

# Deep reinforcement learning for controlled Piecewise Deterministic Markov Process

Application to optimising medical treatment

Orlane Rossini <sup>1</sup>, Alice Cleyen <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

28 Mai 2024



UNIVERSITÉ DE  
MONTPELLIER

INRAE

IMAG  
INSTITUT MONTPELLIERAIN  
ALEXANDER GROTHENDIECK



anr<sup>®</sup>

# Medical context

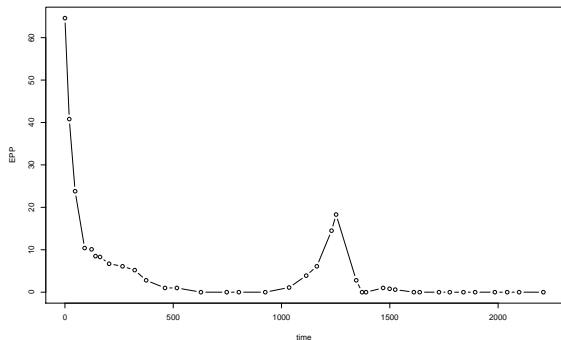


FIGURE: Example of patient data<sup>a</sup>

- Patients who have had **cancer** benefit from **regular monitoring**;
- the concentration of **cancer cells** is measured **over time**;
- The doctor has to make new **decisions** at each visit.

<sup>a</sup>IUCT Oncopole et CRCT, Toulouse, France

# Medical context

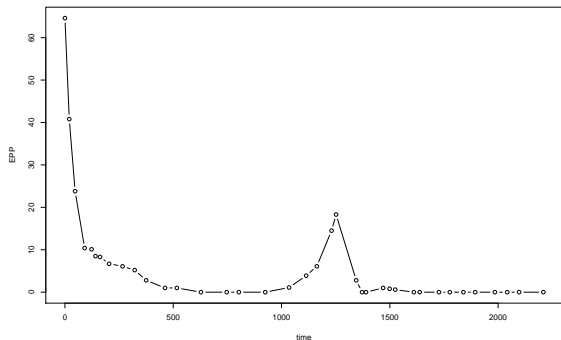


FIGURE: Example of patient data<sup>a</sup>

- Patients who have had **cancer** benefit from **regular monitoring**;
- the concentration of **cancer cells** 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>a</sup>IUCT Oncopole et CRCT, Toulouse, France

**A real-life problem**

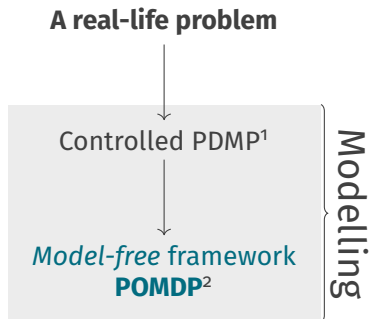


Controlled PDMP<sup>1</sup>

---

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

# Methods

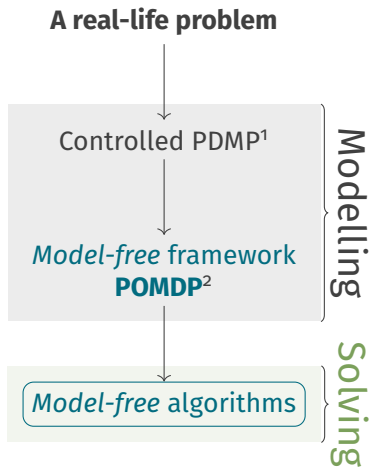


---

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

<sup>2</sup>Partially Observable Markov Decision Process

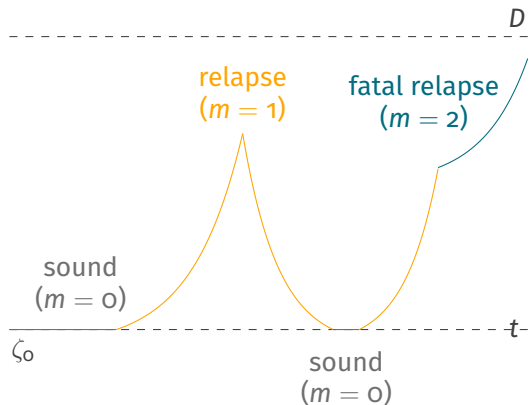
# Methods



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

<sup>2</sup>Partially Observable Markov Decision Process

# Model-free framework<sup>3</sup>



Let the **patient's state**  $s = (m, k, \zeta, u, t, \tau)$ :

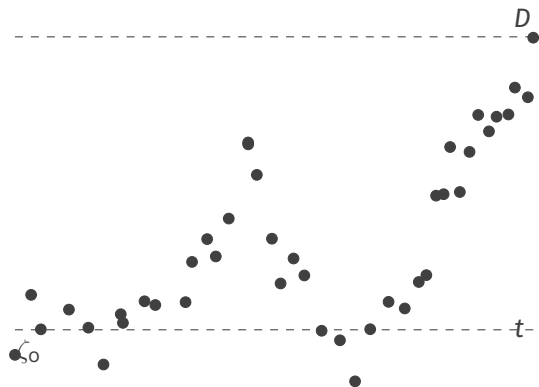
- $m$  the overall state of the patient;
- $k$  relapses number;
- $\zeta$  the biomarker;
- $u$  time since last jump;
- $t$  time since the beginning of follow-up;
- $\tau$  time since a treatment is applied.

Let  $d$  be the **decision** such that:  $d = (\ell, r)$ :

- $\ell$  the treatment (*nothing, chemotherapy, palliative cares*);
- $r$  time before the next visit (15, 30, 60 days).

<sup>3</sup>Partially Observable Markov Decision Process (POMDP)

# Model-free framework<sup>3</sup>



Let the **patient's state**  $s = (m, k, \zeta, u, t, \tau)$ :

- $m$  the overall state of the patient;
- $k$  relapses number;
- $\zeta$  the biomarker;
- $u$  time since last jump;
- $t$  time since the beginning of follow-up;
- $\tau$  time since a treatment is applied.

Let  $d$  be the **decision** such that:  $d = (\ell, r)$ :

- $\ell$  the treatment (*nothing, chemotherapy, palliatif cares*);
- $r$  time before the next visit (15, 30, 60 days).

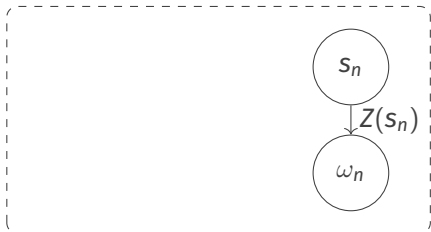
<sup>3</sup>Partially Observable Markov Decision Process (POMDP)



# POMDP<sup>4</sup> Characteristics

Agent

Environment



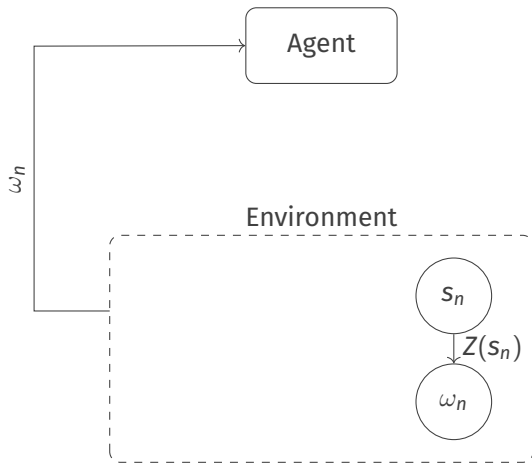
## POMDP DEFINITION

A POMDP is defined by a tuple  $(\mathcal{S}, \mathcal{D}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Patient's state  $s = (m, k, \zeta, u, t, \tau) \in \mathcal{S}$ ;
- Decisions  $d = (\ell, r) \in \mathcal{D}$ ;
- Transition probabilities  $\mathcal{P}(s, d)(s')$ ;
- Observations  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- The observations function  $\mathcal{Z}(s)(\omega)$ ;
- Cost function  $C(s, d, s')$ .

<sup>4</sup>Partially Observable Markov Decision Process

# POMDP<sup>4</sup> Characteristics



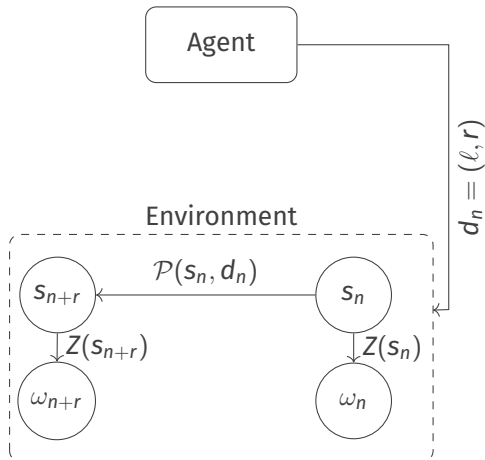
## POMDP DEFINITION

A POMDP is defined by a tuple  $(\mathcal{S}, \mathcal{D}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Patient's state  $s = (m, k, \zeta, u, t, \tau) \in \mathcal{S}$ ;
- Decisions  $d = (\ell, r) \in \mathcal{D}$ ;
- Transition probabilities  $\mathcal{P}(s, d)(s')$ ;
- Observations  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- The observations function  $\mathcal{Z}(s)(\omega)$ ;
- Cost function  $C(s, d, s')$ .

<sup>4</sup>Partially Observable Markov Decision Process

# POMDP<sup>4</sup> Characteristics



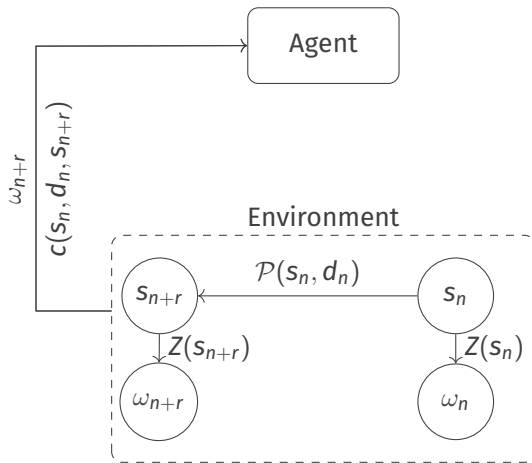
## POMDP DEFINITION

A POMDP is defined by a tuple  $(\mathcal{S}, \mathcal{D}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Patient's state  $s = (m, k, \zeta, u, t, \tau) \in \mathcal{S}$ ;
- Decisions  $d = (\ell, r) \in \mathcal{D}$ ;
- Transition probabilities  $\mathcal{P}(s, d)(s')$ ;
- Observations  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- The observations function  $\mathcal{Z}(s)(\omega)$ ;
- Cost function  $C(s, d, s')$ .

<sup>4</sup>Partially Observable Markov Decision Process

# POMDP<sup>4</sup> Characteristics



## POMDP DEFINITION

A POMDP is defined by a tuple  $(\mathcal{S}, \mathcal{D}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Patient's state  $s = (m, k, \zeta, u, t, \tau) \in \mathcal{S}$ ;
- Decisions  $d = (\ell, r) \in \mathcal{D}$ ;
- Transition probabilities  $\mathcal{P}(s, d)(s')$ ;
- Observations  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- The observations function  $\mathcal{Z}(s)(\omega)$ ;
- Cost function  $C(s, d, s')$ .

## POMDP CHARACTERISTICS

- Discrete observation dates
- Continuous state space
- Partially observable
- Partially known model
- Simulator

<sup>4</sup>Partially Observable Markov Decision Process

# What does it mean *to solve* ?

## Identifying the optimal policy!

$$\underbrace{C(s, d, s')}_{\text{Cost function}} = \underbrace{C_V}_{\text{Visit cost}} + \underbrace{C_D(H - t') \times \left(1 - \frac{1}{2} \mathbb{1}_{m=2 \text{ and } \ell=b}\right) \mathbb{1}_{m'=3}}_{\text{Death cost}} + \underbrace{\beta |\zeta - \zeta_{TH}| \times \left(1 - \frac{1}{2} \mathbb{1}_{\ell=a \text{ or } \ell=b}\right) \mathbb{1}_{\zeta > \zeta_{TH}}}_{\text{Disease cost}} + \underbrace{\kappa_C \times r \left(1 - \frac{1}{2} \mathbb{1}_{m=1 \text{ and } \ell=a}\right) \mathbb{1}_{\ell=a}}_{\text{Chemotherapy cost}} + \underbrace{\kappa_P \times r \left(1 - \frac{1}{2} \mathbb{1}_{m=2 \text{ and } \ell=b}\right) \mathbb{1}_{\ell=b}}_{\text{Palliatif care cost}}$$

# What does it mean *to solve* ?

## Identifying the optimal policy!

$$\underbrace{V(\pi, s)}_{\text{Criterion to optimise}} = \underbrace{\mathbb{E}_s^\pi \left[ \sum_{n=0}^{H-1} c(S_{n-1}, D_n, S_n) \right]}_{\text{Long-term expected cost following the policy } \pi}$$

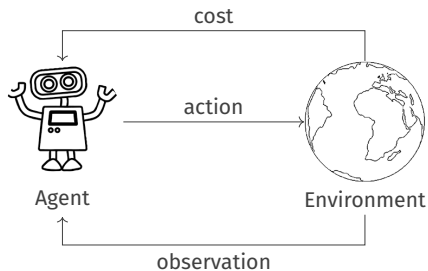
# What does it mean *to solve* ?

## Identifying the optimal policy!

$$\underbrace{V(\pi, s)}_{\text{Criterion to optimise}} = \underbrace{\mathbb{E}_s^\pi \left[ \sum_{n=0}^{H-1} c(S_{n-1}, D_n, S_n) \right]}_{\text{Long-term expected cost following the policy } \pi}$$

$$\underbrace{V^*(s)}_{\text{Value function}} = \underbrace{\min_{\pi \in \Pi} V(\pi, s)}_{\text{Minimising over the all set of policy } \Pi}$$

# Reinforcement learning



Learning optimal policy from all experiences

$$\underbrace{Q^\pi(s, d)}_{\text{Criterion to optimise}} = \underbrace{\mathbb{E}^\pi \left[ \sum_{n=0}^{H-1} c(S_{n-1}, D_n, S_n) \right]}_{\text{Value of an action in a state, following } \pi} | s, d = (\ell, r)$$

$$\underbrace{Q^*(s, d)}_{\text{Q function}} = \min_{\pi \in \Pi} Q^\pi(s, d)$$

$$\underbrace{\pi^*}_{\text{Q function}} = \arg \min_{d \in \mathcal{D}} Q^*(s, d)$$



# Model-free algorithms

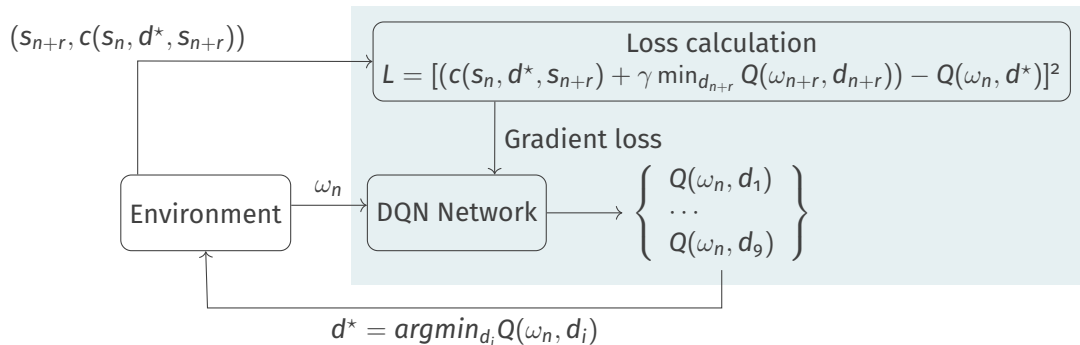


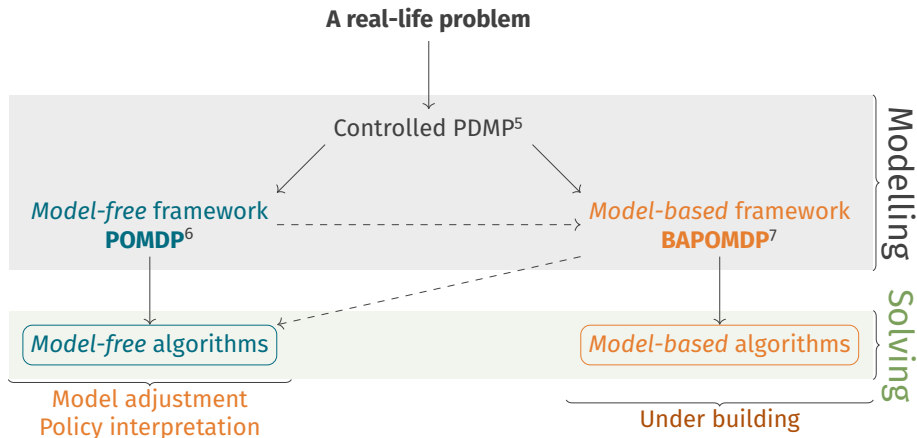
FIGURE: Deep Q-Network

# Results

Policy	Mean Cost	Confidence interval	Survival rate	Relapse rate
<b>OH</b>	119.16	[119.15, 119.17]	95.71%	2.57
<b>Random</b>	158.09	[158.08, 158.10]	18.32%	1.27
<b>Inactive</b>	164.09	[164.08, 164.10]	0.14%	1.00
<b>Threshold</b>	157.27	[157.26, 157.28]	80.52%	1.01
<b>DQN</b>	121.52	[121.51, 121.53]	99.83%	0.59

TABLE: Policy evaluation performance on  $10^5$  simulations

# Conclusion and Future Works



<sup>5</sup>Piecewise Deterministic Markov Process

<sup>6</sup>Partially Observable Markov Decision Process

<sup>7</sup>Bayes Adaptive Partially Observable Markov Decision Process