

Markov Decision Processes with Applications

Ulrich Rieder

University of Ulm



ulm university universität
uulm

Markov Decision Processes with Applications

- Markov Decision Processes

Basic results, Computational aspects

- Partially Observable Markov Decision Processes

Hidden Markov models, Filtered MDPs

Bandit problems, Consumption-Investment problems

- Continuous-Time Markov Decision Processes

Piecewise deterministic MDPs

Parallel queueing model

Markov Decision Processes

$(E, A, D_n, Q_n, r_n, g_N)$ with horizon N

- E state space
- A action space
- $D_n \subset E \times A$ admissible state-action pairs at time n
- Q_n $Q_n(\cdot|x, a)$ transition law at time n
- $r_n : D_n \rightarrow \mathbb{R}$ reward function at time n
- $g_N : E \rightarrow \mathbb{R}$ terminal reward function at time N

decision rule at time n $f_n : E \rightarrow A$ measurable and $f_n(x) \in D_n(x)$ for all $x \in E$

policy $\pi := (f_0, f_1, \dots, f_{N-1})$

For $n = 0, 1, \dots, N$ define the value functions

$$V_{n\pi}(x) := E_x^\pi \left[\sum_{k=n}^{N-1} r_k(X_k, f_k(X_k)) + g_N(X_N) \right]$$

$$V_n(x) := \sup_{\pi} V_{n\pi}(x), \quad x \in E$$

π is called **optimal** if $V_{0\pi}(x) = V_0(x)$ for all $x \in E$.

Integrability Assumption (A_N) :

For $n = 0, 1, \dots, N$

$$\sup_{\pi} E_x^\pi \left[\sum_{k=n}^{N-1} r_k^+(X_k, f_k(X_k)) + g_N^+(X_N) \right] < \infty, \quad x \in E$$

Bertsekas/Shreve (1978), Hernandez-Lerma/Lasserre (1996)...

Puterman (1994), Feinberg/Schwartz (2002) ...

Bäuerle/Rieder (2011)

Let $\mathbb{M}(E) := \{v : E \rightarrow [-\infty, \infty) \mid v \text{ is measurable}\}$ and

define the following operators for $v \in \mathbb{M}(E)$:

$$(L_n v)(x, a) := r_n(x, a) + \int v(x') Q_n(dx' | x, a), \quad (x, a) \in D_n$$

$$(T_{nf_n} v)(x) := (L_n v)(x, f_n(x))$$

$$(T_n v)(x) := \sup_{a \in D_n(x)} (L_n v)(x, a), \quad x \in E \quad \text{Note: } T_n v \notin \mathbb{M}(E)!$$

A decision rule f_n is called a **maximizer** of v at time n if $T_{nf_n} v = T_n v$.

Reward Iteration: $V_{n\pi} = T_{nf_n} V_{n+1, \pi}$, $V_{N\pi} = g_N$.

Bellman Equation: $V_n = T_n V_{n+1}$, $V_N = g_N$.

Verification Theorem: Let $(v_n) \subset \mathbb{M}(E)$ be a solution of the Bellman equation.

a) $v_n \geq V_n$ for $n = 0, 1, \dots, N$.

b) If f_n^* is a maximizer of v_{n+1} for $n = 0, 1, \dots, N - 1$, then $v_n = V_n$ and the policy

$(f_0^*, f_1^*, \dots, f_{N-1}^*)$ is optimal.

Structure Assumption (SA_N): There exist sets $\mathbb{M}_n \subset \mathbb{M}(E)$ of measurable functions and sets Δ_n of decision rules such that for all $n = 0, 1, \dots, N - 1$:

(i) $g_N \in \mathbb{M}_N$.

(ii) If $v \in \mathbb{M}_{n+1}$ then $T_n v$ is well-defined and $T_n v \in \mathbb{M}_n$.

(iii) For all $v \in \mathbb{M}_{n+1}$ there exists a maximizer f_n of v with $f_n \in \Delta_n$.

Structure Theorem:

Assume (SA_N). Then it holds:

a) $V_n \in \mathbb{M}_n$ and (V_n) is a solution of the Bellman equation.

b) $V_n = T_n T_{n+1} \dots T_{N-1} g_N$.

c) For $n = 0, 1, \dots, N - 1$ there exists a maximizer f_n of V_{n+1} with $f_n \in \Delta_n$, and every sequence of maximizers f_n^* of V_{n+1} defines an optimal policy $(f_0^*, f_1^*, \dots, f_{N-1}^*)$ for the N -stage Markov Decision Problem.

$b : E \rightarrow \mathbb{R}_+$ is called an **upper bounding function** if there exist $c_r, c_g, \alpha_b \in \mathbb{R}_+$ such that for all $n = 0, 1, \dots, N - 1$

$$(i) \ r_n^+(x, a) \leq c_r b(x).$$

$$(ii) \ g_N^+(x) \leq c_g b(x).$$

$$(iii) \ \int b(x') Q_n(dx' | x, a) \leq \alpha_b b(x).$$

$$\alpha_b := \sup_{(x,a) \in D} \frac{\int b(x') Q(dx' | x, a)}{b(x)}. \text{ Define } \|v\|_b := \sup_{x \in E} \frac{|v(x)|}{b(x)}.$$

$$\mathbb{B}_b := \{v \in \mathbb{M}(E) \mid \|v\|_b < \infty\}, \ \mathbb{B}_b^+ := \{v \in \mathbb{M}(E) \mid \|v^+\|_b < \infty\}.$$

$b : E \rightarrow \mathbb{R}_+$ is called a **bounding function** if there exist $c_r, c_g, \alpha_b \in \mathbb{R}_+$ such that for all $n = 0, 1, \dots, N - 1$

$$(i) \ |r_n(x, a)| \leq c_r b(x).$$

$$(ii) \ |g_N(x)| \leq c_g b(x).$$

$$(iii) \ \int b(x') Q_n(dx' | x, a) \leq \alpha_b b(x).$$

Theorem: Suppose the N-stage MDP has an upper bounding function b and for all $n = 0, 1, \dots, N - 1$ it holds:

- (i) $D_n(x)$ is compact and $x \rightarrow D_n(x)$ is upper semicontinuous (usc).
- (ii) $(x, a) \rightarrow \int v(x')Q_n(dx'|x, a)$ is usc for all usc $v \in \mathbb{B}_b^+$.
- (iii) $(x, a) \rightarrow r_n(x, a)$ is usc .
- (iv) $x \rightarrow g_N(x)$ is usc.

Then the sets $\mathbb{M}_n := \{v \in \mathbb{B}_b^+ | v \text{ is usc}\}$ and $\Delta_n := \{f_n \text{ decision rule at time } n\}$ satisfy the Structure Assumption (SA_N), in particular: $V_n \in \mathbb{M}_n$ and there exists an optimal policy $(f_0^*, f_1^*, \dots, f_{N-1}^*)$ with $f_n^* \in \Delta_n$.

Markov Decision Processes with Infinite Time Horizon

We consider a stationary MDP with $\beta \in (0, 1]$ and $N = \infty$.

$$J_{\infty\pi}(x) := E_x^\pi \left[\sum_{k=0}^{\infty} \beta^k r(X_k, f_k(X_k)) \right]$$

$$J_\infty(x) := \sup_{\pi} J_{\infty\pi}(x), \quad x \in E.$$

Integrability Assumption (A):

$$\sup_{\pi} E_x^\pi \left[\sum_{k=0}^{\infty} \beta^k r^+(X_k, f_k(X_k)) \right] < \infty, \quad x \in E$$

Convergence Assumption (C):

$$\lim_{n \rightarrow \infty} \sup_{\pi} E_x^\pi \left[\sum_{k=n}^{\infty} \beta^k r^+(X_k, f_k(X_k)) \right] = 0, \quad x \in E$$

Then it holds: $J_{\infty\pi} = \lim_n J_{n\pi}$

limit value function $J := \lim_n J_n \geq J_\infty$. Note: $J \neq J_\infty$ and $J_\infty \notin \mathbb{M}(E)$!

Verification Theorem: Assume (C). Let $v \in \mathbb{M}(E)$ be a fixed point of T such that $v \geq J_\infty$. If f^* is a maximizer of v , then $v = J_\infty$ and the stationary policy (f^*, f^*, \dots) is optimal for the infinite-stage Markov Decision Problem.

Structure assumption (SA):

There exist a set $\mathbb{M} \subset \mathbb{M}(E)$ of measurable functions and a set Δ of decision rules such that:

- (i) $0 \in \mathbb{M}$.
- (ii) If $v \in \mathbb{M}$ then Tv is well-defined and $Tv \in \mathbb{M}$.
- (iii) For all $v \in \mathbb{M}$ there exists a maximizer f of v with $f \in \Delta$.
- (iv) $J \in \mathbb{M}$ and $J = TJ$.

Structure Theorem: Let (C) and (SA) be satisfied. Then it holds:

- a) $J_\infty \in \mathbb{M}$, $J_\infty = TJ_\infty$ and $J_\infty = J$.
- b) There exists a maximizer $f \in \Delta$ of J_∞ , and every maximizer f^* of J_∞ defines an optimal stationary policy (f^*, f^*, \dots) .

Theorem: Suppose the stationary MDP has an upper bounding function b with $\beta\alpha_b < 1$ and it holds:

- (i) $D(x)$ is compact and $x \rightarrow D(x)$ is usc.
- (ii) $(x, a) \rightarrow \int v(x')Q(dx'|x, a)$ is usc for all usc $v \in \mathbb{B}_b^+$.
- (iii) $(x, a) \rightarrow r(x, a)$ is usc.

Then it holds:

- (a) $J_\infty \in \mathbb{B}_b^+$, $J_\infty = TJ_\infty$ and $J_\infty = J$ **(value iteration)**.
- (b) b is usc $\implies J_\infty$ is usc.
- (c) $\emptyset \neq LsD_n^*(x) \subset D_\infty^*(x)$ for all $x \in E$ **(policy iteration)**.
- (d) There exists a decision rule f^* with $f^*(x) \in LsD_n^*(x)$ for all $x \in E$, and the stationary policy (f^*, f^*, \dots) is optimal.

$$\alpha_b := \sup_{(x,a) \in D} \frac{\int b(x')Q(dx'|x,a)}{b(x)}$$

Contracting Markov Decision Processes

Structure Theorem: Let b be a bounding function and $\beta\alpha_b < 1$. If there exists a closed subset $\mathbb{M} \subset \mathbb{B}_b$ and a set Δ of decision rules such that:

- (i) $0 \in \mathbb{M}$.
- (ii) $T : \mathbb{M} \rightarrow \mathbb{M}$.
- (iii) For all $v \in \mathbb{M}$ there exists a maximizer f of v with $f \in \Delta$.

Then it holds:

- a) $J_\infty \in \mathbb{M}$, $J_\infty = TJ_\infty$ and $J_\infty = J$.
- b) J_∞ is the unique fixed point of T in \mathbb{M} .
- c) There exists a maximizer $f \in \Delta$ of J_∞ , and every maximizer f^* of J_∞ defines an optimal stationary policy (f^*, f^*, \dots) .

Howard's Policy Improvement Algorithm

Let J_f be the value function of the stationary policy (f, f, \dots) .

Denote

$$D(x, f) := \{a \in D(x) \mid (LJ_f)(x, a) > J_f(x)\}$$

Let the Markov decision process be contracting.

Then it holds:

a) If for some subset $E_0 \subset E$

$$g(x) \in D(x, f) \text{ for } x \in E_0$$

$$g(x) = f(x) \text{ for } x \notin E_0$$

then $J_g \geq J_f$ and $J_g(x) > J_f(x)$ for $x \in E_0$.

In this case the decision rule g is called an **improvement** of f .

b) If $D(x, f) = \emptyset$ for all $x \in E$, then the stationary policy (f, f, \dots) is optimal.

Remark: (f, f, \dots) is optimal $\iff f$ cannot be improved.

Partially Observable Markov Decision Processes

- $E_X \times E_Y$ state space x observable state, y unobservable state
- A action space
- $D \subset E_X \times A$ admissible state-action pairs, $D(x) \subset A$
- $Q(\cdot|x, y, a)$ transition law
- Q_0 initial distribution (prior distribution) of Y_0
- $r(x, y, a)$ reward function
- $g(x, y)$ terminal reward function
- $\beta \in (0, 1]$ discount factor

Examples : Hidden Markov Model (HMM), Bayesian Decision Model

decision rule at time n $f_n(x_0, a_0, x_1, \dots, x_n) = f_n(h_n)$

policy $\pi = (f_0, f_1, \dots, f_{N-1})$ finite horizon: $N < \infty$

Rieder (1975), Elliott et al. (1995), Bäuerle/Rieder (2011) ...

$$J_{N\pi}(x) := E_x^\pi \left[\sum_{n=0}^{N-1} \beta^n r(X_n, Y_n, f_n(H_n)) + \beta^N g(X_N, Y_N) \right]$$

$$J_N(x) := \sup_{\pi} J_{N\pi}(x), \quad x \in E_X$$

For $n = 0, 1, \dots$ and $C \subset E_Y$ define

$$\mu_n(C | X_0, A_0, X_1, \dots, X_n) := P_x^\pi(Y_n \in C | X_0, A_0, X_1, \dots, X_n)$$

a posteriori-distribution at time n

Filter Equation

$$\mu_0 = Q_0 \text{ and } \mu_{n+1}(\cdot | H_n, A_n, X_{n+1}) = \Phi(X_n, \mu_n(\cdot | H_n), A_n, X_{n+1})$$

where

$$\Phi(x, \rho, a, x')(C) := \frac{\int_C \left[\int q(x', y' | x, y, a) \rho(dy) \right] \nu(dy')}{\int_{E_Y} \left[\int q(x', y' | x, y, a) \rho(dy) \right] \nu(dy')}, \quad C \subset E_Y, \rho \in \mathbb{P}(E_Y)$$

Bayes-Operator

Filtered Markov Decision Process

- $E' := E_X \times \mathbb{P}(E_Y) \ni (x, \rho)$ enlarged state space
- A and $D(x, \rho) := D(x)$
- $Q^X(B|x, \rho, a) := \int Q(B \times E_Y|x, y, a)\rho(dy)$, $B \subset E_X$
 $Q'(B \times C|x, \rho, a) := \int_B 1_C(\Phi(x, \rho, a, x'))Q^X(dx'|x, \rho, a)$, $C \subset \mathbb{P}(E_Y)$
- $r'(x, \rho, a) := \int r(x, y, a)\rho(dy)$
- $g'(x, \rho) := \int g(x, y)\rho(dy)$

Theorem:

a) $J_{N\pi}(x) = J'_{N\pi}(x, Q_0)$ and $J_N(x) = J'_N(x, Q_0)$.

b) Assume (SA_N) . Then the Bellman equation holds, i.e.

$$V'_N(x, \rho) := \beta^N g'(x, \rho)$$

$$V'_n(x, \rho) := \sup_{a \in D(x)} \left\{ r'(x, \rho, a) + \int V'_{n+1}(x', \Phi(x, \rho, a, x')) Q^X(dx'|x, \rho, a) \right\}.$$

Let f'_n be a maximizer of V'_{n+1} for $n = 0, \dots, N - 1$. Then the policy

$\pi^* := (f_0^*, f_1^*, \dots, f_{N-1}^*)$ is optimal for the N -stage POMDP, where

$$f_n^*(h_n) := f'_n(x_n, \mu_n(\cdot|h_n)), \quad h_n = (x_0, a_0, x_1, \dots, x_n).$$

Note that $V'_n(x, \rho) = \beta^n J'_{N-n}(x, \rho)$, $n = 0, \dots, N$

Computational aspects

Kalman Filter

Sufficient Statistics

Bandit Problems

unknown success probabilities $\theta_1 \in [0, 1]$ and $\theta_2 \in [0, 1]$

$Q_0 =$ product of two Uniform-distributions of (θ_1, θ_2)

Aim: maximize the expected number of successes in a finite or infinite number of trials

- $E' := \mathbb{N}_0^2 \times \mathbb{N}_0^2 \ni (m_1, n_1, m_2, n_2) = \rho$
- $A = \{1, 2\}$
- Bayes-Operator $\Phi(\rho, a, \{\text{success}\}) = \rho + e_{2a-1}$
- $r'(\rho, a) := \frac{m_a+1}{m_a+n_a+2}$
- $\beta \in (0, 1]$.

$N < \infty$: There exists an optimal policy.

monotonicity results: stay-on-a-winner property

stopping property if θ_2 is known.

$N = \infty$ and $\beta \in (0, 1)$:

For $K \in \mathbb{R}$ let $J(m, n; K)$ be the unique solution of

$$v(m, n) = \max\{K, \beta(p(m, n)v(m+1, n) + (1-p(m, n))v(m, n+1))\}$$

for $(m, n) \in \mathbb{N}_0^2$ and $p(m, n) := \frac{m+1}{m+n+2}$.

Define the **Gittins-Index**

$$I(m, n) := \min\{K \mid J(m, n; K) = K\}$$

Then it holds:

The stationary Index-policy (f^*, f^*, \dots) is optimal for the infinite-stage Bandit problem where

$$f^*(m_1, n_1, m_2, n_2) = \begin{cases} 1 & \text{if } I(m_1, n_1) \geq I(m_2, n_2) \\ 2 & \text{if } I(m_1, n_1) < I(m_2, n_2). \end{cases}$$

Gittins (1989), Whittle (1980), (1988)

Cox-Ross-Rubinstein Model

- Bond $B_n = (1 + i)^n$
- Stock $S_n = S_0 \cdot \prod_{k=1}^n Y_k$ (Y_k) independent and identically distributed
 $P(Y_k = \mathbf{u}) = \theta = 1 - P(Y_k = \mathbf{d})$ unknown up-probability θ

$Q_0 =$ Uniform-distribution of θ

(NA) : $\mathbf{d} < 1 + i < \mathbf{u}$

$\pi_n =$ amount of money invested in the stock at time n

Then it holds for the wealth process:

$$X_{n+1}^\pi = X_n^\pi (1 + i) + \pi_n (Y_{n+1} - 1 - i), \quad X_0^\pi = x > 0$$

Utility function $U : \mathbb{R}_+ \longrightarrow \mathbb{R}_+$, strictly increasing and concave

$$(P) \left\{ \begin{array}{l} E_x [U(X_N^\pi)] \longrightarrow \max \\ X_N^\pi \geq 0 \\ \pi = (\pi_n) \text{ portfolio-strategy} \end{array} \right.$$

- $E' := \mathbb{R}_+ \times \mathbb{N}_0^2 \ni (x, (m, n)) = (x, \rho)$
- $A = \mathbb{R}, \quad D(x) = \{a \in \mathbb{R} \mid (1+i)x + a(Y - i - 1) \geq 0 \text{ a.s.}\}$
- Bayes-Operator $\Phi(\rho, \mathbf{u}) = (m+1, n)$
- $r' \equiv 0, \quad g'(x, \rho) := U(x)$

$b(x, \rho) := 1 + x$ is a bounding function for the filtered MDP.

Then it holds:

- $J_N(x) = J'_N(x, Q_0)$ is strictly increasing and concave in x .
- There exists an optimal policy $(f_0^*, f_1^*, \dots, f_{N-1}^*)$ for (P) .

Application: $U(x) = \frac{1}{\gamma} x^\gamma$ (**power utility**) $\gamma < 1, \gamma \neq 0$

$$(i) \quad J_N(x, \rho) = J_N(x, m, n) = \frac{1}{\gamma} x^\gamma \cdot d_N(m, n).$$

$$(ii) \quad f_k^*(x, \rho) = f_k^*(x, m, n) = x \cdot \alpha_k(m, n).$$

monotonicity results: $(m, n) \leq (m', n') : \iff m \leq m', n \geq n'$

$$(iii) \quad 0 < \gamma < 1 : \quad \alpha_k(m, n) \geq \alpha_k(\bar{p}) \text{ with } \bar{p} := \frac{m+1}{m+n+2}$$

$$\gamma < 0 : \quad \alpha_k(m, n) \leq \alpha_k(\bar{p})$$

Piecewise Deterministic Markov Decision Processes

- E state space, $E \subset \mathbb{R}^d$

- \mathbb{U} control space

$A := \{ \alpha : \mathbb{R}_+ \longrightarrow \mathbb{U} \text{ measurable} \}$, we write: $\alpha(t) = \alpha_t$

- $\mu(x, u)$ drift between jumps

$\phi_t^\alpha(x)$ (unique) solution of : $dx_t = \mu(x_t, \alpha_t)dt, x_0 = x$

deterministic flow between jumps

- $\lambda > 0$ jump rate (here: λ is independent of (x, u))

$0 := T_0 < T_1 < T_2 < \dots$ jump time points of a Poisson process with rate λ

- $Q(\cdot|x, u)$ distribution of jump goals

- $r(x, u)$ reward rate

- $\beta \geq 0$ discount rate

$\pi = (\pi_t)$ is called a **Markovian policy** (or piecewise open loop policy) if there exists a sequence of measurable functions $f_n : E \rightarrow A$ such that

$$\pi_t = f_n(Z_n)(t - T_n) \text{ for } T_n < t \leq T_{n+1}.$$

We write: $\pi = (\pi_t) = (f_n)$.

piecewise deterministic Markov process

$$X_t = \phi_{t-T_n}^\pi(Z_n) \text{ for } T_n \leq t < T_{n+1}, \quad Z_n = X_{T_n}$$

$$V_\pi(x) := E_x^\pi \left[\int_0^\infty e^{-\beta t} r(X_t, \pi_t) dt \right]$$

$$V_\infty(x) := \sup_\pi V_\pi(x), \quad x \in E$$

- Continuous-time stochastic control: Hamilton-Jacobi-Bellman equation
- Solution via discrete-time MDP

Yuskevich (1987), Davis (1993), Schäl et al. (2004)...

Jacobsen (2006), Guo/Hernandez-Lerma (2009): CTMDP

Discrete-time MDP

- E state space (embedded Markov process)
- A action space
- $Q'(B|x, \alpha) := \lambda \int_0^{\infty} e^{-(\beta+\lambda)t} Q(B|\phi_t^\alpha(x), \alpha_t) dt, B \subset E$
- $r'(x, \alpha) := \int_0^{\infty} e^{-(\beta+\lambda)t} r(\phi_t^\alpha(x), \alpha_t) dt$
- $\beta' = 1$

Note: A is a function space, Q' is substochastic.

$$(Tv)(x) = \sup_{\alpha \in A} \left\{ \int_0^{\infty} e^{-(\beta+\lambda)t} [r(\phi_t^\alpha(x), \alpha_t) + \lambda \int v(z) Q(dz|\phi_t^\alpha(x), \alpha_t)] dt \right\}$$

Theorem:

$$V_\pi(x) = E_x^\pi \left[\sum_{n=0}^{\infty} r'(Z'_n, f_n(Z'_n)) \right] =: J_{\infty\pi}(x)$$

$$V_\infty(x) = \sup_{\pi} J_{\infty\pi}(x) = J_\infty(x), x \in E$$

For a proof of the following result we use the set $\mathcal{R} := \{\alpha : \mathbb{R}_+ \longrightarrow \mathbb{P}(\mathbb{U}) \text{ measurable}\}$ of **relaxed controls** (with the Young topology). Since $\mathcal{R} \supset A$, we have to extend the domain of the data Q' and r' . Then it holds:

$$J_\infty^{\text{rel}}(x) \geq J_\infty(x) = V_\infty(x), \quad x \in E.$$

$b : E \longrightarrow \mathbb{R}_+$ is called an **upper bounding** function for the Piecewise Deterministic Markov Model, if there exist $c_r, c_Q, c_\phi \in \mathbb{R}_+$ such that

$$(i) \quad r^+(x, u) \leq c_r b(x).$$

$$(ii) \quad \int b(x') Q(dx'|x, u) \leq c_Q b(x).$$

$$(iii) \quad \lambda \int_0^\infty e^{-(\lambda+\beta)t} b(\phi_t^\alpha(x)) dt \leq c_\phi b(x).$$

If r is bounded from above, then $b \equiv 1$ is an upper bounding function and $c_Q = 1$ and $c_\phi = \frac{\lambda}{\lambda+\beta}$.

If b is an upper bounding function, then b is an upper bounding function for the MDP' (with and without relaxed controls) and $\alpha_b \leq c_Q c_\phi$.

Theorem: Suppose the Piecewise Deterministic Markov Model has a continuous upper bounding function b with $\alpha_b < 1$ and it holds:

- (i) \mathbb{U} is compact.
- (ii) $(t, x, \alpha) \longrightarrow \phi_t^\alpha(x)$ is continuous.
- (iii) $(x, u) \longrightarrow \int v(z)Q(dz|x, u)$ is usc for all usc $v \in \mathbb{B}_b^+$
- (iv) $(x, u) \longrightarrow r(x, u)$ is usc.

Then it holds:

- a) J_∞^{rel} is upper semi-continuous and $J_\infty^{\text{rel}} = T J_\infty^{\text{rel}}$.
- b) There exists an optimal relaxed policy $\pi^* = (\pi_t^*)$, i.e. π_t^* takes values in $\mathbb{P}(\mathbb{U})$.
- c) If $\phi_t^\alpha(x)$ is independent of α or if \mathbb{U} is convex, $\mu(x, u)$ is linear in u and $u \longrightarrow [r(x, u) + \lambda \int J_\infty^{\text{rel}}(z)Q(dz|x, u)]$ is concave on \mathbb{U} , then there exists an optimal **nonrelaxed** policy $\pi^* = (\pi_t^*)$ such that

$$\pi_t^* = f(X_{T_n}^{\pi^*})(t - T_n), \quad T_n < t \leq T_{n+1} \text{ for a decision rule } f : E \rightarrow A.$$

In particular, π_t^* takes values in \mathbb{U} and $J_\infty^{\text{rel}} = J_\infty = V_\infty$.

Continuous-Time Markov Decision Processes

- CTMDP (X_t) with countable state space E_X and intensities $q_{ij}(u)$
uniformized CTMDP : $\lambda \geq \sum_{j \neq i} q_{ij}(u) = -q_{ii}(u)$, $i \in E_X$, $u \in \mathbb{U}$

- Partially Observable CTMDP:

intensities depend on CTMC (Y_t) with finite state space E_Y , i.e.

$$q_{ij}(y, u) \text{ if } Y_t = y \in E_Y, y \text{ **unobservable state**}$$

Q_0 initial distribution of Y_0

$$\lambda \geq \sum_{j \neq i} q_{ij}(y, u) = -q_{ii}(y, u), i \in E_X, y \in E_Y, u \in \mathbb{U}$$

$$V_\pi(i) := E_i^\pi \left[\int_0^\infty e^{-\beta t} r(X_t^\pi, Y_t, \pi_t) dt \right]$$

$$V_\infty(i) := \sup_{\pi} V_\pi(i), i \in E_X$$

Filter Equation $\mu_t := P_i^\pi(Y_t = \cdot | \mathcal{F}_t^X) \in \mathbb{P}(E_Y)$

$$d\mu_t = b(X_t, \mu_t, \pi_t)dt + H(X_{t-}, \mu_{t-}, X_t, \pi_{t-}), \mu_0 = Q_0$$

Reformulation as filtered PDMDP:

(X_t, μ_t) piecewise deterministic MDP with state space $E_X \times \mathbb{P}(E_Y)$.

Theorem.

a) $V_\pi(i) = J_{\infty\pi}(i, Q_0)$ and $V_\infty(i) = J_\infty(i, Q_0)$, $i \in E_X$.

b) Assumptions! Then the Bellman equation holds, i.e.

$$J_\infty(i, \rho) = \sup_{\alpha \in A} \left\{ \int_0^\infty e^{-(\beta+\lambda)t} (LJ_\infty)(i, \phi_t^\alpha(i, \rho), \alpha_t) dt \right\}$$

where

$$(LJ_\infty)(i, \rho, u) := r(i, \rho, u) + \sum_{j \neq i} (J_\infty(j, \rho + H(i, \rho, j, u)) - J_\infty(i, \rho)) q_{ij}(\rho, u) + \lambda J_\infty(i, \rho)$$

Application: Parallel Queueing Model

two parallel queues and one server

Aim: minimize the expected number of waiting customers

- complete information: μC -rule is optimal
- partial information: $Y_t \equiv Y \in \{\mu_1, \nu_1\} \times \{\mu_2, \nu_2\}$
e.g. two types of customers are in the system and the server can not differ which group is waiting in which queue.

There exists an optimal nonrelaxed policy $f^*(i, \rho) \in \{1, 2\}$.

symmetric case: $Y \in \{(\mu_1, \mu_2), (\mu_2, \mu_1)\}$, $\mu_1 < \mu_2$

$$\mu_t = P_i^\pi(Y = (\mu_1, \mu_2) | \mathcal{F}_t^X) \implies d\mu_t = (\mu_2 - \mu_1)(2\pi_t - 1)\mu_t(1 - \mu_t)dt + \Delta\mu_t$$

$$\Delta\mu_t = \begin{cases} H_1(\mu_{t-}) & \text{if } X_t^1 = X_{t-}^1 - 1 \\ H_2(\mu_{t-}) & \text{if } X_t^2 = X_{t-}^2 - 1 \end{cases} \quad \text{where}$$

$$H_1(\rho) := \frac{\mu_1 \rho}{\mu_1 \rho + \mu_2 (1 - \rho)} - \rho, \quad H_2(\rho) := \frac{\mu_2 \rho}{\mu_2 \rho + \mu_1 (1 - \rho)} - \rho, \quad \rho \in [0, 1]$$

It holds: $H_1(\rho) \leq 0$, $H_2(\rho) \geq 0$.

The stationary policy (f^*, f^*, \dots) is optimal with

$$f^*(i_1, i_2, \rho) = \begin{cases} 1 & i_2 = 0 \\ 2 & i_1 = 0 \\ 1 & \rho \leq \frac{1}{2}, (i_1, i_2) \in \mathbb{N} \times \mathbb{N} \\ 2 & \rho > \frac{1}{2}, (i_1, i_2) \in \mathbb{N} \times \mathbb{N} \end{cases}$$

$$\bar{\mu}_1 := \mu_1 \rho + \mu_2 (1 - \rho), \quad \bar{\mu}_2 := \mu_2 \rho + \mu_1 (1 - \rho)$$

$$\bar{\mu}_1 \geq \bar{\mu}_2 \iff \rho \leq \frac{1}{2}$$

certainty equivalence principle for the μc -rule holds (if $c_1 = c_2$)!

Rieder/Winter (2009), Bäuerle/Rieder (2009)

References

Bäuerle/Rieder (2011) : Markov decision processes with applications to finance.

Bertsekas/Shreve (1978): Stochastic optimal control.

Feinberg/Schwartz (2002): Handbook of Markov decision processes.

Hernandez-Lerma/Lasserre (1996): Discrete-time Markov control processes.

Hinderer(1970): Foundations of non-stationary dynamic programming.

Puterman (1994): Markov decision processes.

Davis (1993): Markov models and optimization.

Elliott/Aggoun/Moore (1995): Hidden Markov models.

Gittins (1989): Bandit processes and dynamic allocation processes.

Guo/Hernandez-Lerma (2009): Continuous-time Markov decision processes.

Jacobsen (2006): Point process theory and applications.

Jeanblanc/Yor/Chesney (2009): Mathematical methods for financial markets.

Rieder/Winter (2009): Optimal control of Markovian jump processes with partial information and applications to a parallel queueing model.
Math. Meth. Operat. Res.70, 567-596.

Bäuerle/Rieder (2009): MDP algorithms for portfolio optimization problems in pure jump markets.
Finance Stoch. 13, 591-611.