

687 2017-10-31

Note Title

10/31/2017

Sarsa (s, a, r, s', a')

$$\left. \begin{aligned} \delta_t &\leftarrow R_t + \gamma q_f(s_{t+1}, A_{t+1}) - q_f(s_t, A_t) \\ q_f(s_t, A_t) &\leftarrow q_f(s_t, A_t) + \alpha \delta_t \end{aligned} \right\} \text{on-policy}$$

converges to $\rightarrow q^{\pi_b}$ under sampling policy π_b

Q-learning (s, a, r, s')

$$\left. \begin{aligned} \delta_t &\leftarrow R_t + \gamma \max_{a'} q_f(s_{t+1}, a') - q_f(s_t, A_t) \\ q_f(s_t, A_t) &\leftarrow q_f(s_t, A_t) + \alpha \delta_t \end{aligned} \right\} \text{off-policy} \longrightarrow q^*$$

Init $q_f(s, a)$ arbitrarily

Repeat (for each episode)

Init $s \sim d_0$

Repeat (for each timestep)

- choose a from s using the policy derived from q
- take action a , observe r, s'
- update $q_f(s, a) \leftarrow q_f(s, a) + \alpha (r + \gamma \max_{a'} q_f(s', a') - q_f(s, a))$

Until $\underline{s \leftarrow s'}$ s is terminal

Course Recap

① Definitions

MDP

S, A, P, R, d_0, γ

v^*, q^*, v^π, q^π

J, π

② Dynamic Programming

Policy evaluation
+ Bellman Eqn.

Policy Improvement \Rightarrow

Policy Iteration

Value Iteration

- Bellman optimality eqn.

- Bellman operator

③ Monte Carlo approximation

⑥ Policy Gradient (∇)

Theorems

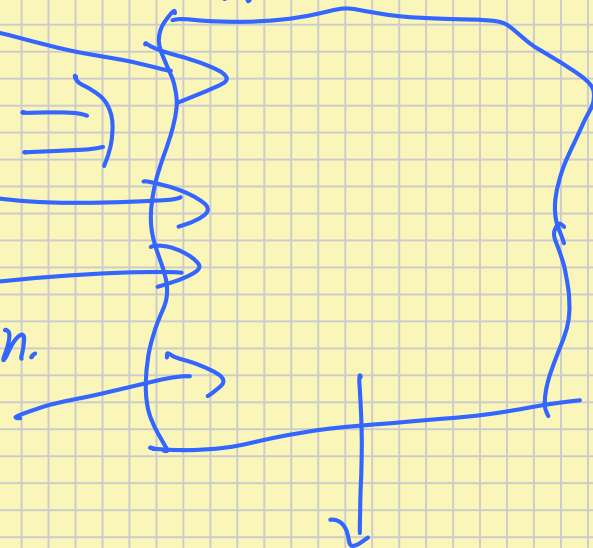
Midterm

⑤ TD(N)

Same

④ TD

Same
+ func. Approx.



Convergence

TD

Sarsa

Q-learning

Tabular

✓

✓

✓

Linear

✓

✓

✗

General

✗

✗

✗

A-Step Returns

- Different potential targets for value function updates.

$G_t^{(n)}$ is the n-step return

$$G_t^{(1)} = R_t + \gamma v(S_{t+1}) = G_t^{(TD)}$$

$$G_t^{(2)} = R_t + \gamma R_{t+1} + \gamma^2 v(S_{t+2})$$

Larger n gives lower bias, higher variance.

$$G_t^{(n)} = \sum_{k=0}^{n-1} \gamma^k R_{t+k} + \gamma^n v(S_{t+n})$$

$$G_t^{(2)} = \text{Monte Carlo} = G_t \text{ or } G_t^{(MC)}$$

Idea: use more than one (in fact, of all); Use a weighted average, where weights sum to 1.

λ -Returns

$$G_t^\lambda \triangleq \left(\lambda^0 G_t^{(1)} + \lambda^1 G_t^{(2)} + \dots + \lambda^i G_t^{(i+1)} + \dots \right) \cdot (1-\lambda)$$

where $\lambda \in (0, 1)$

$$= (1-\lambda) \sum_{i=1}^{\infty} \lambda^{i-1} G_t^{(i)} \text{ if } \lambda < 1$$

and $G_t^{(MC)}$ otherwise

This approach is intuitively appealing to reduce variance - but not why we do it.

Andy Barto says "The math works out." In particular, it works backwards for actual computation.

TD γ ... Complex ... @ NIPS 2011
 Ω -return @ " 2015
ICML 2016

λ -return is MLE-estimate of $v^\pi(S_t)$ if:

- ① $G_t^{(1)}, G_t^{(2)}, \dots$ are independent
- ② Every $G_t^{(i)}$ is normally distributed
- ③ $\text{Var}(G_t^{(i)}) \propto \frac{1}{\lambda^i}$ ($\sigma^2 = \beta/\lambda^i$)
- ④ $\mathbb{E}[G_t^{(i)}] = v^\pi(S_t)$ for all i

These assumptions are false, but the rest of the math works out!

To find best, $\frac{d}{dx}$ + set to 0

$$\sum_{i=1}^{\infty} 2(G_t^{(i)} - x)\lambda^i = 0$$

$$\sum_{i=1}^{\infty} G_t^{(i)} \lambda^i = x \sum_{i=1}^{\infty} \lambda^i \Rightarrow x = \frac{\left(\sum_{i=1}^{\infty} G_t^{(i)} \lambda^i\right) (1-\lambda)}{\lambda} = \left(\sum_{i=1}^{\infty} G_t^{(i)} \lambda^{i-1}\right) (1-\lambda)$$

The λ -return.

Proof:

$$\begin{aligned} L(x | G_t^{(1)}, G_t^{(2)}, \dots) \\ = \Pr(G_t^{(1)}, G_t^{(2)}, \dots | v^\pi(S_t) = x) \\ = \prod_{i=1}^{\infty} \Pr(G_t^{(i)} | v^\pi(S_t) = x) \end{aligned}$$

$$\arg \max_x L(x) = \arg \max_x \ln L(x)$$

$$\downarrow \arg \max_x L(x) = \arg \max_x \sum_{i=1}^{\infty} \ln \Pr(G_t^{(i)} | v^\pi(S_t) = x)$$

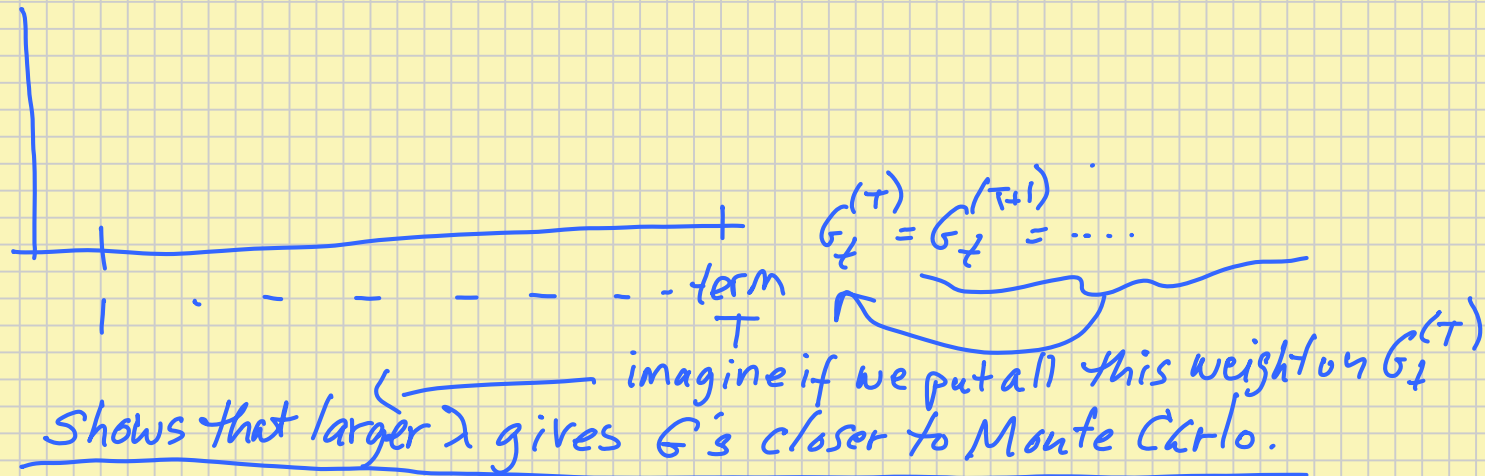
$$= \arg \max_x \sum_{i=1}^{\infty} \ln \left(\frac{1}{\sqrt{2\pi} \beta/\lambda^i} e^{-\frac{(G_t^{(i)} - x)^2}{2\beta/\lambda^i}} \right)$$

$$= \arg \max_x \sum_{i=1}^{\infty} \left[\ln \left(\frac{1}{\sqrt{2\pi} \beta/\lambda^i} \right) - \frac{(G_t^{(i)} - x)^2 \lambda^i}{2\beta} \right]$$

$$= \arg \max_x - \sum_{i=1}^{\infty} (G_t^{(i)} - x)^2 \lambda^i$$

$$= \left(\sum_{i=1}^{\infty} G_t^{(i)} \lambda^{i-1}\right) (1-\lambda)$$

What is the λ -return at a terminal state?



λ -return Algorithm

Use $G_t^{(A)}$ as target. $v(S_t) \leftarrow v(S_t) + \alpha (G_t^{(A)} - v(S_t))$ ← "Forward view," which requires waiting until the end of the episode.