

The existence of optimal control for continuous-time Markov decision processes in random environments

Jinghai Shao

Center for Applied Mathematics, Tianjin University

July 12, 2019

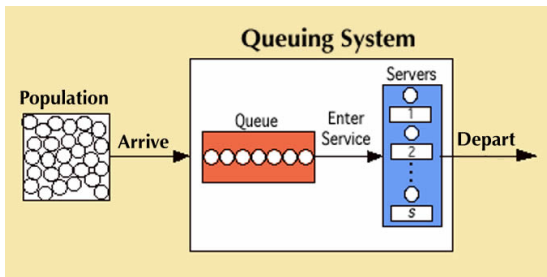
Contents

- ① Continuous-time Markov Decision Processes (CTMDP)
 - ▶ Background
 - ▶ Framework
- ② Optimal Markov control for CTMDP in random environments
 - ▶ Main results
 - ▶ Key steps in the arguments
- ③ Optimal Markov control for CTMDP with delay

Contents

- 1 Continuous-time Markov Decision Processes (CTMDP)
 - Background
 - Framework
- 2 Optimal Markov control for CTMDP in random environments
 - Main results
 - Key points in the arguments
- 3 Optimal control for CTMDP with delay

- CTMDPs have been extensively studied and widely applied in various application fields such as telecommunication, queueing systems, population processes, epidemiology, and so on.
- As an illustrative example, consider the controlled queueing systems:



Control Model

Consider the state space $\mathcal{S} = \{1, 2, \dots\}$, on which there exists a continuous-time Markov chain (Λ_t) with

$$(q_{ij}(a)) \quad \text{for } a \in U, \text{ action space.}$$

Assume

$$U \subset \mathbb{R}^k, \quad \text{compact}; \quad \sum_{j \in \mathcal{S}} q_{ij}(a) = 0, \quad \forall i \in \mathcal{S}, a \in U;$$

$$\sup_{a \in U} \sup_{i \in \mathcal{S}} q_i(a) < \infty.$$

For example, choose **appropriate control policy** to minimize the cost

- **finite-horizon expected cost:**

$$V_T(i, \pi) := \mathbb{E} \left[\int_0^T c(\Lambda_t, \pi_t) dt \right], \text{ where } T > 0.$$

- **infinite-horizon expected discounted cost:**

$$V(i, \pi) := \mathbb{E} \left[\int_0^\infty e^{-\lambda t} c(\Lambda_t, \pi_t) dt \right], \text{ where } \lambda > 0, \text{ discount factor.}$$

Randomized Markov policies: A randomized Markov policy is a real-valued function $\pi_t(C|i)$ that satisfies the following conditions:

- (i) For all $i \in \mathcal{S}$ and $C \in \mathcal{B}(U)$, the mapping $t \mapsto \pi_t(C|i)$ is measurable;
- (ii) For all $i \in \mathcal{S}$, $t \geq 0$, $C \mapsto \pi_t(C|i)$ is a probability measure on $\mathcal{B}(U)$.

stationary : if $\pi_t(C|i) \equiv \pi(C|i)$.

deterministic : if $\pi_t(C|i) = \delta_{u_t}(C|i)$, Dirac measure.

♣ Π : the set of all randomized Markov policies.

- * X.P. Guo, Hernandez-Lerma, Springer-Verlag, Berlin, 2009.
- * X.P. Guo, X. Huang, Y. Huang, *Finite-horizon optimality for CTMDPs with unbounded transition rates*, Adv. Appl. Prob. 2015.
- * X.P. Guo, U. Rieder, *Average optimality for CTMDPs in Polish spaces*, Ann. Appl. Probab. 2006.
- * A. Piunovskiy, Y. Zhang, *Discounted CTMDPs with unbounded rates: the convex analytic approach*, SIAM J. Control Optim. 2011.

An existing method

Consider

$$J_\lambda(i, \pi) := \mathbb{E} \left[\int_0^\infty e^{-\lambda t} c(\Lambda_t, \pi_t) dt \right],$$

and the corresponding value function

$$J_\lambda^*(i) := \inf_{\pi \in \Pi} J_\lambda(i, \pi).$$

Key point: The function J_λ^* satisfies the HJB equation

$$J_\lambda^*(i) = \inf_{a \in U} \left\{ \frac{c(i, a)}{\lambda + q_i(a)} + \frac{1}{\lambda + q_i(a)} \sum_{j \neq i} J_\lambda^*(j) q_{ij}(a) \right\}, \quad i \in \mathcal{S}.$$

Let

$$\varphi_{ij}^{(n)}(a) := \begin{cases} \frac{\delta_{ij}}{\lambda + q_i(f)} & n = 1, \\ \frac{1}{\lambda + q_i(f)} [\delta_{ij} + \sum_{k \neq i} q_{ik}(f) \varphi_{kj}^{(n-1)}(f)] & n = 2. \end{cases}$$

Then

$$\begin{aligned} J_\lambda(i, f) &= \sum_{j \in \mathcal{S}} \int_0^\infty e^{-\lambda t} c(j, f) P_f(0, i, t, j) dt \\ &= \sum_{j \in \mathcal{S}} c(j, f) \left[\lim_{n \rightarrow \infty} \varphi_{ij}^{(n)}(f) \right]. \end{aligned}$$

Framework

Let us consider further a diffusion process satisfying SDE:

$$dX_t = b(X_t, \Lambda_t)dt + \sigma(X_t, \Lambda_t)dB_t,$$

where (B_t) is a d -dimension B.M., $b : \mathbb{R}^d \times \mathcal{S} \rightarrow \mathbb{R}^d$, and $\sigma : \mathbb{R}^d \times \mathcal{S} \rightarrow \mathbb{R}^{d \times d}$.

The optimal control problem:

$$\inf_{\Pi} \mathbb{E} \left[\int_0^T f(t, X_t, \Lambda_t, \mu_t) dt + g(X_T, \Lambda_T) \right],$$

where Π is the set of admissible control policies which will be given later.

Some notations

- ① Let $\psi : [0, T] \rightarrow [0, \infty)$ be an increasing function such that

$$\lim_{r \rightarrow 0} \psi(r) = 0 \quad \forall r \in [0, T].$$

- ② $\mathcal{P}(U)$: all the probab. measures over U , endowed with the L^1 -Wasserstein distance, becoming a Polish space.
- ③ $\mathcal{D}([0, T]; \mathcal{P}(U))$: measurable maps $[0, T] \mapsto (\mathcal{P}(U), W_1)$, càdlàg.
- ④ Endow $\mathcal{D}([0, T]; \mathcal{P}(U))$ with the pseudopath topology, which makes it being a Polish space.
- ⑤ For $\mu : [0, T] \rightarrow \mathcal{P}(U)$ in $\mathcal{D}([0, T]; \mathcal{P}(U))$, put

$$w_\mu([a, b]) = \sup\{W_1(\mu_t, \mu_s); s, t \in [a, b]\}, \quad a, b \in [0, T], a < b;$$

$$w_\mu''(\delta) = \sup \min \{W_1(\mu_t, \mu_{t_1}), W_1(\mu_t, \mu_{t_2})\},$$

where the supremum is taken over t_1, t , and t_2 satisfying

$$t_1 \leq t \leq t_2, \quad t_2 - t_1 \leq \delta.$$

The process (X_t) is determined by the following SDE:

$$dX_t = b(X_t, \Lambda_t)dt + \sigma(X_t, \Lambda_t)dB_t, \quad (1)$$

where (B_t) is a Brownian motion; (Λ_t) is a continuous-time Markov process on \mathcal{S} associated with the q -pair $(q_i(u), q_{ij}(u))$ satisfying

$$\mathbb{P}(\Lambda_{t+\delta} = j | \Lambda_t = i, \mu_t = \mu) = \begin{cases} q_{ij}(\mu)\delta + o(\delta) & i \neq j, \\ 1 - q_i(\mu)\delta + o(\delta), & i = j, \end{cases} \quad (2)$$

provided $\delta > 0$. The decision-maker still tries to minimize the cost through controlling the transition rates of the Markov chain (Λ_t) , but now the cost function may depend on the diffusion process (X_t) .

Definition

A ψ -relaxed control is a term $\alpha = (\Omega, \mathcal{F}, \mathcal{F}_t, \mathbb{P}, B_t, X_t, \Lambda_t, \mu_t, s, x, i)$ such that

- (1) $(s, x, i) \in [0, T] \times \mathbb{R}^d \times \mathcal{S}$;
- (2) $(\Omega, \mathcal{F}, \mathbb{P})$ is a probability space with the filtration $\{\mathcal{F}_t\}_{t \in [0, T]}$;
- (3) (B_t) is a d -dim B.M. on $(\Omega, \mathcal{F}, \mathcal{F}_t, \mathbb{P})$, and (X_t, Λ_t) is a stochastic process on $\mathbb{R}^d \times \mathcal{S}$ satisfying (1) and (2) with $X_s = x, \Lambda_s = i$;
- (4) $\mu_t \in \mathcal{P}(U)$ is adapted to the σ -field generated by Λ_t , $t \mapsto \mu_t$ is in $\mathcal{D}([0, T]; \mathcal{P}(U))$ almost surely, and for every $i' \in \mathcal{S}$ the curve $t \mapsto \nu_t(\cdot, i') := \mu_t(\cdot | \Lambda_t = i')$ satisfies

$$w_v''(\delta) \leq \psi(\delta), \quad \delta \in (0, T];$$

- The collection of all ψ -relaxed control with initial value (s, x, i) is denoted by $\tilde{\Pi}_{s, x, i}$.

Actually, the randomized policy can be viewed as a Markov feedback control.

$$\begin{aligned}\mu_t(C) &= \sum_{i \in \mathcal{S}} \pi_t(C|i) \mathbf{1}_{\Lambda_t=i} \\ &= \pi_t(C|\Lambda_t), \quad t \geq 0.\end{aligned}$$

Contents

- 1 Continuous-time Markov Decision Processes (CTMDP)
 - Background
 - Framework
- 2 Optimal Markov control for CTMDP in random environments
 - Main results
 - Key points in the arguments
- 3 Optimal control for CTMDP with delay

Assumptions

- (H1) $U \subset \mathbb{R}^k$ is a compact set for some $k \in \mathbb{N}$.
- (H2) $\forall u \in U$, $(q_{ij}(u))$ is conservative. $M := \sup_{u \in U} \sup_{i \in \mathcal{S}} q_i(u) < \infty$.
- (H3) $\forall i, j \in \mathcal{S}$, $u \mapsto q_{ij}(u)$ is continuous on U .
- (H4) \exists a compact function $\Phi : \mathcal{S} \rightarrow [1, \infty)$, a compact set $B_0 \in \mathcal{B}(\mathcal{S})$, constants $\lambda > 0$ and $\kappa_0 < \infty$ such that

$$Q_u \Phi(i) := \sum q_{ij}(u) \Phi(j) \leq \lambda \Phi(i) + \kappa_0 \mathbf{1}_{B_0}(i), \quad i \in \mathcal{S}, u \in U.$$

- (H5) \exists a constant $C_1 > 0$ such that

$$|b(x, i) - b(y, i)|^2 + \|\sigma(x, i) - \sigma(y, i)\|^2 \leq C_1 |x - y|^2, \quad x, y \in \mathbb{R}^d, i \in \mathcal{S},$$

where $|x|^2 = \sum_{k=1}^d x_k^2$, $\|\sigma\|^2 = \text{tr}(\sigma\sigma')$.

- (H6) $\exists C_2 > 0$ such that $|b(x, i)|^2 + \|\sigma(x, i)\|^2 \leq C_2(1 + |x|^2)$, $x \in \mathbb{R}^d, i \in \mathcal{S}$.

Theorem 1

Assume that (H1)-(H6) hold, and f, g are lower semi-continuous functions bounded from below. Then for every $s \in [0, T)$, $x \in \mathbb{R}^d$, $i \in \mathcal{S}$, there exists an optimal ψ -relaxed control $\alpha^* \in \tilde{\Pi}_{s,x,i}$, i.e.

$$\begin{aligned} V(s, x, i) &= J(s, x, i, \alpha^*) \\ &= \inf_{\alpha \in \tilde{\Pi}_{s,x,i}} \mathbb{E} \left[\int_s^T f(t, X_t, \Lambda_t, \mu_t) dt + g(X_T, \Lambda_T) \right]. \end{aligned}$$

Theorem 2

Suppose (H1)-(H6) hold. Assume that f and g are continuous functions and there exists a positive constant C_3 such that

$$\begin{aligned} |f(t, x, i, u) - f(t, x', i, u)| + |g(x, i) - g(x', i)| &\leq C_3|x - x'|, \\ |f(t, x, i, u)| + |g(x, i)| &\leq C_3, \end{aligned}$$

for every $t \in [0, T]$, $x, x' \in \mathbb{R}^d$, $i \in \mathcal{S}$ and $u \in U$. Then $V(s, x, i)$ is **continuous** on $[0, T] \times \mathbb{R}^d \times \mathcal{S}$.

Theorem 3 (Dynamic programming principle)

Assume all the conditions of Theorem 2 are still valid. Then for $s < t < T$,

$$V(s, x, i) = \inf \left\{ \mathbb{E}_\alpha \left[\int_s^t f(r, X_r, \Lambda_r, \mu_r) dr + V(t, X_t, \Lambda_t) \right]; \alpha \in \tilde{\Pi}_{s,x,i} \right\}.$$

Key steps to prove Theorem 1

To simplify the proof, transform the relaxed controls into the canonical path space. Let

$$\mathcal{U} = \{\nu \in \mathcal{D}([0, T]; \mathcal{P}(U)); w''_\nu(\delta) \leq \psi(\delta)\}$$
$$\mathcal{Y} = C([0, T]; \mathbb{R}^d) \times \mathcal{D}([0, T]; \mathcal{S}) \times \mathcal{U}.$$

Denote by $\tilde{\mathcal{D}}, \tilde{\mathcal{U}}$ the Borel measurable sets, and $\tilde{\mathcal{D}}_t, \tilde{\mathcal{U}}_t$ the σ -fields up to time t . Each ψ -relaxed control $\alpha = (\Omega, \mathcal{F}, \mathcal{F}_t, \mathbb{P}, B_t, X_t, \Lambda_t, \mu_t, s, x, i)$ can be transformed into \mathcal{Y} via the map

$$\Psi(\omega) = (X_t(\omega), \Lambda_t(\omega), \mu_t(\omega))_{t \in [0, T]},$$

with $X_r := x, \Lambda_r := i, \mu_r := \mu_s, \forall r \in [0, s]$.

We can use $R := \mathbb{P} \circ \Psi^{-1}$ to represent the control α in canonical space \mathcal{Y} .

Key steps to prove Theorem 1

Consider the nontrivial case $V(0, x, i) < \infty$, and $\exists (R_n)_{n \geq 1}$ such that

$$\lim_{n \rightarrow \infty} J(0, x, i, R_n) = V(0, x, i). \quad (\text{e1})$$

- 1 Prove the tightness of the distributions of $(X_t)_{t \in [0, T]}$, $(\Lambda_t)_{t \in [0, T]}$ and $(\mu_t)_{t \in [0, T]}$ under the sequence of probab. measures R_n , $n \geq 1$.
 - Taking a subsequence if necessary, using Skorokhod's representation theorem, \exists a probab. space $(\Omega', \mathcal{F}', \mathbb{P}')$ and $(X_t^{(n)}, \Lambda_t^{(n)}, \mu_t^{(n)})_{t \in [0, T]}$ taking values in \mathcal{Y} with the distribution R_n , such that

$$(X_t^{(n)}, \Lambda_t^{(n)}, \mu_t^{(n)})_{t \in [0, T]} \longrightarrow (X_t^{(0)}, \Lambda_t^{(0)}, \mu_t^{(0)})_{t \in [0, T]}, \text{ a.s. } n \rightarrow \infty.$$

Key steps to prove Theorem 1

- 2 Prove that $(X_t^{(0)}, \Lambda_t^{(0)}, \mu_t^{(0)})$ satisfies

▶
$$X_t^{(0)} = x + \int_0^t b(X_s^{(0)}, \Lambda_s^{(0)}) ds + \int_0^t \sigma(X_s^{(0)}, \Lambda_s^{(0)}) dB_s.$$

▶

$$\mathbb{P}(\Lambda_{t+\delta}^{(0)} = j | \Lambda_t^{(0)} = i', \mu_t^{(0)} = \mu) = \begin{cases} q_{i'j}(\mu)\delta + o(\delta) & i' \neq j, \\ 1 - q_{i'}(\mu)\delta + o(\delta), & i' = j. \end{cases}$$

▶ $\mu_t^{(0)}$ is adapted to $\sigma(\Lambda_t^{(0)})$.

- 3 the control α^* associated with $(X_t^{(0)}, \Lambda_t^{(0)}, \mu_t^{(0)})$ is an optimal ψ -relaxed control.

Contents

- 1 Continuous-time Markov Decision Processes (CTMDP)
 - Background
 - Framework
- 2 Optimal Markov control for CTMDP in random environments
 - Main results
 - Key points in the arguments
- 3 Optimal control for CTMDP with delay

Let $\psi : [0, T] \rightarrow [0, \infty)$ be increasing, $\lim_{r \rightarrow 0} \psi(r) = 0$.

For $r_0 \in (0, T)$, define a shift operator $\theta_{r_0} : \mathcal{D}([0, T]; \mathcal{S}) \rightarrow \mathcal{D}([0, T]; \mathcal{S})$ by

$$(\theta_{r_0} \lambda)(t) = \lambda_{(t-r_0) \vee 0}, \quad t \in [0, T].$$

Moreover, $\theta_{r_0}^k \lambda(t) := \lambda_{(t-kr_0) \vee 0}$ for $\lambda \in \mathcal{D}([0, T]; \mathcal{S})$, $k \geq 0$.

$m \geq 1$ is a fixed integer. A functional $h : [0, T] \times \mathcal{S}^{m+1} \rightarrow \mathcal{P}(U)$ is said to be in the class Υ_ψ if for every $i_0, \dots, i_m \in \mathcal{S}$, $t \mapsto \tilde{\mu}(t) := h(t, i_0, \dots, i_m)$ satisfies

$$w_{\tilde{\mu}}([t_1, t_2]) \leq \psi(|t_2 - t_1|), \quad t_1, t_2 \in [0, T].$$

Definition: For $s \in [0, T)$ and $i \in \mathcal{S}$, a *history-dependent control* is a term $\alpha = (\Lambda_t, \mu_t)$ such that

(1) (Λ_t) is an \mathcal{F}_t -adapted jumping process satisfying

$$\mathbb{P}(\Lambda_{t+\delta} = j | \Lambda_t = i, \mu_t = \mu) = \begin{cases} q_{ij}(\mu)\delta + o(\delta), & \text{if } i \neq j, \\ 1 + q_{ii}(\mu)\delta + o(\delta), & \text{otherwise,} \end{cases}$$

with initial value $\Lambda_s = i$ for $s \in [0, T)$ and $i \in \mathcal{S}$.

(2) There exists $h \in \Upsilon_\psi$ such that

$$\mu_t = h(t, \theta_{r_0}^0 \Lambda(t), \dots, \theta_{r_0}^m \Lambda(t)).$$

The collection of all history-dependent α with initial value (s, i) is denoted by $\Pi_{s,i}$. Let $f : [0, T] \times \mathcal{S} \times \mathcal{P}(U) \rightarrow [0, \infty)$, $g : \mathcal{S} \rightarrow [0, \infty)$ be two lower semi-continuous functions. The expected cost for the history-dependent control $\alpha \in \Pi_{s,i}$ is defined by

$$J(s, i, \alpha) = \mathbb{E} \left[\int_s^T f(t, \Lambda_t, \mu_t) dt + g(\Lambda_T) \right],$$

and the value function is defined by

$$V(s, i) = \inf_{\alpha \in \Pi_{s,i}} J(s, i, \alpha).$$

A history-dependent control $\alpha^* \in \Pi_{s,i}$ is said to be *optimal*, if

$$V(s, i) = J(s, i, \alpha^*).$$

- ① $\mu_t = h(\Lambda_t)$ for some $h : \mathcal{S} \rightarrow \mathcal{P}(U)$. In this situation, α is corresponding to the stationary randomized Markov policy studied by many works.
- ② $\mu_t = h(\Lambda_{(t-r_0) \vee 0})$ for some $h : \mathcal{S} \rightarrow \mathcal{P}(U)$. Now the control policies are purely determined by the jumping process with a positive delay. This kind of controls is very natural to be used in the realistic application.

Assumptions:

(A1) $\mu \mapsto q_{ij}(\mu)$ is continuous $\forall i, j \in \mathcal{S}$, and $M := \sup_{i \in \mathcal{S}} \sup_{\mu \in \mathcal{P}(U)} q_i(\mu) < \infty$.

(A2) There exists a compact function $\Phi : \mathcal{S} \rightarrow [1, \infty)$, a compact set $B_0 \subset \mathcal{S}$, constants $\lambda_0 > 0$ and $\kappa_0 \geq 0$ such that

$$Q_\mu \Phi(i) := \sum_{j \neq i} q_{ij}(\mu) (\Phi(j) - \Phi(i)) \leq \lambda_0 \Phi(i) + \kappa_0 \mathbf{1}_{B_0}(i).$$

(A3) There exists a $K \in \mathbb{N}$ such that for every $i \in \mathcal{S}$ and $\mu \in \mathcal{P}(U)$, $q_{ij}(\mu) = 0$, if $|j - i| > K$.

Theorem 4

Assume (A1)-(A3) hold. Then for every $s \in [0, T)$, $i \in \mathcal{S}$, there exists an optimal control $\alpha^* \in \Pi_{s,i}$.

- X.P. Guo, X.X. Huang, Y.H. Huang, Adv. Appl. Prob. 2015.

Further work

Consider the following SDE:

$$dX_t = b(X_t, \Lambda_t, \mu_t)dt + \sigma(X_t, \Lambda_t, \mu_t)dB_t,$$

where $b : \mathbb{R}^d \times \mathcal{S} \times \mathcal{P}(U) \rightarrow \mathbb{R}^d$, $\sigma : \mathbb{R}^d \times \mathcal{S} \times \mathcal{P}(U) \rightarrow \mathbb{R}^{d \times d}$, and $(B_t)_{t \geq 0}$ is a d -dimensional \mathcal{F}_t -Brownian motion. Here $(\Lambda_t)_{t \geq 0}$ is a continuous-time jumping process on \mathcal{S} satisfying

$$\mathbb{P}(\Lambda_{t+\delta} = j | \Lambda_t = i, X_t = x, \nu_t = \nu) = \begin{cases} q_{ij}(x, \nu)\delta + o(\delta), & \text{if } j \neq i, \\ 1 + q_{ii}(x, \nu)\delta + o(\delta), & \text{otherwise,} \end{cases}$$

provided $\delta > 0$ for every $x \in \mathbb{R}^d$, $\nu \in \mathcal{P}(U)$, $i, j \in \mathcal{S}$.

Thank You For Your Attention !

EMAIL: shaojh@tju.edu.cn