

Apprentissage par renforcement pour le contrôle de processus de Markov déterministe par morceaux

Orlane Rossini ¹, Alice Cleyen ^{1,2}, Benoîte de Saporta ¹, Régis Sabbadin ³
et Merixell Vinyals ³

¹IMAG, Univ Montpellier, CNRS, Montpellier, France

²John Curtin School of Medical Research, The Australian National University,
Canberra, ACT, Australia

³Univ Toulouse, INRAE-MIAT, Toulouse, France

July 1st 2024



UNIVERSITÉ DE
MONTPELLIER

INRAE

IMAG
INSTITUT MONTPELLIERAIN
ALEXANDER GROTHENDIECK



anr[®]

Le contexte médical

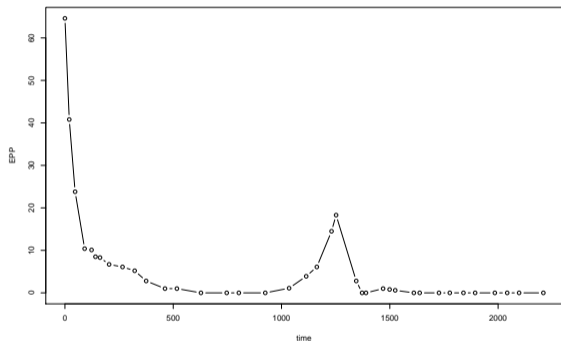


FIGURE: Exemple de donnée d'un patient^a

- 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.

^aIUCT Oncopole et CRCT, Toulouse, France

Le contexte médical

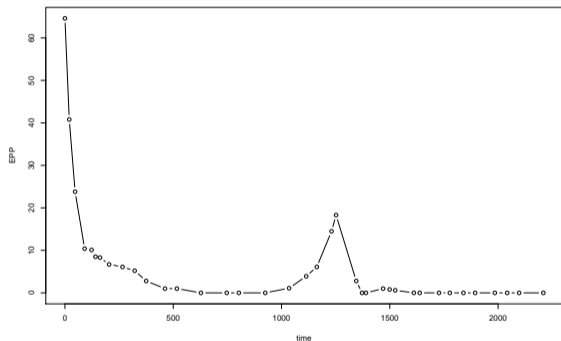
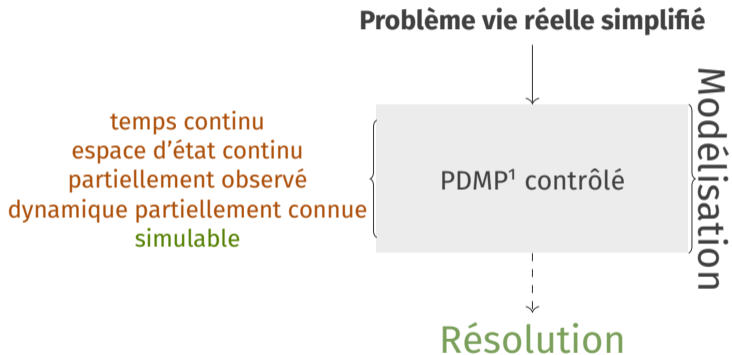


FIGURE: Exemple de donnée d'un patient^a

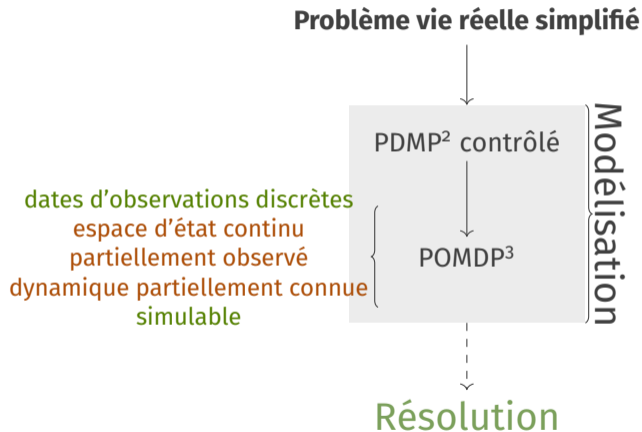
- Des patients ayant eu un **cancer** bénéficient d'un **suivi régulier**;
- La concentration d'**immunoglobuline clonale** est mesurée **dans le temps**;
- Le médecin doit prendre de nouvelles **décisions** à chaque visite.

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

^aIUCT Oncopole et CRCT, Toulouse, France



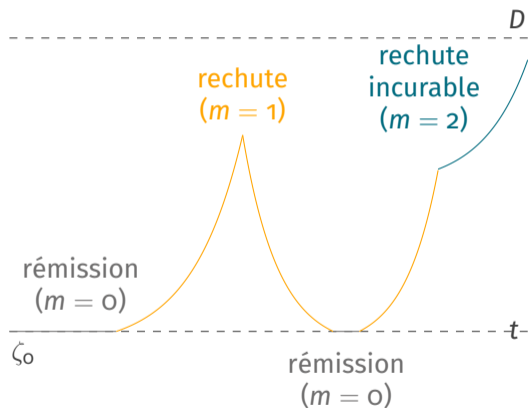
¹Processus Markovien Déterministe par Morceaux



²Processus Markovien Déterministe par Morceaux

³Processus de Décision Markovien Partiellement Observé

Le modèle POMDP⁴



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

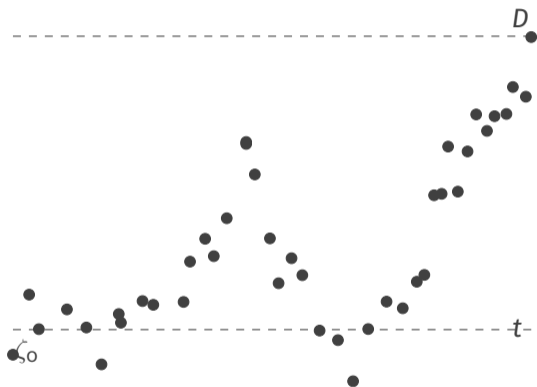
- m état général du patient;
- k nombre de rechute;
- ζ biomarqueur;
- u temps depuis le dernier saut;
- t temps écoulé depuis le début du suivi;
- τ temps depuis l'application d'un traitement.

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

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

⁴Processus de Décision Markovien Partiellement Observé

Le modèle POMDP⁴



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

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

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

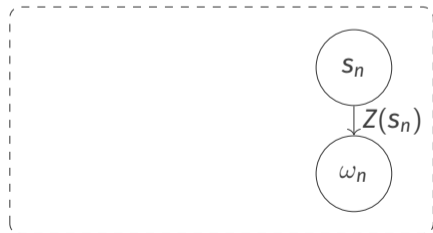
- ℓ traitement (*rien, chimiothérapie*);
- r temps avant la prochaine visite (*15, 30, 60 jours*).

⁴Processus de Décision Markovien Partiellement Observé

Caractéristiques d'un POMDP⁵

Agent

Environnement



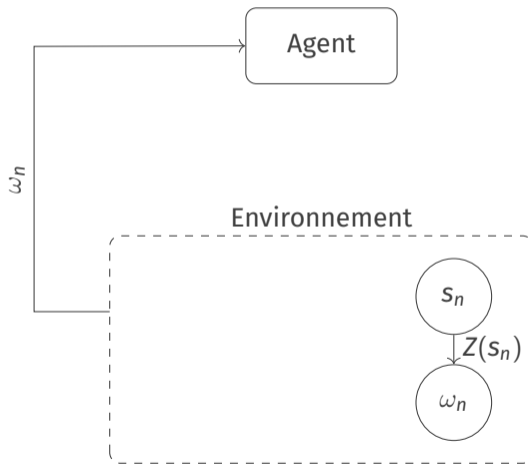
POMDP DEFINITION

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 $C(s, d, s')$.

⁵Processus de Décision Markovien Partiellement Observé

Caractéristiques d'un POMDP⁵



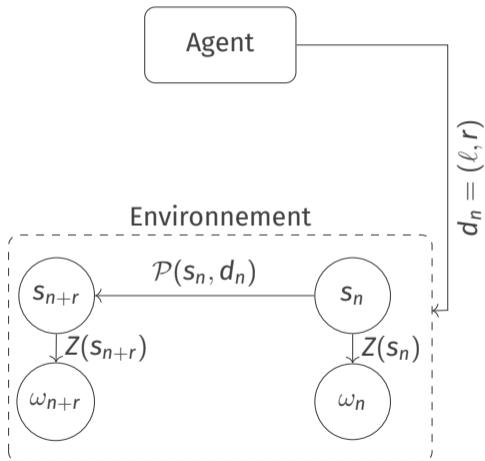
POMDP DEFINITION

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 $C(s, d, s')$.

⁵Processus de Décision Markovien Partiellement Observé

Caractéristiques d'un POMDP⁵



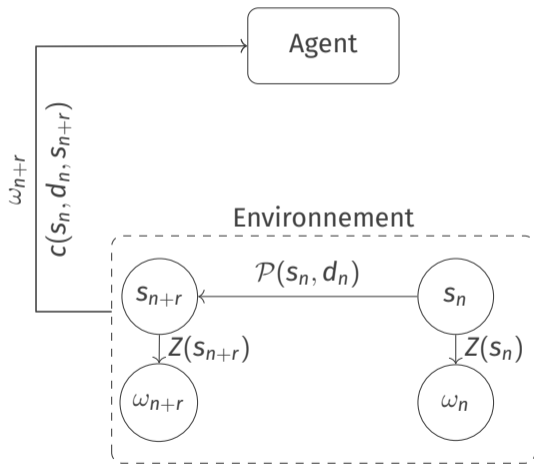
POMDP DEFINITION

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 $C(s, d, s')$.

⁵Processus de Décision Markovien Partiellement Observé

Caractéristiques d'un POMDP⁵



POMDP DEFINITION

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

- Etat du patient $s = (m, k, \zeta, u, t, \tau) \in \mathcal{S}$;
- 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!

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

⁶Processus de Décision Markovien Partiellement Observé

Identifier une politique optimale!

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

⁶Processus de Décision Markovien Partiellement Observé

Résoudre un POMDP⁶

Identifier une politique optimale!

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

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

⁶Processus de Décision Markovien Partiellement Observé

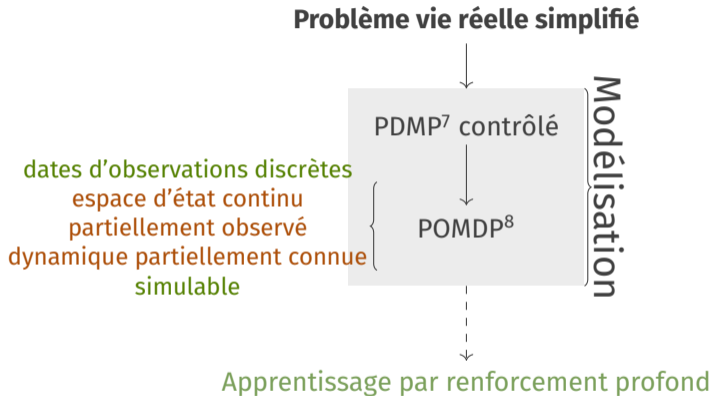
Identifier une politique optimale!

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

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

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

⁶Processus de Décision Markovien Partiellement Observé

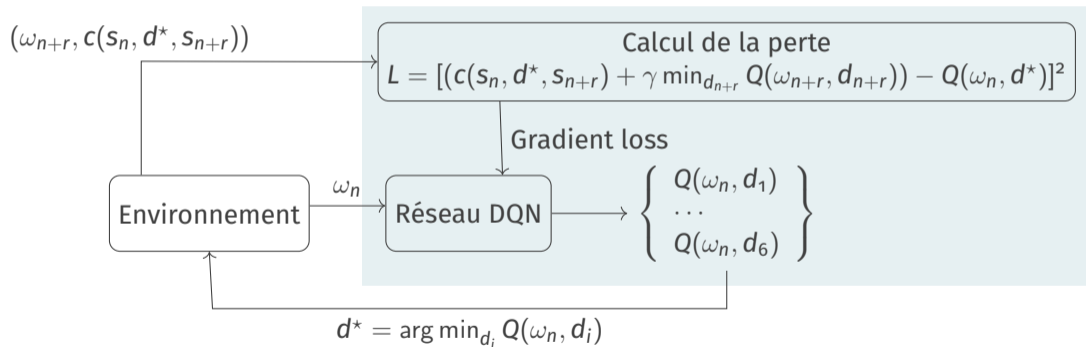


⁷Processus Markovien Déterministe par Morceaux

⁸Processus de Décision Markovien Partiellement Observé

⁹Univ Toulouse, INRAE-MIAT, Toulouse, France

Algorithme DQN¹⁰



¹⁰Deep Q-Network

Résultats

Politique	Coût moyen (log)	Interval de confiance
OH	8.79	[7.89, 9.69]
Random	11.82	[10.80, 12.84]
Inactive	12.49	[11.54, 13.45]
Threshold	9.89	[8.94, 10.83]
DQN	12.49	[11.54, 13.45]
R2D2¹¹	8.47	[7.61, 9.33]

TABLE: Policy evaluation performance on 10^5 simulations

¹¹R2D2 \approx DQN + LSTM

Résultats

Politique	Coût moyen (log)	IC	Taux de survie	Rechutes
OH	8.79	[7.89, 9.69]	93.45%	2.14
Random	11.82	[10.80, 12.84]	27.45%	1.01
Inactive	12.49	[11.54, 13.45]	0.01%	1.00
Threshold	9.89	[8.94, 10.83]	78.95%	1.01
DQN	12.49	[11.54, 13.45]	0.02%	1
R2D2¹¹	8.47	[7.61, 9.33]	96.95%	0.65

TABLE: Policy evaluation performance on 10^5 simulations

¹¹R2D2 \approx DQN + LSTM

Conclusion

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

- Politiques peu conforme à la réalité
- Algorithmes sensibles à la fonction de coût et sa paramétrisation
- Historique important dans l'apprentissage

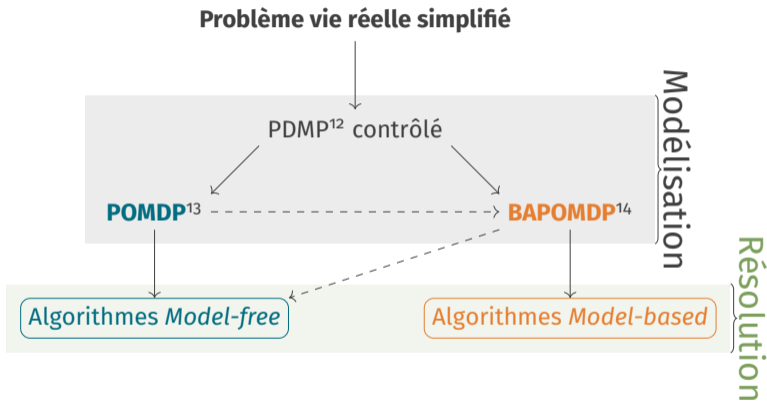
Conclusion

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

- Politiques peu conforme à la réalité
- Algorithmes sensibles à la fonction de coût et sa paramétrisation
- Historique important dans l'apprentissage

Nécessite beaucoup de données pour apprendre la politique optimale !

Perspectives



¹²Processus Markovien Déterministe par Morceaux

¹³Processus de Décision Markovien Partiellement Observé

¹⁴Processus de Décision Markovien Partiellement Observé Bayes Adaptif

- Politiques difficilement interprétables
- Très sensibles à la fonction de coût
- L'historique améliore les performances
- Nécessite beaucoup de données pour apprendre la politique optimale !