

LINMA2725 Stochastic Optimal Control and Reinforcement Learning

Part III

Course 3: Actor-Critic Methods

Guillaume Berger

November 25, 2024

Reference: [1], Chapter 10.

Any questions or feedback are welcome.

Stochastic Optimal Control and Reinforcement Learning

Part III: Stochastic Systems

Guillaume Berger

Table of content

- Stochastic systems and stochastic control (1 course)
- Learning techniques for stochastic control (1-2 courses)
 - Variance reduction techniques
 - Actor-critic methods
- Online learning techniques for stochastic control (1-2 courses)

Problem setting and reminders

Consider a given (randomized) policy $\check{\phi}(u|x)$

Let ϖ be the steady-state distribution of the closed-loop system

Let $Q(x, u)$ be its Q -function and $h(x)$ its value function

Cost can be total, discounted or average cost

Objective: approximate Q or similar quantities

Assume: $Q^\theta(x, u) = \theta^\top \psi(x, u)$ (linear template)

For theoretical analysis: assume X and U finite

Advantage function

Advantage function: $V(x, u) = Q(x, u) - h(x)$

Motivation: $h = \arg \min_{g: X \rightarrow \mathbb{R}} \mathbb{E}_{\omega} \left[\{Q(\Phi(n)) - g(X(n))\}^2 \right]$

Hence, L_2 -norm of V is smaller than Q

But optimal policy of V and Q are the same!

Note that $\mathbb{E}_{\omega} [V(\Phi(n)) \mid X(n)] = 0$

The idea is that if we “offset” $Q(x, u)$ with a function $g(x)$, then this does not change the associated optimal policy, i.e.,

$$\arg \min_{u \in \mathcal{U}} Q(x, u) = \arg \min_{u \in \mathcal{U}} \{Q(x, u) - g(x)\}.$$

We choose the offset $g(x)$ that minimizes the L^2 -norm of $Q(x, u) - g(x)$, which is a sensible idea since having quantities with a smaller magnitude (norm) involved in the recursive algorithm will generally induce a smaller variance of the algorithm. The optimal $g(x)$ is given by $h(x) \triangleq \sum_{u \in \mathcal{U}} \phi(u|x)$.

Advantage TD(λ)-learning

TD(λ) algorithm (on-policy for advantage)

For initialization $\omega_0, \zeta_0 \in \mathbb{R}^m$, the sequence of estimates are defined recursively:

$$\begin{aligned}\omega_{n+1} &= \omega_n + \alpha_{n+1} \zeta_n \mathcal{D}_{n+1} \\ \mathcal{D}_{n+1} &= (-H^\omega(\Phi(n)) + c_n + \gamma H^\omega(\Phi(n+1))) \Big|_{\omega=\omega_n} \\ \zeta_{n+1} &= \lambda \gamma \zeta_n + \tilde{\psi}^\circ(\Phi(n+1)), \quad H^\omega(\Phi(n)) = \omega^\top \tilde{\psi}^\circ(\Phi(n))\end{aligned}\tag{9.68}$$

where $\underline{\psi}(x) = \sum_u \psi(x, u) \check{\phi}(u | x)$ and $\underline{\tilde{\psi}}(x, u) = \psi(x, u) - \underline{\psi}(x)$,

Same as TD(λ)-learning but with a special basis!

Motivation of advantage TD(λ)-learning

Remember $\lambda = 1$: θ^* solves

$$\theta^* = \arg \min_{\theta} \|H^\theta - Q\|_{\omega}^2 \stackrel{\text{def}}{=} \sum_{x \in X, u \in U} (H^\theta(x, u) - Q(x, u))^2 \omega(x, u)$$

New result: the projection of V on $\{\theta^\top \psi + \xi^\top \underline{\psi}\}$ is equal to the projection of Q on $\{\theta^\top \tilde{\psi}\}$

Hence, H^{θ^*} is the projection of V on $\{\theta^\top \psi + \xi^\top \underline{\psi}\}$

Proof of the result. Let \hat{Q} be the projection of Q on $\{\theta^\top \tilde{\psi} : \theta \in \mathbb{R}^d\}$. We show that $V - \hat{Q}$ is orthogonal to ψ and $\underline{\psi}$. This will imply that \hat{Q} is the projection of V on $\{\theta^\top \psi + \xi^\top \underline{\psi} : \theta, \xi \in \mathbb{R}^d\}$. We show the orthogonality only for ψ since the proof for $\underline{\psi}$ is easier. For that, observe that

$$\begin{aligned} \mathbf{E}_\varpi \left[(V - \hat{Q})\psi \right] &= \mathbf{E}_\varpi \left[(Q - h - \hat{Q})(\tilde{\psi} + \underline{\psi}) \right] \\ &= \mathbf{E}_\varpi \left[(Q - \hat{Q})\tilde{\psi} \right] + \mathbf{E}_\varpi \left[Q\underline{\psi} \right] - \mathbf{E}_\varpi \left[\hat{Q}\underline{\psi} \right] - \mathbf{E}_\varpi \left[h\tilde{\psi} \right] - \mathbf{E}_\varpi \left[h\underline{\psi} \right] \end{aligned}$$

the first term is zero by definition of \hat{Q} ,

the third and fourth terms are both zero because $\mathbf{E}_\varpi \left[\tilde{\psi}(\Phi(n)) \mid X(n) \right] = 0$

the second is equal to $\mathbf{E}_\varpi \left[h\underline{\psi} \right]$ since $\mathbf{E}_\varpi \left[Q(\Phi(n)) \mid X(n) \right] = h(X(n))$

$$= \mathbf{E}_\varpi \left[h\underline{\psi} \right] - \mathbf{E}_\varpi \left[h\underline{\psi} \right] = 0.$$

This concludes the proof. □

Advantage TD(λ)-learning as a LS system

LSTD(1) for advantage with weighting (on-policy)

With weighting function $w: \mathcal{X} \rightarrow (0, \infty)$, initialization $\zeta_0 \in \mathbb{R}^m$, $\hat{\Sigma}_0 \in \mathbb{R}^{m \times m}$, and time horizon N ,

$$\omega_N = \hat{\Sigma}_N^{-1} \bar{\psi}_N^Q \quad (10.16a)$$

$$\text{with } \hat{\Sigma}_N = \frac{1}{N} \left(\hat{\Sigma}_0 + \sum_{n=1}^N w_n \tilde{\psi}_{(n)} \tilde{\psi}_{(n)}^\top \right) \quad (10.16b)$$

$$\bar{\psi}_N^Q = \frac{1}{N} \sum_{n=1}^N c_n \zeta_n \quad (10.16c)$$

$$\zeta_n = \gamma \zeta_{n-1} + w_n \tilde{\psi}_{(n)}, \quad w_n = w(X(n)), \quad \tilde{\psi}_{(n)} = \tilde{\psi}^\circ(\Phi(n)), \quad 1 \leq n \leq N \quad (10.16d)$$

Regeneration

Let $z^\bullet = (x^\bullet, u^\bullet) \in X \times U$ denote any state with positive steady-state probability: $\omega(z^\bullet) > 0$.

$$\tilde{Q}_\gamma(z) = \mathbb{E}\left[\sum_{k=0}^{\tau_\bullet-1} \gamma^k \tilde{c}(\Phi(k)) \mid \Phi(0) = z\right] + \mathbb{E}\left[\gamma^{\tau_\bullet} \sum_{k=0}^{\infty} \gamma^k \tilde{c}(\Phi(k + \tau_\bullet)) \mid \Phi(0) = z\right]$$

$$\tau_\bullet = \min\{k \geq 1 \mid \Phi(k) = z^\bullet\}$$

Result:
$$\tilde{Q}_\gamma(z) = \underbrace{\mathbb{E}\left[\sum_{k=0}^{\tau_\bullet-1} \gamma^k \tilde{c}(\Phi(k)) \mid \Phi(0) = z\right]}_{\text{focus on this to reduce variance!}} + \underbrace{\tilde{Q}_\gamma(z^\bullet) \mathbb{E}[\gamma^{\tau_\bullet} \mid \Phi(0) = z]}_{\text{tends to constant when } \gamma \rightarrow 1}$$

focus on this to reduce variance!

tends to constant when $\gamma \rightarrow 1$

Focusing on the first term (orange) is another way of removing an offset without impacting the associated optimal policy because the second term (blue) is assumed to be constant (independent of z).

When $\gamma = 1$, regeneration is also crucial to have boundedness of the eligibility vectors ζ_n in the TD(1)-learning algorithm, as we will see next.

Regeneration for average cost

Define $H_3(z) = \mathbb{E} \left[\sum_{k=0}^{\tau_{\bullet}-1} \tilde{c}(\Phi(k)) \mid \Phi(0) = z \right]$

Goal: approximate H_3 with H^θ

Motivations:

- Average cost is important for stochastic systems!
- Other costs like mean discounted cost can be formulated as average cost

The average cost can be used as a metric to compare policies between them because it is a scalar. By contrast, the discounted cost $h_\gamma(x)$ cannot be used as a metric because it is a function. Alternatively, one can use the “mean discounted cost”: given a probability distribution μ on \mathbf{X} , define the metric $\langle \mu, h_\gamma \rangle \triangleq \sum_{x \in \mathbf{X}} \mu(x) h_\gamma(x)$.

The mean discounted cost of a given system can be formulated as the average cost of another system constructed from the original one; see [1, § 10.4.1] for such a construction.

Regenerative TD(λ)-learning for average cost

Let $z^\bullet = (x^\bullet, u^\bullet) \in X \times U : \omega(z^\bullet) > 0$.

Regenerative TD(λ) algorithm for average cost (on-policy)

For initialization $\theta_0, \zeta_0 \in \mathbb{R}^d$, the sequence of estimates are defined recursively:

$$\begin{aligned}\theta_{n+1} &= \theta_n + \alpha_{n+1} \zeta_n \mathcal{D}_{n+1} \\ \mathcal{D}_{n+1} &= \left(-H^\theta(\Phi(n)) + \tilde{c}_n + \mathbb{1}\{\Phi(n+1) \neq z^\bullet\} H^\theta(\Phi(n+1)) \right) \Big|_{\theta=\theta_n} \\ \zeta_{n+1} &= \lambda \mathbb{1}\{\Phi(n+1) \neq z^\bullet\} \zeta_n + \psi_{(n+1)} \\ \eta_{n+1} &= \eta_n + \tilde{c}_n / (n+1), \quad \tilde{c}_n = c(\Phi(n)) - \eta_n \quad n \geq 0\end{aligned} \tag{10.28}$$

Analysis of regenerative TD(λ)-learning

This is a linear SA algorithm

$$\begin{aligned} A_{n+1} &= \zeta_n [-\psi_{(n)} + \mathbb{1}\{\Phi(n+1) \neq z^\bullet\} \psi_{(n+1)}]^\top \\ b_{n+1} &= -\zeta_n [c_n - \eta_n] \end{aligned}$$

The associated ODE is thus linear, with vector field

$$\bar{f}(\vartheta) = A\vartheta + b, \quad A = \mathbb{E}_\omega[A_n], \quad b = -\mathbb{E}_\omega[\zeta_n \tilde{c}(\Phi(n))]$$

Theorem 10.14. (*L_2 Optimality of TD(1) for average cost*) Consider the algorithm (10.28) with linear function approximation $H^\theta = \theta^\top \psi$. Assume that Φ is uni-chain and that $\omega(z^\bullet) > 0$.

Then, in the special case $\lambda = 1$,

(i) $A = -R(0)$

(ii) Any solution to $0 = \bar{f}(\theta^*) = \mathbb{E}_\omega[\zeta_n \mathcal{D}_{n+1}]$ solves the minimum norm problem:

$$\theta^* \in \arg \min_{\theta} \|H^\theta - H_3\|_\omega^2 = \arg \min_{\theta} \mathbb{E}_\omega [(H^\theta(\Phi(n)) - H_3(\Phi(n)))^2] \quad (10.29)$$

See [1, Theorem 10.14]. Reminder: $R(0) \triangleq \mathbb{E}_\infty [\psi(\Phi(n))\psi(\Phi(n))^\top]$ is positive definite (under mild assumption of linear independence of ψ) so that A is Hurwitz.

Regenerative relative TD(λ)-learning

Regenerative relative TD(λ) algorithm for average cost (on-policy)

For initialization $\theta_0, \zeta_0 \in \mathbb{R}^d$, the sequence of estimates are defined recursively:

$$\begin{aligned}\theta_{n+1} &= \theta_n + \alpha_{n+1} \zeta_n \mathcal{D}_{n+1} \\ \mathcal{D}_{n+1} &= (-H^\theta(\Phi(n)) - \delta \langle \mu, H^\theta \rangle + c_n + H^\theta(\Phi(n+1))) \Big|_{\theta=\theta_n} \\ \zeta_{n+1} &= \lambda \mathbb{1}\{\Phi(n+1) \neq z^\bullet\} \zeta_n + \psi_{(n+1)}\end{aligned}\tag{10.32}$$

Variance expected to be reduced compared to non-relative version

ODE approximation $\bar{f}(\theta) = A\theta + b$ can be shown to be invertible if δ is small enough

Conditions to ensure that A is Hurwitz not available; hence, better to use LSTD(λ)

See [1, Theorem 10.15] and the discussion below it.

Actors and critics

Parametrized (randomized) policy: $\check{\phi}^\theta(u|x)$ ←



Actor

Denote $c_\theta(x) = \sum_u \check{\phi}^\theta(u|x)c(x,u)$

$$P_\theta(x, x') = \sum_u \check{\phi}^\theta(u|x)P_u(x, x'),$$

$$T_\theta(z, z') = P_u(x, x')\check{\phi}^\theta(u'|x')$$



Critic

Actor-Critic Objective.

$$\Gamma(\theta) = \sum_{x \in \mathcal{X}} c_\theta(x) \pi_\theta(x) = \sum_{z \in \mathcal{Z}} c(z) \omega_\theta(z)$$

= average cost

Gradient and score function

Score function: $\Lambda^\theta(x, u) = \nabla_\theta \log[\check{\Phi}^\theta(u | x)]$

Note that $\nabla_\theta c_\theta(x) = \sum_u \check{\Phi}^\theta(u | x) \Lambda^\theta(x, u) c(x, u)$

$$\nabla_\theta P_\theta(x, x') = \sum_u \check{\Phi}^\theta(u | x) \Lambda^\theta(x, u) P_u(x, x'), \quad x, x' \in X, \theta \in \mathbb{R}^d$$

$$\underline{\Lambda^\theta(z')} = \nabla_\theta \log(T_\theta(z, z')), \quad z, z' \in Z$$

Gradient of critic and sensitivity theorem

Denote $Q_\theta(x, u) = Q$ -function of system controlled with $\check{\phi}^\theta$

It holds that $\nabla\Gamma(\theta) = E_{\omega_\theta} [\Lambda^\theta(\Phi(k))Q_\theta(\Phi(k))]$ (sensitivity theorem)

Proof: use Poisson equation $T_\theta Q_\theta = Q_\theta - c + \Gamma(\theta)$.

Proof. Expand the Poisson equation:

$$\sum_{z' \in \mathbf{X} \times \mathbf{U}} T_\theta(z, z') Q_\theta(z') = Q_\theta(z) - c(z) + \Gamma(\theta).$$

Take gradient on both sides:

$$\sum_{z' \in \mathbf{X} \times \mathbf{U}} \nabla_\theta T_\theta(z, z') Q_\theta(z') + T_\theta(z, z') \nabla_\theta Q_\theta(z') = \nabla_\theta Q_\theta(z) + \nabla_\theta \Gamma(\theta)$$

or

$$\sum_{z' \in \mathbf{X} \times \mathbf{U}} T_\theta(z, z') \frac{\nabla_\theta T_\theta(z, z')}{T_\theta(z, z')} Q_\theta(z') + T_\theta(z, z') \nabla_\theta Q_\theta(z') = \nabla_\theta Q_\theta(z) + \nabla_\theta \Gamma(\theta)$$

Use the score function:

$$\sum_{z' \in \mathbf{X} \times \mathbf{U}} T_\theta(z, z') \Lambda^\theta(z') Q_\theta(z') + T_\theta(z, z') \nabla_\theta Q_\theta(z') = \nabla_\theta Q_\theta(z) + \nabla_\theta \Gamma(\theta)$$

Take the expectation with respect to ϖ_θ on both sides:

$$\sum_{z \in \mathbf{X} \times \mathbf{U}} \varpi^\theta(z) \sum_{z' \in \mathbf{X} \times \mathbf{U}} T_\theta(z, z') \Lambda^\theta(z') Q_\theta(z') + T_\theta(z, z') \nabla_\theta Q_\theta(z') = \sum_{z \in \mathbf{X} \times \mathbf{U}} \varpi^\theta(z) \nabla_\theta Q_\theta(z) + \nabla_\theta \Gamma(\theta)$$

Use that $\sum_{z \in \mathbf{X} \times \mathbf{U}} \varpi^\theta(z) T_\theta(z, z') = \varpi^\theta(z')$:

$$\sum_{z' \in \mathbf{X} \times \mathbf{U}} \varpi^\theta(z) \Lambda^\theta(z') Q_\theta(z') + \varpi_\theta(z') \nabla_\theta Q_\theta(z') = \sum_{z \in \mathbf{X} \times \mathbf{U}} \varpi^\theta(z) \nabla_\theta Q_\theta(z) + \nabla_\theta \Gamma(\theta)$$

Change z' into z :

$$\sum_{z \in \mathbf{X} \times \mathbf{U}} \varpi^\theta(z) \Lambda^\theta(z) Q_\theta(z) = \nabla_\theta \Gamma(\theta).$$

This concludes the proof. □

Actor-critic method (ideal case)

Stochastic gradient method:

$$\theta_{n+1} = \theta_n - \alpha_{n+1} \check{\nabla}_{\Gamma}(n), \quad \check{\nabla}_{\Gamma}(n) \stackrel{\text{def}}{=} \Lambda^{\theta_n}(\Phi(n)) Q_{\theta_n}(\Phi(n))$$

Challenges:

- Q_{θ} is not known
 - Variance can be large
- ... addressed next

Actor-critic method

Idea: $H_\theta^\omega \approx Q_\theta$ where $H_\theta^\omega = \omega^\top \psi_\theta$, $\omega \in \mathbb{R}^{d'}$, $\theta \in \mathbb{R}^d$

Actor-Critic Algorithm

For initialization $\theta_0 \in \mathbb{R}^d$ and $\omega_0, \zeta_0 \in \mathbb{R}^{d'}$,

$$\theta_{n+1} = \theta_n - \alpha_{n+1} \check{\nabla}_\Gamma(n), \quad \check{\nabla}_\Gamma(n) \stackrel{\text{def}}{=} \Lambda^{\theta_n}(\Phi(n)) H_{\theta_n}^{\omega_n}(\Phi(n)) \quad (10.44a)$$

$$\Phi(n+1) \sim T_{\theta_n}(z, \cdot), \quad \text{with } z = \Phi(n) \quad (10.44b)$$

$$\left. \begin{aligned} \mathcal{D}_{n+1} &= \left\{ -H_\theta^\omega(\Phi(n)) + \tilde{c}_n + \mathbb{1}\{\Phi(n+1) \neq z^\bullet\} H_\theta^\omega(\Phi(n+1)) \right\} \Big|_{\substack{\theta=\theta_n \\ \omega=\omega_n}} \\ \omega_{n+1} &= \omega_n + \beta_{n+1} \zeta_n \mathcal{D}_{n+1} \\ \zeta_{n+1} &= \lambda \mathbb{1}\{\Phi(n+1) \neq z^\bullet\} \zeta_n + \psi_{\theta_{n+1}}(\Phi(n+1)) \\ \eta_{n+1} &= \eta_n + \beta_{n+1} \tilde{c}_n, \quad \tilde{c}_n = c(\Phi(n)) - \eta_n \end{aligned} \right\} \begin{array}{l} \text{Regenerative TD}(\lambda) \\ (10.44c) \end{array}$$

Analysis of actor-critic method

Proposition 10.17. *Suppose that $\lambda = 1$ and the step-size sequences satisfy $\lim_{n \rightarrow \infty} \frac{\beta_n}{\alpha_n} = \infty$*

Assume the following consistency condition holds:

$$\text{for each } \theta \in \mathbb{R}^n \text{ there is a } \omega_\theta^* \in \mathbb{R}^d \text{ satisfying } H_\theta^{\omega_\theta^*} = Q_\theta. \quad (10.45)$$

Then, the ODE approximation of (10.44) is gradient descent $\frac{d}{dt}\vartheta = -\nabla\Gamma(\vartheta)$. □

Consistency assumption (10.45) can be restrictive; remedy on next slide

See [1, Proposition 10.17].

Relaxing the consistency assumption

Assume: Compatible Features Property (CFP)

$$\Lambda_i^\theta \in \mathcal{H}^\theta \quad \text{for each } \theta \in \mathbb{R}^d \text{ and } 1 \leq i \leq d$$

(can always be done by design of \mathcal{H}^θ)

Proposition 10.18. *let $\hat{Q} \in \arg \min\{\|H - Q_\theta\|_{\omega_\theta}^2 : H \in \mathcal{H}^\theta\}$*

Then,

$$\nabla \Gamma(\theta) = \mathbb{E}_{\omega_\theta} [\Lambda^\theta(\Phi(k)) Q_\theta(\Phi(k))] = \mathbb{E}_{\omega_\theta} [\Lambda^\theta(\Phi(k)) \hat{Q}(\Phi(k))]$$

Proof. Since \hat{Q} is the projection of Q_θ on \mathcal{H}^θ , $\hat{Q} - Q_\theta$ is orthogonal to any function $H \in \mathcal{H}^\theta$, i.e.,

$$\mathbb{E}_{\varpi_\theta} \left[(\hat{Q} - Q_\theta)H \right] = 0.$$

By the CFP $\Lambda_i^\theta \in \mathcal{H}^\theta$, concluding the proof. □

Relaxing the consistency assumption

(continued)

Proposition 10.19. *Suppose that the assumptions of Prop. 10.17 hold, but with (10.45) replaced with the compatible features assumption (10.47).*

Then, the ODE approximation of (10.44) is unchanged: $\frac{d}{dt}\vartheta = -\nabla\Gamma(\vartheta)$. □

Hence, same conclusion but with realistic assumption (CFP)

See [1, Proposition 10.19].

Variance reduction through advantage

Result: For any function $g: X \rightarrow \mathbb{R}$, $0 = \mathbb{E}[g(X(k))\Lambda^\theta(\Phi(k))]$

Hence, we can subtract $G_n(X(n))$ from $H_{\theta_n}^{\omega_n}(\Phi(n))$ without impact

For instance, use $V_\theta^\omega(x, u) \stackrel{\text{def}}{=} H_\theta^\omega(x, u) - \underline{H}_\theta^\omega(x) = \omega^\top \tilde{\psi}_\theta(x, u)$

This gives: $\theta_{n+1} = \theta_n - \alpha_{n+1} \check{\nabla}_\Gamma^v(n)$, $\check{\nabla}_\Gamma^v(n) \stackrel{\text{def}}{=} \Lambda^{\theta_n}(\Phi(n)) \underline{V_{\theta_n}^{\omega_n}(\Phi(n))}$

The idea is the same as for the TD(λ)-learning algorithm applied to the advantage function: offset the Q-function $Q_\theta(x, u)$ by some function $g_\theta(x)$ because this will not change the value of $\mathbb{E}_{\varpi_\theta} [\Lambda^\theta(\Phi(n))Q_\theta(\Phi(n))]$ but is expected to reduce the variance of the approximated gradient.

Proof of the result. Observe that

$$\mathbb{E}_{\varpi_\theta} [g(X(n))\Lambda^\theta(\Phi(n))] = \sum_{x \in \mathbf{X}} \pi_\theta(x)g(x) \sum_{u \in \mathbf{U}} \check{\phi}^\theta(u|x)\Lambda^\theta(x, u).$$

It holds that

$$\sum_{u \in \mathbf{U}} \check{\phi}^\theta(u|x)\Lambda^\theta(x, u) = \sum_{u \in \mathbf{U}} \nabla_\theta \check{\phi}^\theta(u|x) = \nabla_\theta \sum_{u \in \mathbf{U}} \check{\phi}^\theta(u|x) = \nabla_\theta(1) = 0,$$

concluding the proof. □

Variance reduction through advantage

(continued)

Using the advantage instead of the Q -function is expected to induce a smaller variance of the method

Even better: if computable, use

$$\underline{\Lambda} H_{\theta}^{\omega}(x) = \mathbb{E}[\Lambda^{\theta}(\Phi(n)) H_{\theta}^{\omega}(\Phi(n)) \mid \mathcal{F}_n^{-}; X(n) = x]$$

The latter has smaller asymptotic variance than $\Lambda^{\theta_n}(\Phi(n)) V_{\theta_n}^{\omega_n}(\Phi(n))$

Caveats: variance of actor-critic methods

The theoretical results were obtained for $\lambda = 1$

However, despite the proposed remedies, the variance of the actor-critic method can remain large for $\lambda = 1$

In practice, sometimes better to use $\lambda < 1$ to tame the variance

Especially for continuous systems, $\lambda = 1$ may be unapplicable

Remember that

$$\zeta_n = \sum_{k=0}^n \lambda^k \mathbf{1} \left\{ \bigwedge_{0 \leq i \leq k} \Phi(n-i) \neq z^\bullet \right\} \psi(\Phi(n-k)).$$

If z^\bullet has a small steady-state probability (e.g., for continuous systems, this probability is zero in general), and $\lambda = 1$, then ζ_n can grow unbounded as $n \rightarrow \infty$. This will induce large variance in the method. Hence, even though the theoretical results (in terms of quality of the approximation of the limit point) were mostly obtained for $\lambda = 1$, in practice, it may be better to use $\lambda < 1$ to have a smaller variance.

Next course

- Online learning techniques for stochastic control
 - Bandit problem
 - Regret minimization
 - Exploration vs. exploitation trade-off

References

- [1] Sean Meyn. *Control systems and reinforcement learning*. Cambridge University Press, 2022.