

# Apprentissage par renforcement pour le contrôle de processus Markovien déterministe par morceaux

Application à l'optimisation d'un traitement médical

Orlane Rossini <sup>1</sup>, Alice Cleyen <sup>1,2</sup>, Benoîte de Saporta <sup>1</sup> et Régis Sabbadin <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

June 2024



UNIVERSITÉ DE  
MONTPELLIER

INRAE

IMAG  
INSTITUT MONTPELLIERAIN  
ALEXANDER GROTHENDIECK



anr<sup>®</sup>

# Le contexte médical

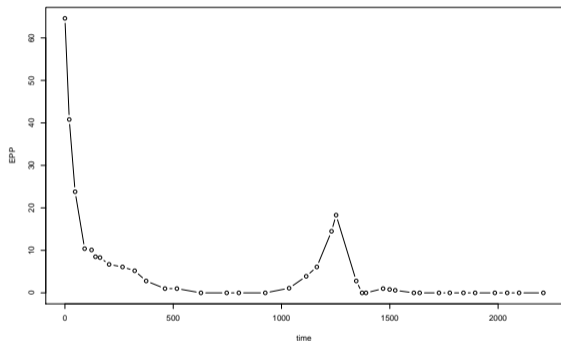


FIGURE: Exemple de donnée d'un patient<sup>a</sup>

- Des patients ayant eu un **cancer** bénéficiant d'un **suivi régulier**;
- La concentration d'**immunoglobuline clonale** est mesurée **dans le temps**;
- Le médecin doit prendre de nouvelles **décisions** à chaque visite.

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

# Le contexte médical

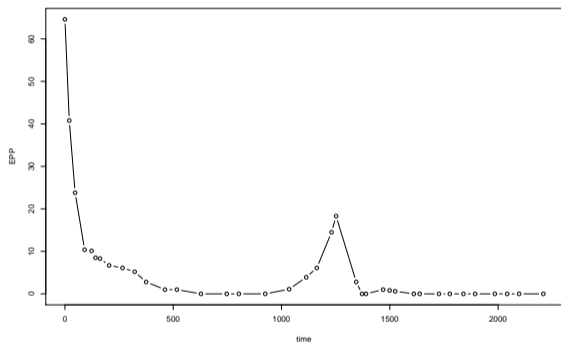
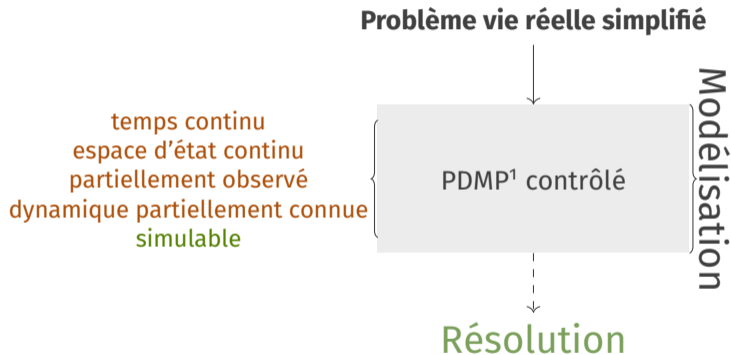


FIGURE: Exemple de donnée d'un patient<sup>a</sup>

- Des patients ayant eu un **cancer** bénéficient d'un **suivi régulier**;
- La concentration d'**immunoglobuline clonale** est mesurée **dans le temps**;
- Le médecin doit prendre de nouvelles **décisions** à chaque visite.

⇒ **Optimiser la prise de décision pour assurer la qualité de vie du patient**

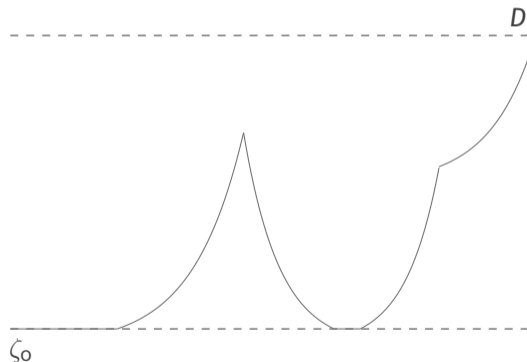
<sup>a</sup>IUCT Oncopole et CRCT, Toulouse, France



<sup>1</sup>Processus Markovien Déterministe par Morceaux

# Le modèle PDMP<sup>2</sup> contrôlé

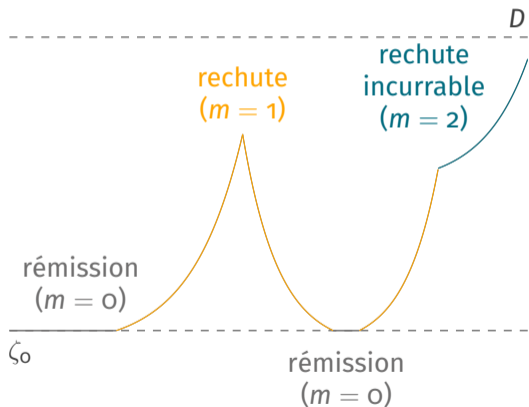
On passe aléatoirement d'un régime déterministe à un autre.



<sup>2</sup>Processus Markovien Déterministe par Morceaux

# Le modèle PDMP<sup>2</sup> contrôlé

On passe aléatoirement d'un régime déterministe à un autre.



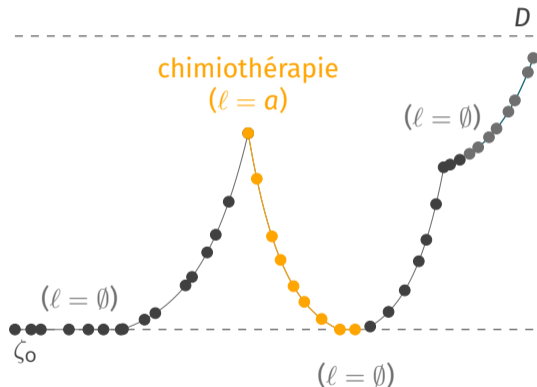
Soit l'état du patient  $x = (m, k, \zeta, u)$ :

- $m$  l'état du patient;
- $k$  le nombre de rechute;
- $\zeta$  le biomarqueur;
- $u$  le temps depuis le dernier saut.

<sup>2</sup>Processus Markovien Déterministe par Morceaux

# Le modèle PDMP<sup>2</sup> contrôlé

On passe aléatoirement d'un régime déterministe à un autre.



Soit l'état du patient  $x = (m, k, \zeta, u)$ :

- $m$  l'état du patient;
- $k$  le nombre de rechute;
- $\zeta$  le biomarqueur;
- $u$  le temps depuis le dernier saut.

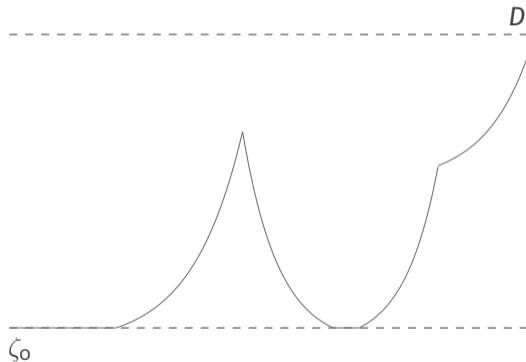
Soit  $d$  la décision telle que:  $d = (\ell, r)$ :

- $\ell$  le traitement;
- $r$  le temps avant la prochaine visite.

<sup>2</sup>Processus Markovien Déterministe par Morceaux

# Caractéristiques d'un PDMP<sup>3</sup>

Un PDMP se définit par trois caractéristiques locales.

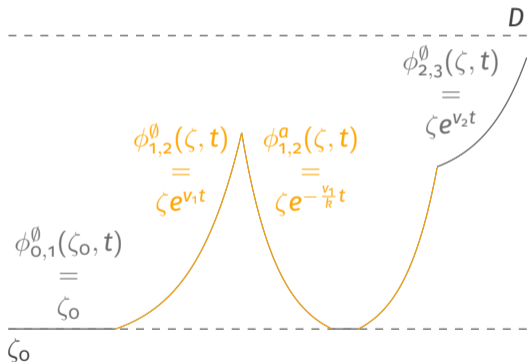


<sup>3</sup>Processus Markovien Déterministe par Morceaux



# Caractéristiques d'un PDMP<sup>3</sup>

Un PDMP se définit par trois caractéristiques locales.



## LE FLOT

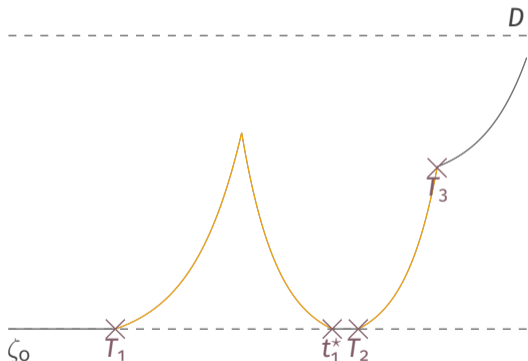
Description de la partie déterministe du processus.

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

<sup>3</sup>Processus Markovien Déterministe par Morceaux

# Caractéristiques d'un PDMP<sup>3</sup>

Un PDMP se définit par trois caractéristiques locales.



## L'INTENSITÉ DE SAUT

Description des mécanismes de saut du processus.

- Saut à la frontière (déterministe)

$$t^*(x) = t_m^{\ell^*}(\zeta) = \inf\{t > 0 : \phi_{m,k}^{\ell}(\zeta, t) \in \{\zeta_0, D\}\}$$

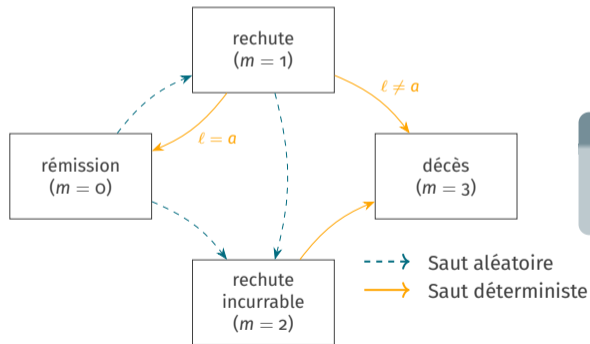
- Saut aléatoire

$$\mathbb{P}(T > t) = e^{-\int_0^t \lambda_m^{\ell}(\phi(x,s)) ds}$$

<sup>3</sup>Processus Markovien Déterministe par Morceaux

# Caractéristiques d'un PDMP<sup>3</sup>

Un PDMP se définit par trois caractéristiques locales.

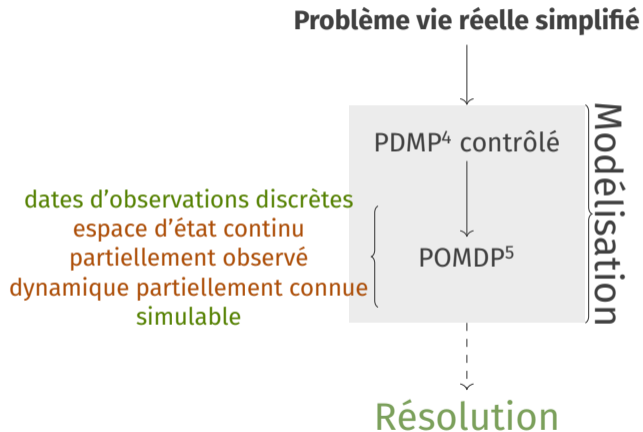


## LE NOYAU

Description de l'état du processus après chaque saut.

$$\mathbb{P}(X' \in A | X = x) = \int_A Q_m^d(\Phi^\ell(x, T), dx')$$

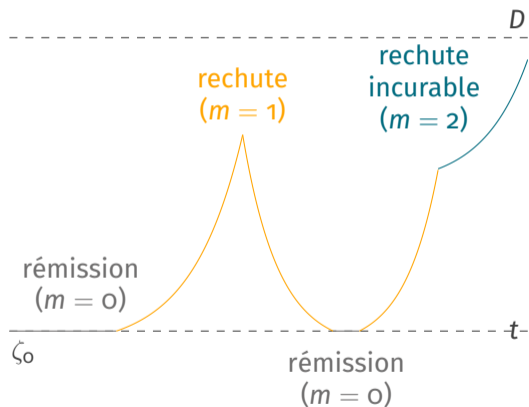
<sup>3</sup>Processus Markovien Déterministe par Morceaux



<sup>4</sup>Processus Markovien Déterministe par Morceaux

<sup>5</sup>Processus de Décision Markovien Partiellement Observé

# Le modèle POMDP<sup>6</sup>



Soit  $s = (m, k, \zeta, u, t, \tau)$  l'état du patient:

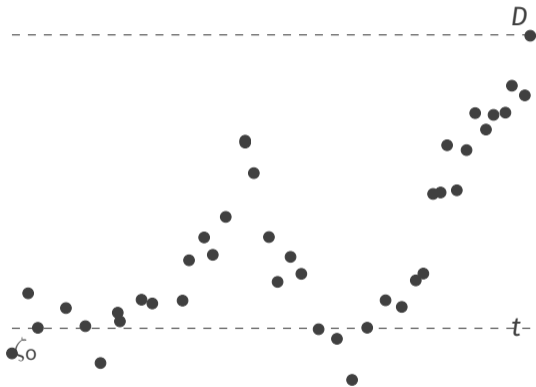
- $m$  état général du patient;
- $k$  nombre de rechute;
- $\zeta$  biomarqueur;
- $u$  temps depuis le dernier saut;
- $t$  temps écoulé depuis le début du suivi;
- $\tau$  temps depuis l'application d'un traitement.

Soit  $d$  la **décision** telle que:  $d = (\ell, r)$ :

- $\ell$  traitement (*rien, chimiothérapie*);
- $r$  temps avant la prochaine visite (*15, 30, 60 jours*).

<sup>6</sup>Processus de Décision Markovien Partiellement Observé

# Le modèle POMDP<sup>6</sup>



Soit  $s = (m, k, \zeta, u, t, \tau)$  l'état du patient:

- $m$  état général du patient;
- $k$  nombre de rechute;
- $\zeta$  biomarqueur;
- $u$  temps depuis le dernier saut;
- $t$  temps écoulé depuis le début du suivi;
- $\tau$  temps depuis l'application d'un traitement.

Soit  $d$  la décision telle que:  $d = (\ell, r)$ :

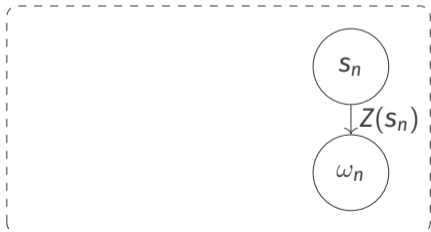
- $\ell$  traitement (*rien, chimiothérapie*);
- $r$  temps avant la prochaine visite (*15, 30, 60 jours*).

<sup>6</sup>Processus de Décision Markovien Partiellement Observé

# Caractéristiques d'un POMDP<sup>7</sup>

Agent

Environnement



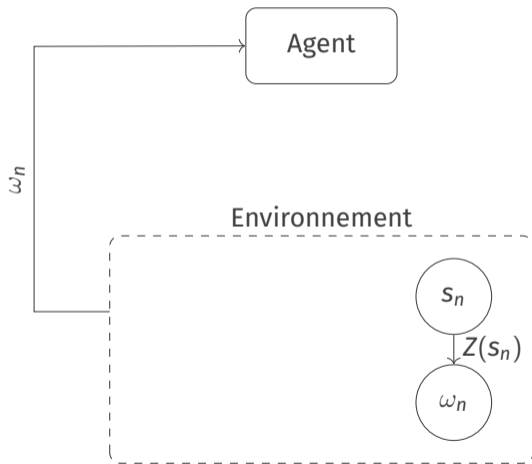
## POMDP DEFINITION

Un POMDP se définit par un tuple  $(S, \mathcal{D}, \mathcal{K}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Etat du patient  $s = (m, k, \zeta, u, t, \tau) \in S$ ;
- Décisions  $d = (\ell, r) \in \mathcal{D}$ ;
- $\mathcal{K}(s) \subseteq \mathcal{D}$  l'espace des décisions admissibles dans l'état  $s$ ;
- Probabilité de transition  $\mathcal{P}(s, d)(s')$ ;
- Observation  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Fonction d'observation  $\mathcal{Z}(s)(\omega)$ ;
- Fonction coût  $C(s, d, s')$ .

<sup>7</sup>Processus de Décision Markovien Partiellement Observé

# Caractéristiques d'un POMDP<sup>7</sup>



## POMDP DEFINITION

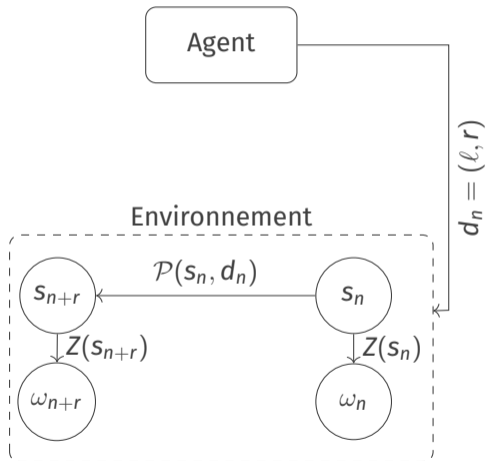
Un POMDP se définit par un tuple  $(S, \mathcal{D}, \mathcal{K}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Etat du patient  $s = (m, k, \zeta, u, t, \tau) \in S$ ;
- Décisions  $d = (\ell, r) \in \mathcal{D}$ ;
- $\mathcal{K}(s) \subseteq \mathcal{D}$  l'espace des décisions admissibles dans l'état  $s$ ;
- Probabilité de transition  $\mathcal{P}(s, d)(s')$ ;
- Observation  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Fonction d'observation  $\mathcal{Z}(s)(\omega)$ ;
- Fonction coût  $C(s, d, s')$ .

<sup>7</sup>Processus de Décision Markovien Partiellement Observé



# Caractéristiques d'un POMDP<sup>7</sup>



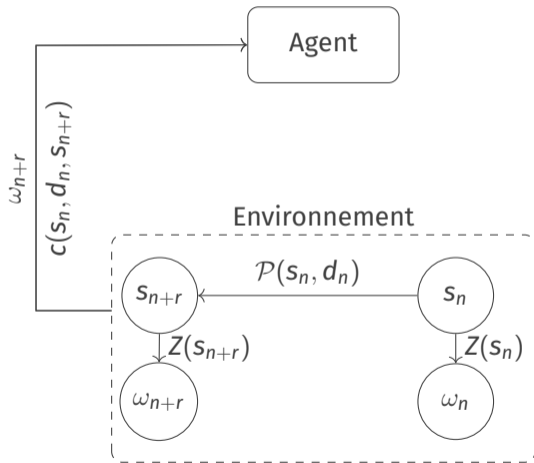
## POMDP DEFINITION

Un POMDP se définit par un tuple  $(S, \mathcal{D}, \mathcal{K}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Etat du patient  $s = (m, k, \zeta, u, t, \tau) \in S$ ;
- Décisions  $d = (\ell, r) \in \mathcal{D}$ ;
- $\mathcal{K}(s) \subseteq \mathcal{D}$  l'espace des décisions admissibles dans l'état  $s$ ;
- Probabilité de transition  $\mathcal{P}(s, d)(s')$ ;
- Observation  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Fonction d'observation  $\mathcal{Z}(s)(\omega)$ ;
- Fonction coût  $C(s, d, s')$ .

<sup>7</sup>Processus de Décision Markovien Partiellement Observé

# Caractéristiques d'un POMDP<sup>7</sup>



## POMDP DEFINITION

Un POMDP se définit par un tuple  $(S, \mathcal{D}, \mathcal{K}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Etat du patient  $s = (m, k, \zeta, u, t, \tau) \in S$ ;
- Décisions  $d = (\ell, r) \in \mathcal{D}$ ;
- $\mathcal{K}(s) \subseteq \mathcal{D}$  l'espace des décisions admissibles dans l'état  $s$ ;
- Probabilité de transition  $\mathcal{P}(s, d)(s')$ ;
- Observation  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Fonction d'observation  $\mathcal{Z}(s)(\omega)$ ;
- Fonction coût  $C(s, d, s')$ .

## Identifier une politique optimale!

$$\underbrace{C(s, d, s')}_{\text{Fonction de coût}} = \underbrace{C_V}_{\text{coût de la visite}} + \underbrace{C_D(H - t') \times \mathbb{1}_{m'=3}}_{\text{coût de la mort}} + \underbrace{\kappa_C \times r \times \mathbb{1}_{\ell=a}}_{\text{coût de la chimiothérapie}}$$

---

<sup>8</sup>Processus de Décision Markovien Partiellement Observé

## Identifier une politique optimale!

$$\underbrace{V(\pi, s)}_{\text{Critère à optimiser}} = \underbrace{\mathbb{E}_s^\pi \left[ \sum_{n=0}^{H-1} c(S_{n-1}, D_n, S_n) \right]}_{\text{Coût attendu à long terme suite à la politique menée } \pi}$$

---

<sup>8</sup>Processus de Décision Markovien Partiellement Observé

# Résoudre un POMDP<sup>8</sup>

## Identifier une politique optimale!

$$\underbrace{V(\pi, s)}_{\text{Critère à optimiser}} = \underbrace{\mathbb{E}_s^\pi \left[ \sum_{n=0}^{H-1} c(S_{n-1}, D_n, S_n) \right]}_{\text{Coût attendu à long terme suite à la politique menée } \pi}$$

$$\underbrace{V^*(s)}_{\text{Fonction valeur}} = \underbrace{\min_{\pi \in \Pi} V(\pi, s)}_{\text{Minimisation sur l'ensemble des politiques } \Pi.}$$

<sup>8</sup>Processus de Décision Markovien Partiellement Observé

## Identifier une politique optimale!

En réalité on observe pas l'espace d'état !

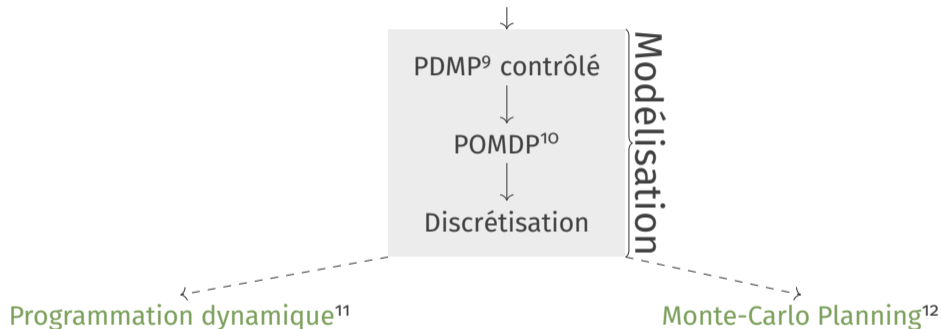
Soit l'historique  $h = (\omega_0, d_0, \omega_1, d_1, \dots, \omega_n)$

$$\underbrace{V^*(h)}_{\text{Fonction valeur}} = \underbrace{\min_{\pi \in \Pi} V(\pi, h)}_{\text{Minimisation sur l'ensemble des politiques } \Pi}.$$

---

<sup>8</sup>Processus de Décision Markovien Partiellement Observé

## Problème vie réelle simplifié

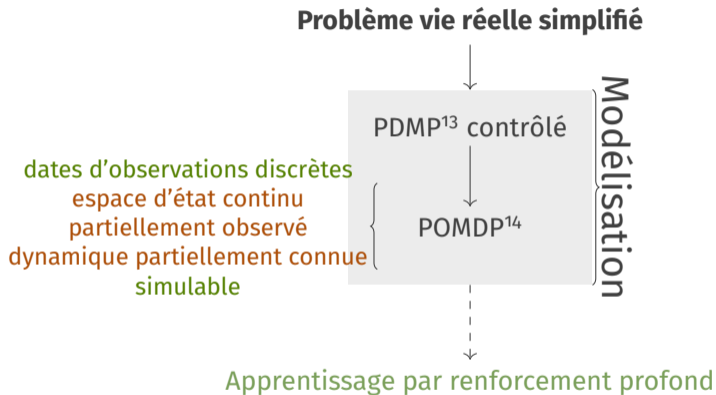


<sup>9</sup>Processus Markovien Déterministe par Morceaux

<sup>10</sup>Processus de Décision Markovien Partiellement Observé

<sup>11</sup>A. Cleyne and B. de Saporta. Numerical method to solve impulse control problems for partially observed piecewise deterministic Markov processes. 2023

<sup>12</sup>A. Cleyne, B. de Saporta, A. Thierry D'Argenlieu, and R. Sabbadin. Medical follow-up optimization : A Monte-Carlo planning strategy. 2024



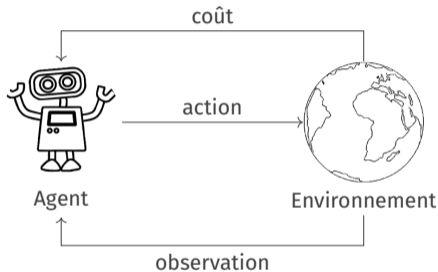
<sup>13</sup>Processus Markovien Déterministe par Morceaux

<sup>14</sup>Processus de Décision Markovien Partiellement Observé

<sup>15</sup>Univ Toulouse, INRAE-MIAT, Toulouse, France



# Apprentissage par renforcement



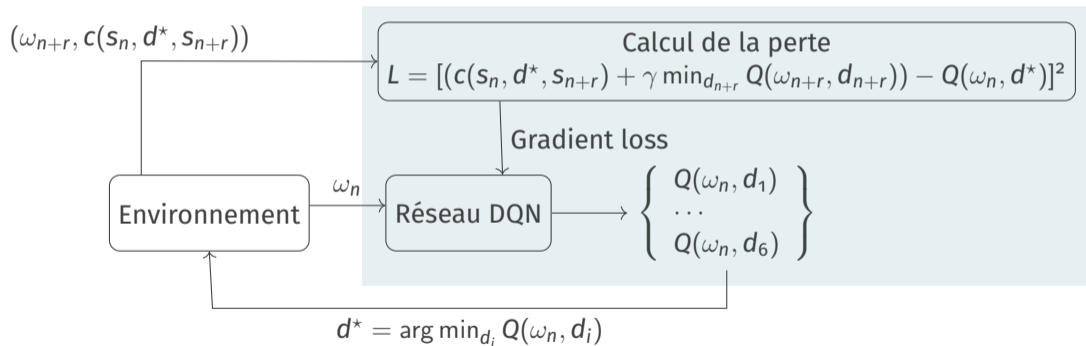
La politique optimale est obtenue à partir  
des expériences  $\langle \omega, d, \omega', c \rangle$

$$\underbrace{Q^\pi(s, d)}_{\text{Critère à optimiser}} = \underbrace{\mathbb{E}^\pi \left[ \sum_{n=0}^{H-1} c(S_{n-1}, D_n, S_n) \mid s, d = (\ell, r) \right]}_{\text{Valeur d'une action dans un état suivant la politique } \pi}$$

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

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

# Algorithme DQN<sup>16</sup>



# Résultats

Politique	Coût moyen (log)	Interval de confiance
<b>OH</b>	8.79	[7.89, 9.69]
<b>Random</b>	11.82	[10.80, 12.84]
<b>Inactive</b>	12.49	[11.54, 13.45]
<b>Threshold</b>	9.89	[8.94, 10.83]
<b>DQN</b>	12.49	[11.54, 13.45]
<b>R2D2</b>	8.47	[7.61, 9.33]

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

# Résultats

Politique	Coût moyen (log)	IC	Taux de survie	Rechutes
<b>OH</b>	8.79	[7.89, 9.69]	93.45%	2.14
<b>Random</b>	11.82	[10.80, 12.84]	27.45%	1.01
<b>Inactive</b>	12.49	[11.54, 13.45]	0.01%	1.00
<b>Threshold</b>	9.89	[8.94, 10.83]	78.95%	1.01
<b>DQN</b>	12.49	[11.54, 13.45]	0.02%	1
<b>R2D2</b>	8.47	[7.61, 9.33]	96.95%	0.65

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

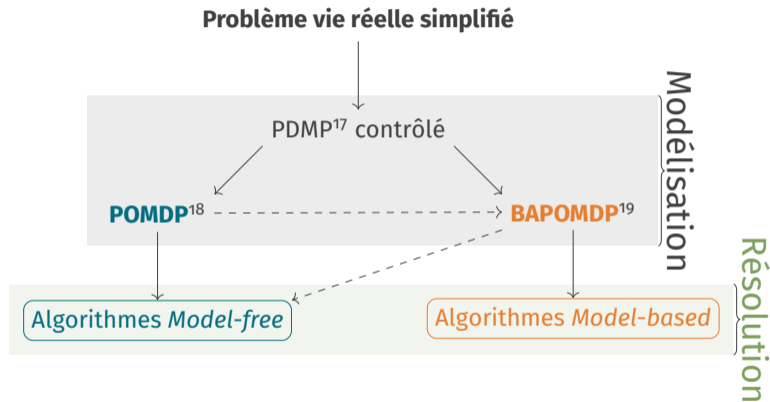
# Résultats

Politique	Coût moyen (log)	IC	Taux de survie	Rechutes
<b>OH</b>	8.79	[7.89, 9.69]	93.45%	2.14
<b>Random</b>	11.82	[10.80, 12.84]	27.45%	1.01
<b>Inactive</b>	12.49	[11.54, 13.45]	0.01%	1.00
<b>Threshold</b>	9.89	[8.94, 10.83]	78.95%	1.01
<b>DQN</b>	12.49	[11.54, 13.45]	0.02%	1
<b>R2D2</b>	8.47	[7.61, 9.33]	96.95%	0.65

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

**Nécessite beaucoup de données pour apprendre la politique optimale !**

# Conclusion et futurs travaux



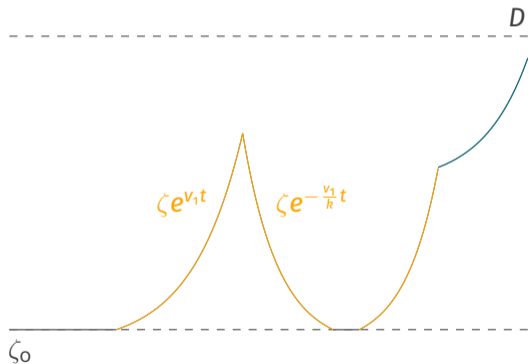
<sup>17</sup>Processus Markovien Déterministe par Morceaux

<sup>18</sup>Processus de Décision Markovien Partiellement Observé

<sup>19</sup>Processus de Décision Markovien Partiellement Observé Bayes Adaptif

# Une dynamique partiellement connue

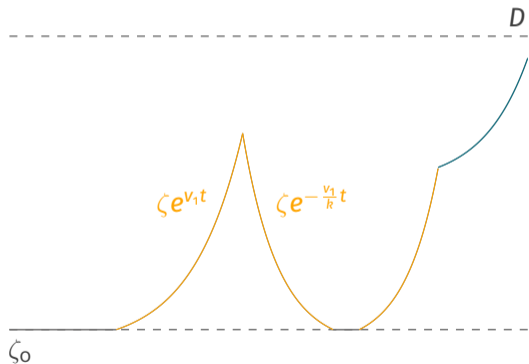
On ne connaît pas le paramètre de pente  $v_1$  de la maladie.



Hypothèse:  
 $v_1 \sim \text{log-normale} (\mu, \sigma^{-2})$ .

# Une dynamique partiellement connue

On ne connaît pas le paramètre de pente  $v_1$  de la maladie.



Hypothèse:

$$v_1 \sim \text{log-normale}(\mu, \sigma^{-2}).$$

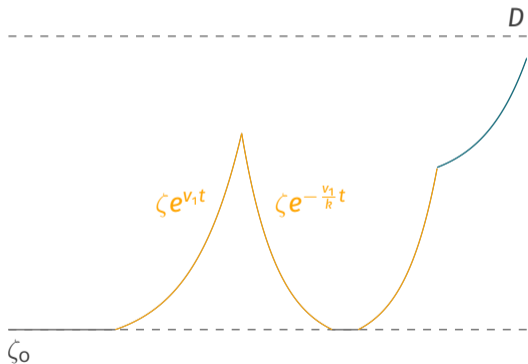
Inférence bayésienne:

$$(\mu, \sigma^{-2}) \sim \text{gamma-log-normale}(\alpha, \beta, \kappa, \nu).$$



# Une dynamique partiellement connue

On ne connaît pas le paramètre de pente  $v_1$  de la maladie.



Hypothèse:

$v_1 \sim \text{log-normale}(\mu, \sigma^{-2})$ .

Inférence bayésienne:

$(\mu, \sigma^{-2}) \sim \text{gamma-log-normale}(\alpha, \beta, \kappa, \nu)$ .

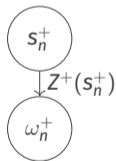
## MISE À JOUR DES HYPERPARAMÈTRES

- $\alpha_{n+1} = \frac{\beta_n \alpha_n + \log(\hat{v}_1)}{\beta_{n+1}}$
- $\beta_{n+1} = \beta_n + 1$
- $\kappa_{n+1} = \kappa_n + \frac{1}{2}$
- $\nu_{n+1} = \nu_n + \frac{\beta_n (\log(\hat{v}_1 - \alpha_n))^2}{2(\beta_{n+1})}$

# Un BAPOMDP<sup>20</sup> partiellement observé

Agent

Environnement



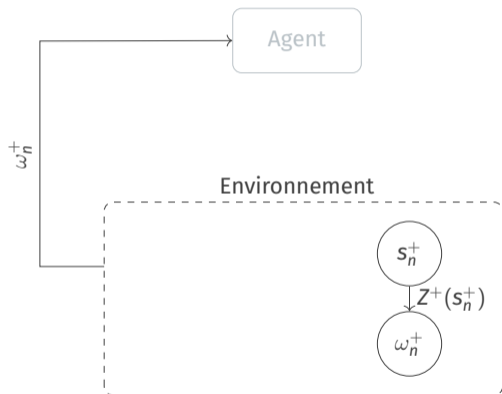
## BAMDP PO

Un BAMDP-PO se définit par un tuple  $(\mathcal{S}^+, \mathcal{D}, \mathcal{K}, \mathcal{P}^+, \Omega^+, \mathcal{Z}^+, \mathcal{C})$ .

- L'hyper-état du patient  $s^+ = (m, k, \zeta, u, \tau, t, \alpha, \beta, \kappa, \nu)$ ;
- Les décisions restent inchangées;
- $\mathcal{K}(\omega) \subseteq \mathcal{D}$  l'espace des décisions admissibles selon l'observation  $\omega$ ;
- La probabilité de transition  $\mathcal{P}(s^+, d)(s')$ ;
- Les observations  $\omega^+ = (z, F(\zeta), \tau, t, \tilde{\alpha}, \tilde{\beta}, \tilde{\kappa}, \tilde{\nu})$ ;
- La fonction d'observations  $\mathcal{Z}(s^+)(\omega^+)$ ;
- La fonction de coût  $\mathcal{C} : \mathcal{D} \times \mathcal{S} \rightarrow \mathbb{R}$ .

<sup>20</sup>Processus de décision Markovien Partiellement Observé Bayes adaptatif

# Un BAPOMDP<sup>20</sup> partiellement observé



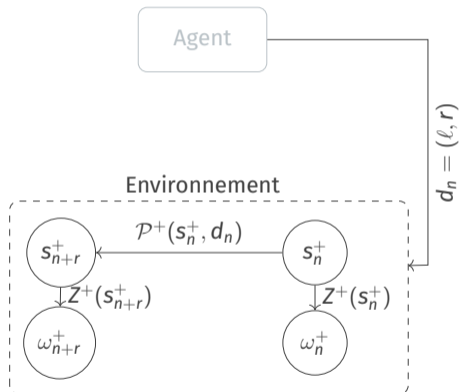
## BAMDP PO

Un BAMDP-PO se définit par un tuple  $(\mathcal{S}^+, \mathcal{D}, \mathcal{K}, \mathcal{P}^+, \Omega^+, \mathcal{Z}^+, \mathcal{C})$ .

- L'hyper-état du patient  $s^+ = (m, k, \zeta, u, \tau, t, \alpha, \beta, \kappa, \nu)$ ;
- Les décisions restent inchangées;
- $\mathcal{K}(\omega) \subseteq \mathcal{D}$  l'espace des décisions admissibles selon l'observation  $\omega$ ;
- La probabilité de transition  $\mathcal{P}(s^+, d)(s')$ ;
- Les observations  $\omega^+ = (z, F(\zeta), \tau, t, \tilde{\alpha}, \tilde{\beta}, \tilde{\kappa}, \tilde{\nu})$ ;
- La fonction d'observations  $\mathcal{Z}(s^+)(\omega^+)$ ;
- La fonction de coût  $\mathcal{C} : \mathcal{D} \times \mathcal{S} \rightarrow \mathbb{R}$ .

<sup>20</sup>Processus de décision Markovien Partiellement Observé Bayes adaptatif

# Un BAPOMDP<sup>20</sup> partiellement observé



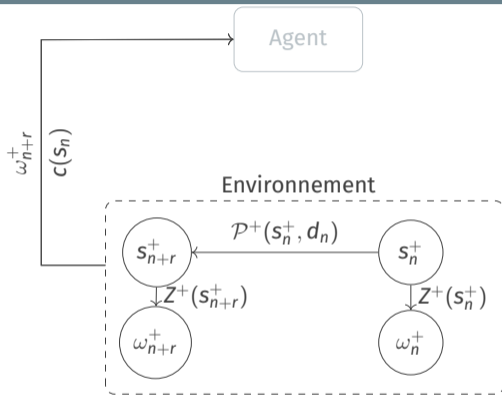
## BAMDP PO

Un BAMDP-PO se définit par un tuple  $(\mathcal{S}^+, \mathcal{D}, \mathcal{K}, \mathcal{P}^+, \Omega^+, \mathcal{Z}^+, \mathcal{C})$ .

- L'hyper-état du patient  $s^+ = (m, k, \zeta, u, \tau, t, \alpha, \beta, \kappa, \nu)$ ;
- Les décisions restent inchangées;
- $\mathcal{K}(\omega) \subseteq \mathcal{D}$  l'espace des décisions admissibles selon l'observation  $\omega$ ;
- La probabilité de transition  $\mathcal{P}(s^+, d)(s')$ ;
- Les observations  $\omega^+ = (z, F(\zeta), \tau, t, \tilde{\alpha}, \tilde{\beta}, \tilde{\kappa}, \tilde{\nu})$ ;
- La fonction d'observations  $\mathcal{Z}(s^+)(\omega^+)$ ;
- La fonction de coût  $\mathcal{C} : \mathcal{D} \times \mathcal{S} \rightarrow \mathbb{R}$ .

<sup>20</sup>Processus de décision Markovien Partiellement Observé Bayes adaptatif

# Un BAPOMDP<sup>20</sup> partiellement observé



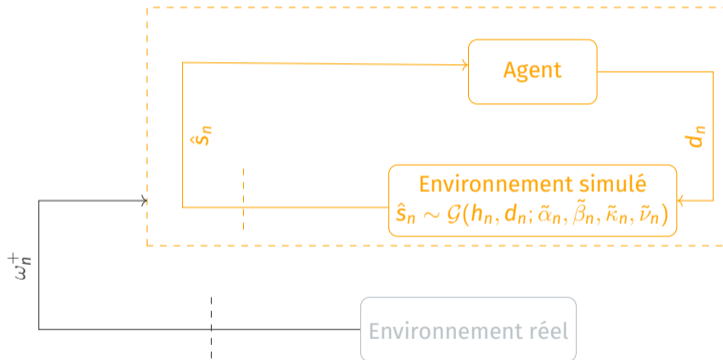
## BAMDP PO

Un BAMDP-PO se définit par un tuple  $(\mathcal{S}^+, \mathcal{D}, \mathcal{K}, \mathcal{P}^+, \Omega^+, \mathcal{Z}^+, C)$ .

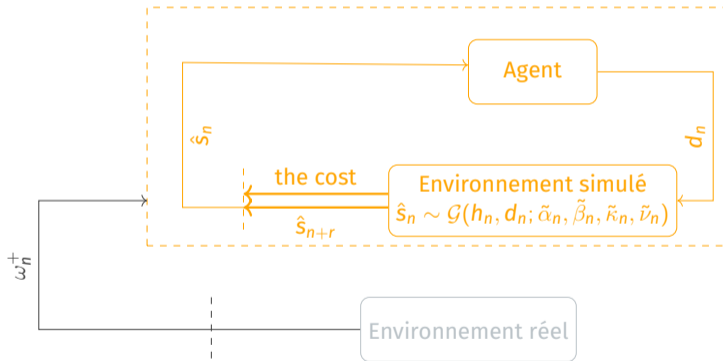
- L'hyper-état du patient  $s^+ = (m, k, \zeta, u, \tau, t, \alpha, \beta, \kappa, \nu)$ ;
- Les décisions restent inchangées;
- $\mathcal{K}(\omega) \subseteq \mathcal{D}$  l'espace des décisions admissibles selon l'observation  $\omega$ ;
- La probabilité de transition  $\mathcal{P}(s^+, d)(s')$ ;
- Les observations  $\omega^+ = (z, F(\zeta), \tau, t, \tilde{\alpha}, \tilde{\beta}, \tilde{\kappa}, \tilde{\nu})$ ;
- La fonction d'observations  $\mathcal{Z}(s^+)(\omega^+)$ ;
- La fonction de coût  $C : \mathcal{D} \times \mathcal{S} \rightarrow \mathbb{R}$ .

<sup>20</sup>Processus de décision Markovien Partiellement Observé Bayes adaptatif

# Une suggestion de résolution



# Une suggestion de résolution



# Une suggestion de résolution

