

687 2017-10-19

Note Title

10/19/2017

First Visit MC:

If each state s is visited infinitely often then $\hat{v}_k \xrightarrow{\text{a.s.}} v^\pi$.

If each state-action pair is visited infinitely often, then $\hat{q}_k \xrightarrow{\text{a.s.}} q^\pi$.

Doing this: MC Estimation of Action-values:

$\hat{q}(s, a)$ = average return from the first time action a taken in state s in each episode.

Problem: What if π never chooses a in state s ?

Is $\hat{q}(s, a)$ still defined? YES
 $\hat{q}(s, a)$ will not be estimated!

What to do?

Solution 1: Exploring starts:

Randomize s_0 and a_0 s.t. every (s, a) pair has a non-zero probability.
→ Not possible for all systems!

Solution 2: Stochastic Policy π

Non-zero probability for every action in each state.

Monte Carlo Control with Exploring Starts

- Avoid generating ^{infinite #} of episodes in the evaluation step by improving after a single episode.
- Accumulate returns over all episodes.

Pseudocode:

Init: For all $s \in S, a \in A$,
 $q(s, a) \leftarrow$ arbitrary value
 $\pi(s) \leftarrow$ "action"
 $\text{Returns}(s, a) \leftarrow$ empty list

Repeat Forever

- ① Generate an episode using exploring starts using π
- ② For each (s, a) in the episode,
 $G \leftarrow$ return following first occurrence
of (s, a) in the episode
Append G to $\text{Returns}(s, a)$
 $q(s, a) \leftarrow \text{mean}(\text{Returns}(s, a))$
- ③ For each s in the episode, set
 $\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} q(s, a)$

If this converges, it does so to an optimal π .

Proof: Suppose it converges to a suboptimal π . Then q has converged to the actual values for π . Then policy improvement must move π to a better policy.

Note: It can fail to converge.

Monte Carlo Control with Stochastic Policies:

ϵ -greedy MC Control:

Init: $q(s,a) \leftarrow \text{arb. values}$

$\pi(s,a) \leftarrow \text{arb. s.t. } \pi(s,a)$

Returns(s,a) \leftarrow empty

Repeat Forever:

① Generate an episode using π

② For each (s,a) in the episode:

$G \leftarrow$ Return following first occurrence
of (s,a) in the episode

Append G to Returns(s,a)

$q(s,a) \leftarrow \text{mean}(\text{Returns}(s,a))$

③ For each s in the episode

$a^* \in \arg \max_a q(s,a)$

$$\pi(s,a) \leftarrow \begin{cases} 1-\epsilon + \frac{\epsilon}{|a|} & \text{if } a = a^* \\ \frac{\epsilon}{|a|} & \text{if } a \neq a^* \end{cases}$$

Properties:

- By variations on the Policy Improvement Thm (Sutton and Barto, Vol. I, Sect. 5.4), converges to an ϵ -greedy policy.
- In fact to the optimal ϵ -greedy policy (not necessarily optimal among all policies).

TD (Temporal Difference) Learning:

Properties:

- Like MC, learns directly from experiences: does not need P and R .
- Like DP (dynamic programming): update estimates in terms of other estimates.

But first, another MC algorithm:

Consider updates: $f(x) \leftarrow f(x) + \alpha (Y - f(x))$

"step size"

target (usually an estimate)

$$f_w(x) \quad L(w) = \mathbb{E} \left[\frac{1}{2} (Y - f_w(x))^2 \right]$$

↑ weights

$$w \leftarrow w + \alpha \mathbb{E} \left[(Y - f_w(x)) \cdot \frac{\partial f_w(x)}{\partial w} \right]$$

↓ gradient of $L(w)$

gradient descent

$$w \leftarrow w + \alpha (Y - f_w(x)) \frac{\partial f_w(x)}{\partial w} = \begin{bmatrix} 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{bmatrix} \quad \begin{matrix} \text{I only in the} \\ x^{\text{th}} \text{ spot} \end{matrix}$$

estimate from one sample:
stochastic gradient descent

$$\nabla v(s_t) \leftarrow v(s_t) + \alpha \sum_{K=0}^{\infty} \gamma^K R_{t+K}$$

stochastic gradient descent on
 $L(v) = \mathbb{E} [(G_t - v(s_t))^2]$

But have to wait until the end
of the episode!

use the Bellman Equation:

$$v'(s) = \mathbb{E} [R_t + \gamma v''(S_{t+1}) | S_t = s]$$

$$v(s_t) \leftarrow v(s_t) + \alpha (R_t + \gamma v(S_{t+1}) - v(s_t))$$

Called the TD Update

$$\text{so } f(x) \leftarrow f(x) + \alpha (Y - f(x))$$

The TD Error:
 δ_t

$$v(s) \leftarrow v(s) + \alpha \delta_t$$

"Reward Prediction Error"

MC is a stochastic gradient descent.

Is TD stochastic gradient descent?

If it is: $f(v) = \mathbb{E} \left[\frac{1}{2} (R_t + \gamma v(S_{t+1}) - v(S_t))^2 \right]$

$$= \mathbb{E} \left[\frac{1}{2} \delta_t^2 \right]$$

$$\nabla f(v) = \delta_t \cdot \frac{\partial \delta_t}{\partial v} \quad (\text{ignoring expectation})$$

$$= (R_t + \gamma v(S_{t+1}) - v(S_t)) \left(\gamma \frac{\partial v(S_{t+1})}{\partial v} - \frac{\partial v(S_t)}{\partial v} \right)$$

\downarrow

$$S_{t+1} \rightarrow \begin{bmatrix} 0 \\ i \\ 0 \end{bmatrix} \qquad \begin{bmatrix} 0 \\ j \\ 0 \end{bmatrix} \leftarrow S_t$$

$$v(S_t) \leftarrow v(S_t) + \alpha \delta_t$$

$$v(S_{t+1}) \leftarrow v(S_{t+1}) - \alpha \gamma \delta_t$$

Seems close to s.g.d., but it is not properly sampled. (Will cover later in the course.)

Some properties:

- TD offers an update at each step.
- Converges a.s. to correct v .