## Decisions Under Uncertainty: Reinforcement Learning for Impulse Control Piecewise Deterministic Markov Processes

Orlane Rossini <sup>1</sup>, Alice Cleynen <sup>1,2</sup>, Benoîte de Saporta <sup>1</sup>, Régis Sabbadin <sup>3</sup> and Meritxell Vinyals <sup>3</sup>

<sup>1</sup>IMAG, Univ Montpellier, CNRS, Montpellier, France <sup>2</sup>John Curtin School of Medical Research, The Australian National University, Canberra, ACT, Australia <sup>3</sup>Univ Toulouse, INRAE-MIAT, Toulouse, France

May 2025





### Medical context



FIGURE: Example of patient data<sup>*a*</sup>

- Patients who have had cancer benefit from regular follow-up;
- The concentration of clonal immunoglobulin is measured over time;
- The doctor has to make new decisions at each visit.



<sup>&</sup>lt;sup>a</sup>IUCT Oncopole and CRCT, Toulouse, France

### Medical context



FIGURE: Example of patient data<sup>*a*</sup>

- Patients who have had cancer benefit from regular follow-up;
- The concentration of clonal immunoglobulin is measured over time;
- The doctor has to make new decisions at each visit.

### ⇒ Optimising decision-making to ensure the patient's quality of life

<sup>&</sup>lt;sup>a</sup>IUCT Oncopole and CRCT, Toulouse, France

## Controlled PDMP<sup>1</sup>

### We switch randomly from one deterministic regime to another.



<sup>1</sup>Piecewise Deterministic Markov Processes

## Controlled PDMP<sup>1</sup>

We switch randomly from one deterministic regime to another.



Let  $x = (m, \ell, k, \zeta, u)$  the patient's condition:

- *m* the patient's condition;
- $\ell$  the current treatment;
- *k* the number of treatments;
- $\zeta$  the biomarker;
- *u* the time since the last jump.

### A PDMP is defined by three local characteristics.



#### FLOW

Description of the deterministic part of the process.

$$\Phi^{\ell}(x,t) = (m,k,\ell,\phi^{\ell}_{m,k}(\zeta,t),u+t)$$

## Local Characteristics of a PDMP<sup>2</sup>

### A PDMP is defined by three local characteristics.



#### UMP INTENSITY

Description of the process jump mechanisms.

• Boundary jump (deterministic)

$$t^*(\mathbf{x}) = t_{m,k}^{\ell \star}(\zeta) = \inf\{t > \mathsf{o} : \phi_{m,k}^{\ell}(\zeta, t) \in \{\zeta_{\mathsf{o}}, D\}\}$$

• Random jump

$$\mathbb{P}(T > t) = e^{-\int_0^t \lambda_{m,k}^\ell(\Phi(x,s)) \, \mathrm{d}s}$$

## Local Characteristics of a PDMP<sup>2</sup>

A PDMP is defined by three local characteristics.



#### Markov kernel

Description of the state of the process after each jump.

$$\mathbb{P}(X' \in A | X = x) = \int_A Q^d_{m,k}(\Phi^{\ell}(x,T), \mathrm{d}x')$$

<sup>2</sup>Piecewise Deterministic Markov Processes

Identify an  $\epsilon$ -optimal strategy  $S = (\tau_n, \chi_n)_{n \geq 1}$ 

$$\underbrace{\mathcal{V}(\mathcal{S}, \mathbf{X})}_{\text{Expected cost of strategy}\mathcal{S}} = \mathbb{E}_{\mathbf{X}}^{\mathcal{S}} \left[ \int_{0}^{+\infty} e^{-\gamma t} \underbrace{c_{\mathbf{R}}(X_{t})}_{\text{current trajectory cost}} dt + \sum_{n=1}^{\infty} \underbrace{c_{\mathbf{I}}}_{\text{impulse cost}} (X_{\tau_{n}}, X_{\tau_{n}^{+}}) \right],$$

<sup>&</sup>lt;sup>3</sup>Piecewise Deterministic Markov Processes

## Solving impulse control for PDMP<sup>3</sup>

Identify an  $\epsilon$ -optimal strategy  $S = (\tau_n, \chi_n)_{n \geq 1}$ 



$$\mathcal{V}^{\star}(\mathbf{X}) = \inf_{\mathcal{S} \in \mathsf{S}} \mathcal{V}(\mathcal{S}, \mathbf{X})$$

<sup>&</sup>lt;sup>3</sup>Piecewise Deterministic Markov Processes

### Difficulties

### Partially known dynamics



Hypothesis:  $v_1 \sim$  Log-Normal  $(\mu, \sigma^{-2})$ , with  $\mu$  and  $\sigma$  unknown.

### **Partial observation**



### Methods



<sup>4</sup>Piecewise Deterministic Markov Processes

### Methods



<sup>5</sup>Piecewise Deterministic Markov Processes <sup>6</sup>Bayes-Adaptive Partially Observed Markov Decision Process



# Agent



#### MDP DEFINITION

A MDP is defined by a tuple ( $\mathbb{S}$ ,  $\mathbb{A}$ , P, c).

- Patient condition  $s = (m, k, \zeta, u) \in S$ ;
- Actions  $a = (\ell, r) \in \mathbb{A}$ ;
- Transition function *P*(s'|s, a);
- Cost function  $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$ .





#### MDP DEFINITION

A MDP is defined by a tuple ( $\mathbb{S}$ ,  $\mathbb{A}$ , P, c).

- Patient condition  $s = (m, k, \zeta, u) \in S$ ;
- Actions  $a = (\ell, r) \in \mathbb{A}$ ;
- Transition function *P*(s'|s, a);
- Cost function  $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$ .



#### MDP DEFINITION

A MDP is defined by a tuple ( $\mathbb{S}$ ,  $\mathbb{A}$ , P, c).

- Patient condition  $s = (m, k, \zeta, u) \in S$ ;
- Actions  $a = (\ell, r) \in \mathbb{A}$ ;
- Transition function *P*(*s*'|*s*, *a*);
- Cost function  $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$ .



<sup>8</sup>Partially Observed Markov Decision Process



<sup>8</sup>Partially Observed Markov Decision Process



<sup>8</sup>Partially Observed Markov Decision Process

#### POMDP DEFINITION

A POMDP is defined by a tuple (S, A, P,  $\Omega$ , Z, c).

- Patient condition  $s = (m, k, \zeta, u) \in S$ ;
- Actions  $a = (\ell, r) \in \mathbb{A}$ ;
- Transition function P(s'|s, a);
- **Observation**  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- **Observation function**  $Z(\omega|s)$ ;
- Cost function  $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$ .



POMDP DEFINITION

A POMDP is defined by a tuple ( $\mathbb{S}$ ,  $\mathbb{A}$ , P,  $\Omega$ , Z, c).

- Patient condition  $s = (m, k, \zeta, u) \in S$ ;
- Transition function P(s'|s, a);
- **Observation**  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega;$
- **Observation function**  $Z(\omega|s)$ ;
- Cost function  $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$ .

<sup>8</sup>Partially Observed Markov Decision Process









### Generate transition from prior





<sup>9</sup>Bayes Adaptive Partially observed Markov decision process



<sup>9</sup>Bayes Adaptive Partially observed Markov decision process



#### **BAPOMDP** DEFINITION

Un BAPOMDP se définit par un tuple ( $\mathbb{S}^+$ ,  $\mathbb{A}$ ,  $P^+$ ,  $\Omega$ , Z, c).

- Space of hyperstate  $\mathbb{S}^+ = \mathbb{S} \times \Theta$ ;
- Actions  $a = (\ell, r) \in \mathbb{A}$ ;
- **Transition function**  $P^+(s', \theta'|s, a, \theta)$ ;
- Observation  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega;$
- Observation function  $Z(\omega|s)$ ;
- Cost function  $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$ .

<sup>9</sup>Bayes Adaptive Partially observed Markov decision process



#### BAPOMDP DEFINITION

Un BAPOMDP se définit par un tuple ( $\mathbb{S}^+$ ,  $\mathbb{A}$ ,  $P^+$ ,  $\Omega$ , Z, c).

- Space of hyperstate  $\mathbb{S}^+ = \mathbb{S} \times \Theta$ ;
- Actions  $a = (\ell, r) \in \mathbb{A}$ ;
- **Transition function**  $P^+(s', \theta'|s, a, \theta)$ ;
- Observation  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega;$
- Observation function  $Z(\omega|s)$ ;
- Cost function  $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$ .

<sup>&</sup>lt;sup>9</sup>Bayes Adaptive Partially observed Markov decision process



<sup>&</sup>lt;sup>10</sup>Bayes Adaptative Partially Observable Markov Decision Process



<sup>&</sup>lt;sup>10</sup>Bayes Adaptative Partially Observable Markov Decision Process



<sup>10</sup>Bayes Adaptative Partially Observable Markov Decision Process

### In reality, we do not observe state space!

Let  $h_n = (\omega_0, a_0, \omega_1, a_1, \dots, \omega_n)$  be the history



<sup>10</sup>Bayes Adaptative Partially Observable Markov Decision Process

## **Reinforcement Learning**



The optimal policy is obtained from the experiments  $< \omega, a, \omega', c >$ , generate from  $P^+$  transition function

$$\underbrace{\pi^{\star}}_{Q \text{ function}} = \arg\min_{a \in \mathbb{A}} Q^{\star}(s, a)$$



<sup>11</sup>Deep Q-Network

## Conclusion and future work



<sup>12</sup>Piecewise Deterministic Markov Processes
<sup>13</sup>Partially Observed Markov Decision Process
<sup>14</sup>Bayes Adaptative Partially Observed Markov Decision Process