$ as number of visits $\rightarrow \infty$

Blackjack

Objective: maximize sum constraint to sum ≤ 21

Rules : Face card = 10, ace = 1 or 11, number cards = number

2 cards dealt to both dealer and player, one of dealer's card is face up

Scoring : natural \leftarrow has 21 at beginning

Actions : hit, stick, going bust
↑ changes turns

Win : sum closest to 21 unless someone goes bust

Exercise 5.1:

- value function jumps as probability of winning from those states is very high
- dip in left now as getting usable aces is hard and unusable aces merely add 1 to the sum
- higher sum with usable aces and hence a greater chance of winning

Exercise 5.2:

- results should be same as both converge to $q_{\pi}(s)$

Note: Estimates for each state are independent

Compute required to get value of a single state is independent of the number of states

5.2 MC Estimation of Action Values

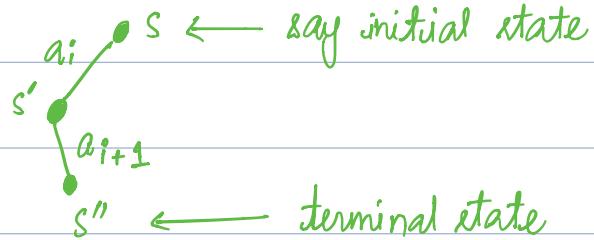
model-free \leftarrow action value is used

task is to estimate $q_{\pi}(s, a)$

two types $\begin{cases} \rightarrow \text{every visit MC} \\ \rightarrow \text{first visit MC} \end{cases}$

exploring starts: every $s-a$ pair has non-zero probability of being selected at the start

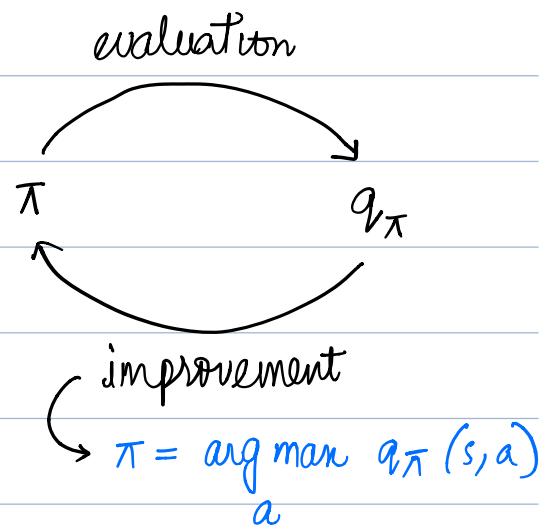
Exercise 5.3 : Backup diagram for MC estimation of q_{π}



5.3 : MC Control

GPI (generalized policy iteration)

↳ has both π and q_{π}



assumptions for convergence :

- have exploring starts
- infinite number of episodes

↳ two ways to overcome this :

1. measurements, approximations are made to obtain bounds on mag.

and prob. of error and accordingly sufficient steps are taken

- this approach however requires too many steps for non-trivial problems

2. in MC policy iteration eval. and improvement is alternated

on an episode-by-episode basis. Policy is improved at all states visited during the episode after the end of the episode

Exercise 5.4:

$$R^t(s) \leftarrow R^t(s) + \frac{1}{t+1} [R - R^t(s)]$$

5.4 : MC Control w/o Exploring Starts

- overcoming the first assumption made for convergence
- on-policy v/s off-policy
 - ↳ evaluate/improve policy that is used to make decisions

eq: MC ES

- for on-policy methods, policy is usually *soft*
 - ↳ $\pi(a|s) > 0 \quad \forall a, s$

$$\left\{ \begin{array}{l} A^* \leftarrow \arg \max_a q_\pi(s, a) \\ \pi(a|s) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |A(s)| & \text{if } a = A^* \\ \varepsilon / |A(s)| & \text{otherwise} \end{cases} \end{array} \right\} \quad \forall a \in A(s)$$

policy improvement theorem: assures that any ε -greedy policy wrt. q_π is an improvement over ε -soft policy π

$$\begin{aligned} q_{\pi'}(s, \pi'(s)) &= \sum_a \pi'(a|s) q_\pi(s|a) \\ &= \frac{\varepsilon}{|A(s)|} \sum_a q_\pi(s|a) + (1-\varepsilon) \max_a q_\pi(s|a) \\ &\geq \frac{\varepsilon}{|A(s)|} \sum_a q_\pi(s|a) + (1-\varepsilon) \sum_a \frac{\pi(a|s) - \varepsilon / |A(s)|}{1-\varepsilon} q_\pi(s|a) \\ &= \sum_a \pi(a|s) q_\pi(s|a) = v_{\pi'}(s) \end{aligned}$$

Thus, $\pi' \succ \pi$ (ie $v_{\pi'}(s) \geq v_{\pi}(s) \quad \forall s \in S$)

5.5 : Off-policy Prediction via Importance Sampling

dilemma : action values are learned conditional to subsequent optimal behaviour

↳ this restricts exploration

solution : have 2 policies, target and behaviour

"off-policy" methods

↑
optimize

↑ facilitate
exploration

- learning happens from data "off" target policy

off policy methods :

- slow convergence
- high variance
- more powerful, generally applicable

prediction problem : both target and behaviour policies are fixed
ie, estimate v_π, a_π given b (behaviour policy)

assumption of coverage :

$\pi(a|s) > 0$ implies that $b(a|s) > 0$

π may be deterministic

b is stochastic for states where it is non-identical to π

importance sampling:

- general method of estimating expected values under one distib. given samples from another
- importance sampling ratio: returns are weighed according to the relative trajectories occurring under the off-policy

For initial state s_t , the probability of subsequent state action trajectory $A_t, s_{t+1}, A_{t+1}, \dots, s_T$ under policy π is:

$$\begin{aligned} & \Pr \{ A_t, s_{t+1}, A_{t+1}, \dots, s_T \mid s_t ; A_{t:T-1} \sim \pi \} \\ &= \pi(A_t | s_t) p(s_{t+1} | s_t, A_t) \pi(A_{t+1} | s_{t+1}) \dots p(s_T | s_{T-1}, A_{T-1}) \\ &= \prod_{k=t}^{T-1} \pi(A_k | s_k) p(s_{k+1} | s_k, A_k) \end{aligned}$$

$$\begin{aligned} \therefore P_{t:T-1} &= \frac{\Pr \{ \dots ; A_{t:T-1} \sim \pi \}}{\Pr \{ \dots ; A_{t:T-1} \sim b \}} \\ &= \prod_{k=t}^{T-1} \frac{\pi(A_k | s_k)}{b(A_k | s_k)} \end{aligned}$$

(This is importance sampling ratio)

now, $E[G_t | s_t] = v_b(s)$ is transformed into:

$$v_\pi(s) = E[P_{t:T-1} G_t | s_t]$$

Let $\Gamma(s) \leftarrow$ set of all time steps when s is visited

$T(t) \leftarrow$ first time of termination following t

$G_t \leftarrow$ return following t till $T(t)$

\therefore under ordinary importance sampling,

$$V(s) = \frac{\sum_{t \in \Gamma(s)} p_{t:T(t)-1} G_t}{|\Gamma(s)|}$$

\therefore under weighted importance sampling,

$$V(s) = \frac{\sum_{t \in \Gamma(s)} p_t : T-1 G_t}{\sum_{t \in \Gamma(s)} p_t : T-1}$$

- ordinary importance sampling is unbiased whereas weighted importance sampling is biased.
- bias converges to zero asymptotically
- variance of OIS is unbounded, for WIS it is bounded to 1

Exercise 5.5: (?)

$$G_t = \{9, 8, 7, 6, 5, 4, 3, 2, 1, 0\}$$

$$\Pr\{\text{trajectory}\} = p^k (1-p)^{10-k}$$

$\Gamma(s) \leftarrow$ differs for every visit and first visit

$$|\Gamma(s)| = 10$$

$$|\Gamma(s)| = 8 \text{ or } 9$$

