

An example of medical treatment optimization under model uncertainty

Orlane Rossini ¹, Aymar Thierry d'Argenlieu ¹, Alice Cleyne ^{1,2}, Benoîte de Saporta ¹ and Régis Sabbadin 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

September 7, 2023



UNIVERSITÉ DE
MONTPELLIER

INRAE

IMAG

INSTITUT MONTPELLIERAIN
ALEXANDER GROTHENDIECK



anr[®]

Contents

- ▶ Introduction
- ▶ Mathematical Model Introduction
- ▶ A Framework for Partial Observability
- ▶ A Framework for Unknown Transitions
- ▶ Conclusion and Perspectives



A medical context

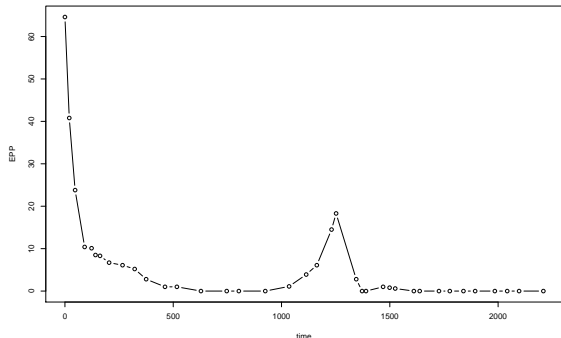
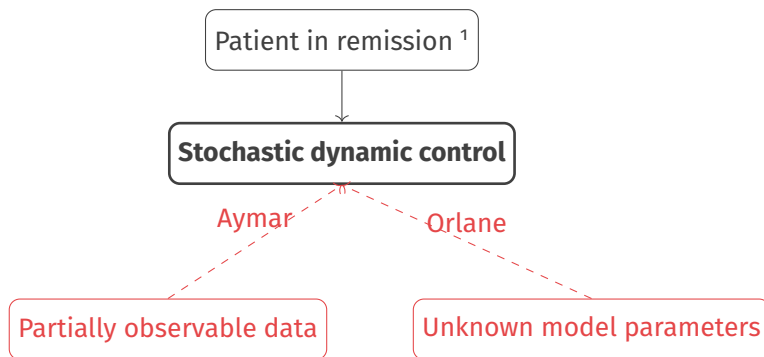


Figure: Patient Data^a

^aData from IUC Oncopole, Toulouse, and CRCT, Toulouse, France

- Patients who have had **cancer** are **regularly monitored**;
- **Clonal immunoglobulin** concentration is monitored **over time**;
- The doctor has to make new **decisions** at each visit.

A medical context



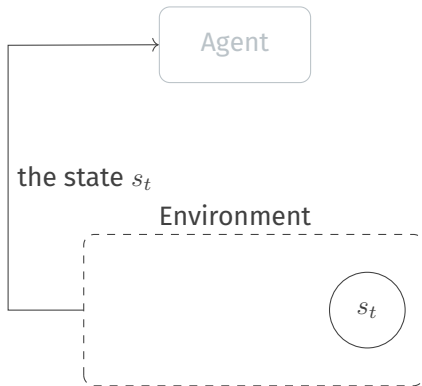
¹Data from IUC Oncopole, Toulouse, and CRCT, Toulouse, France

Contents

- ▶ Introduction
- ▶ **Mathematical Model Introduction**
- ▶ A Framework for Partial Observability
- ▶ A Framework for Unknown Transitions
- ▶ Conclusion and Perspectives



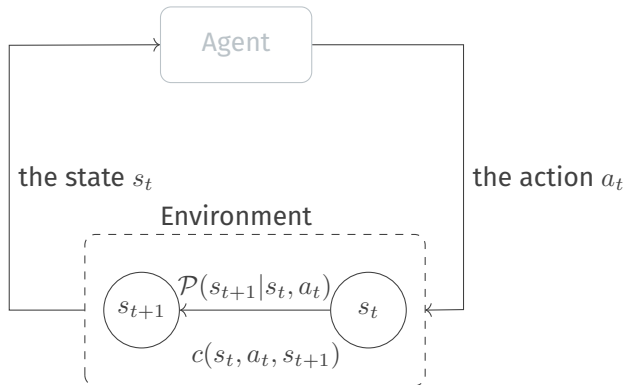
Markov Decision Process (MDP²)



- $s \in \mathcal{S}$ the state space
- $a \in \mathcal{A}$ the action space
- \mathcal{P} the transition matrix
- $c(s_t, a_t)$ the cost function

²ML Puterman (1994). "Finite-horizon Markov decision processes". In: *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. New York: Wiley-Interscience, pp. 78–9.

Markov Decision Process (MDP²)



- $s \in \mathcal{S}$ the state space
- $a \in \mathcal{A}$ the action space
- \mathcal{P} the transition matrix
- $c(s_t, a_t)$ the cost function

²ML Puterman (1994). "Finite-horizon Markov decision processes". In: *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. New York: Wiley-Interscience, pp. 78–9.

Solving a MDP

Minimizing a cost

Policy π

Let $f : \mathcal{S} \rightarrow \mathcal{A}$ for all $s \in \mathcal{S}$ is a decision rule.

A sequence of decision rules $\pi = (f_0, f_1, \dots, f_{H-1})$ is a policy.

Let Π be the set of all eligible policies.

Policy cost and value function

$$J_{\pi}(s_0) = \mathbb{E}\left[\sum_{t=0}^{H-1} c(S_t, A_t) \mid \pi(S_t), S_0 = s_0\right]$$

Let π^* the optimal policy such that:

$$V(s_0) = J_{\pi^*}(s_0) = \min_{\pi \in \Pi} J_{\pi}(s_0)$$

Optimization criterion

$$V^*(s_t) = \min_{a_t \in \mathcal{A}} [c(s_t, a_t) + \sum_{s_{t+1} \in \mathcal{S}} \mathcal{P}(s_{t+1} \mid s_t, a_t) V^*(s_{t+1})]$$

A model-free method

Q-learning^{3,4} algorithm

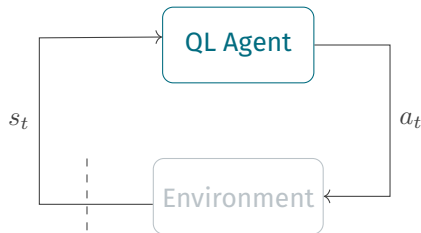


³Christopher J. C. H. Watkins and Peter Dayan (May 1992). “Q-learning”. In: *Mach. Learn.* 8.3, pp. 279–292. ISSN: 1573-0565. DOI: [10.1007/BF00992698](https://doi.org/10.1007/BF00992698).

⁴VP Vivek and Dr. Shalabh Bhatnagar (Aug. 2022). “Finite Horizon Q-learning: Stability, Convergence, Simulations and an application on Smart Grids”. In: [arXiv:2110.15093v3](https://arxiv.org/abs/2110.15093v3). DOI: [10.48550/arXiv.2110.15093](https://doi.org/10.48550/arXiv.2110.15093). eprint: [2110.15093v3](https://arxiv.org/abs/2110.15093v3).

A model-free method

Q-learning^{3,4} algorithm



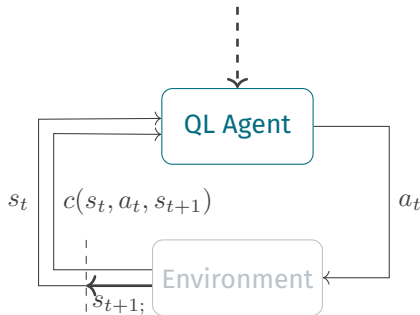
³Christopher J. C. H. Watkins and Peter Dayan (May 1992). “Q-learning”. In: *Mach. Learn.* 8.3, pp. 279–292. ISSN: 1573-0565. DOI: [10.1007/BF00992698](https://doi.org/10.1007/BF00992698).

⁴VP Vivek and Dr. Shalabh Bhatnagar (Aug. 2022). “Finite Horizon Q-learning: Stability, Convergence, Simulations and an application on Smart Grids”. In: [arXiv:2110.15093v3](https://arxiv.org/abs/2110.15093v3). DOI: [10.48550/arXiv.2110.15093](https://doi.org/10.48550/arXiv.2110.15093). eprint: [2110.15093v3](https://arxiv.org/abs/2110.15093v3).

A model-free method

Q-learning^{3,4} algorithm

$$Q_n(s_t, a_t) = (1 - \alpha)Q_{n-1}(s_t, a_t) + \alpha[c(s_t, a_t) + \min_{a_{t+1} \in \mathcal{A}} Q_{n-1}(s_{t+1}, a_{t+1})]$$



³Christopher J. C. H. Watkins and Peter Dayan (May 1992). "Q-learning". In: *Mach. Learn.* 8.3, pp. 279–292. ISSN: 1573-0565. DOI: [10.1007/BF00992698](https://doi.org/10.1007/BF00992698).

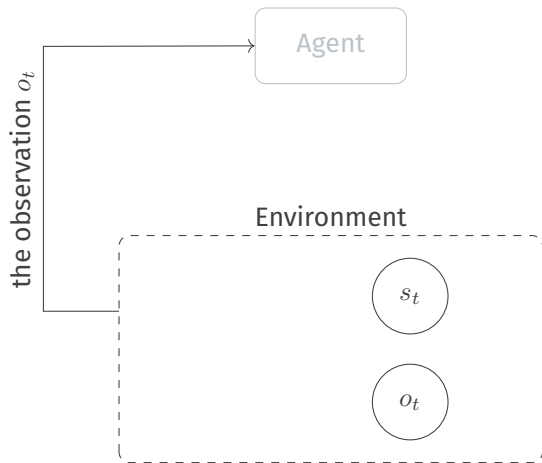
⁴VP Vivek and Dr. Shalabh Bhatnagar (Aug. 2022). "Finite Horizon Q-learning: Stability, Convergence, Simulations and an application on Smart Grids". In: [arXiv:2110.15093v3](https://arxiv.org/abs/2110.15093v3). DOI: [10.48550/arXiv.2110.15093](https://doi.org/10.48550/arXiv.2110.15093). eprint: [2110.15093v3](https://arxiv.org/abs/2110.15093v3).

Contents

- ▶ Introduction
- ▶ Mathematical Model Introduction
- ▶ **A Framework for Partial Observability**
- ▶ A Framework for Unknown Transitions
- ▶ Conclusion and Perspectives

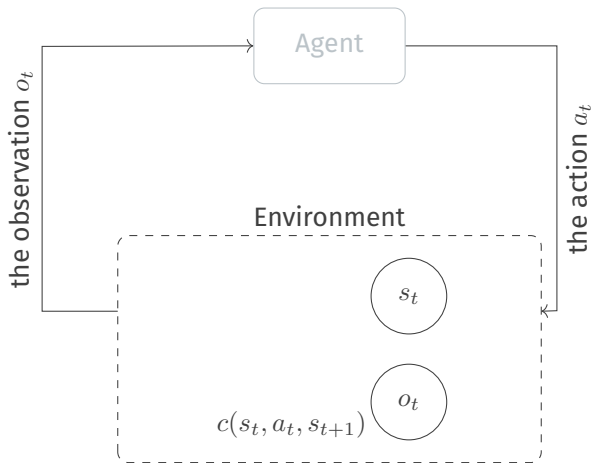


Partially observable Markov Decision Process (POMDP)



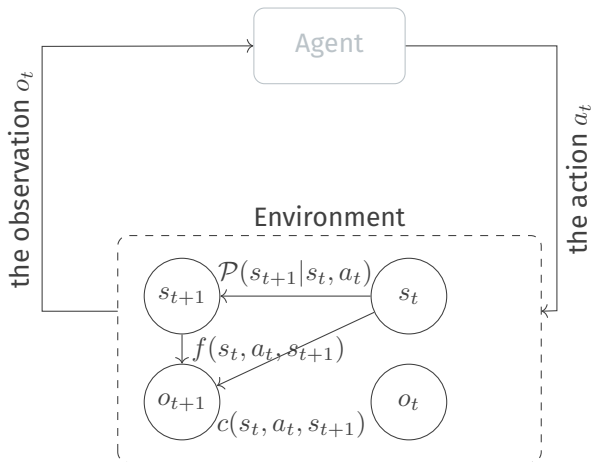
- $s \in \mathcal{S}$ the state space
- $o \in \mathcal{O}$ the observation space
- $a \in \mathcal{A}$ the action space
- \mathcal{P} the transition matrix
- f a measurable function
- $c(s_t, a_t, s_{t+1})$ the cost function

Partially observable Markov Decision Process (POMDP)



- $s \in \mathcal{S}$ the state space
- $o \in \mathcal{O}$ the observation space
- $a \in \mathcal{A}$ the action space
- \mathcal{P} the transition matrix
- f a measurable function
- $c(s_t, a_t, s_{t+1})$ the cost function

Partially observable Markov Decision Process (POMDP)



- $s \in \mathcal{S}$ the state space
- $o \in \mathcal{O}$ the observation space
- $a \in \mathcal{A}$ the action space
- \mathcal{P} the transition matrix
- f a measurable function
- $c(s_t, a_t, s_{t+1})$ the cost function

Solving a POMDP

Minimizing a cost

The *history* is defined as a sequence of actions and observations.

A history

$$h_t = \{o_0, a_0, o_1, a_1, \dots, o_{t-1}, a_{t-1}, o_t\}$$

Solving a POMDP

Minimizing a cost

The *history* is defined as a sequence of actions and observations.

A history

$$h_t = \{o_0, a_0, o_1, a_1, \dots, o_{t-1}, a_{t-1}, o_t\}$$

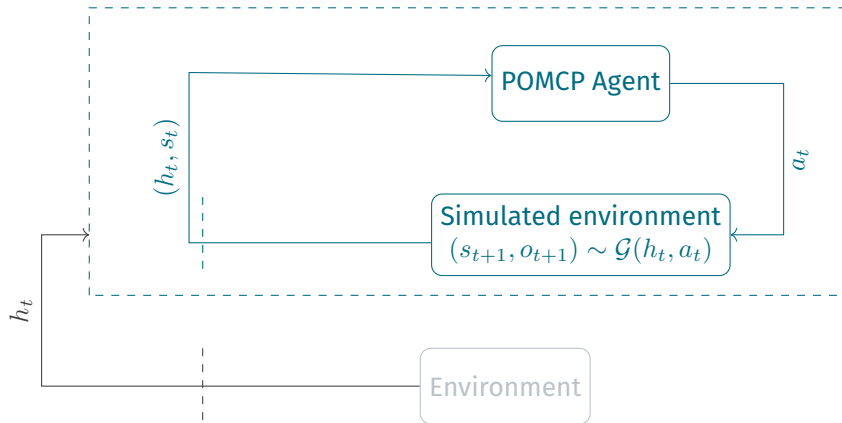
Optimization criterion

$$V^*(h) = \min_{a_t \in \mathcal{A}} [c(s_t, a_t) + \sum_{o_{t+1} \in \mathcal{O}} \mathcal{P}(o_{t+1} | h_{t+1}, a_t) V^*(h_{t+1})]$$

A model-based method

Partially Observable Monte-Carlo Planning (POMCP⁵)

with $h_t = (o_0, a_0, o_1, \dots, o_{t-1}, a_{t-1}, o_t)$,

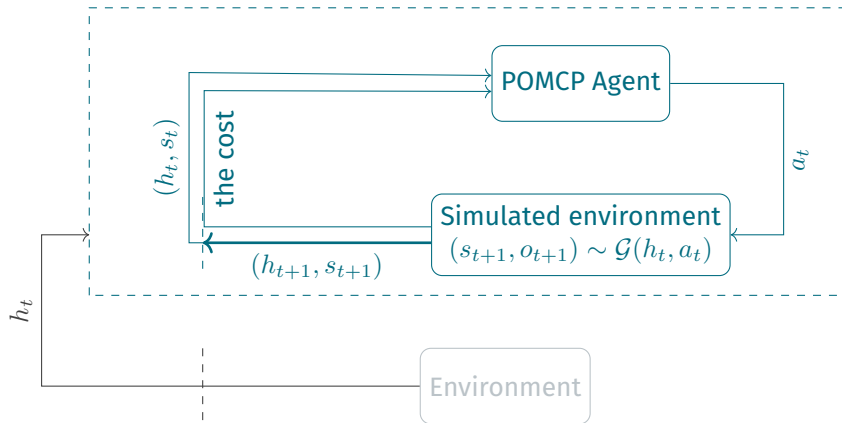


⁵David Silver and Joel Veness (2010). "Monte-Carlo Planning in Large POMDPs". In: *Advances in Neural Information Processing Systems* 23. URL: https://papers.nips.cc/paper_files/paper/2010/hash/edfbe1afcf9246bb0d40eb4d8027d90f-Abstract.html

A model-based method

Partially Observable Monte-Carlo Planning (POMCP⁵)

with $h_t = (o_0, a_0, o_1, \dots, o_{t-1}, a_{t-1}, o_t)$,

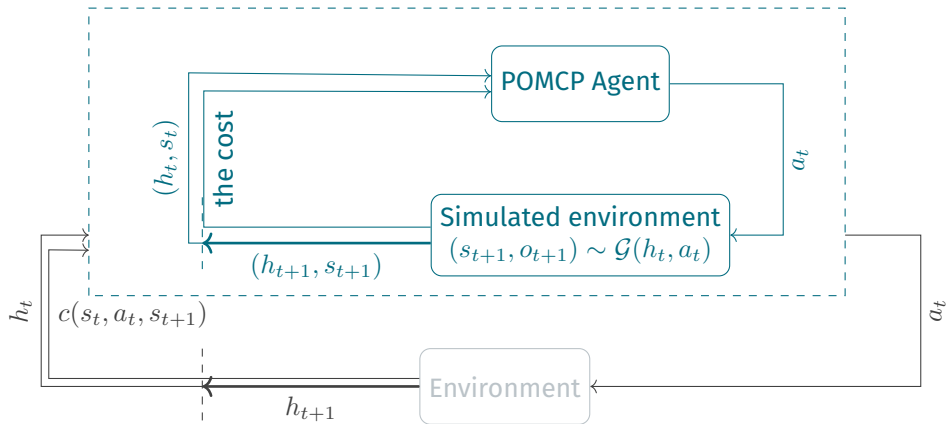


⁵David Silver and Joel Veness (2010). "Monte-Carlo Planning in Large POMDPs". In: *Advances in Neural Information Processing Systems 23*. URL: https://papers.nips.cc/paper_files/paper/2010/hash/edfbe1afcf9246bb0d40eb4d8027d90f-Abstract.html

A model-based method

Partially Observable Monte-Carlo Planning (POMCP⁵)

with $h_t = (o_0, a_0, o_1, \dots, o_{t-1}, a_{t-1}, o_t)$,



⁵David Silver and Joel Veness (2010). "Monte-Carlo Planning in Large POMDPs". In: *Advances in Neural Information Processing Systems* 23. URL: https://papers.nips.cc/paper_files/paper/2010/hash/edfbe1afcf9246bb0d40eb4d8027d90f-Abstract.html

Contents

- ▶ Introduction
- ▶ Mathematical Model Introduction
- ▶ A Framework for Partial Observability
- ▶ A Framework for Unknown Transitions**
- ▶ Conclusion and Perspectives



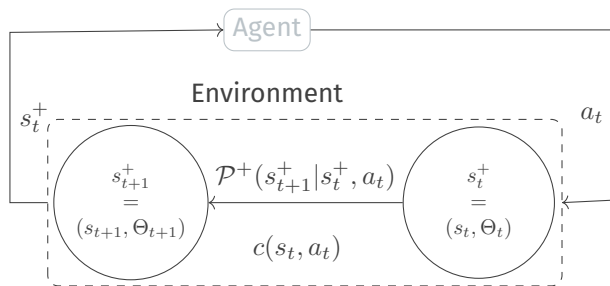
A bayesian approach

$s_t \backslash s_{t+1}$	(0, 0, 0)	(1, 0, 1)	(1, 0, 2)	(1, 1, 1)	(1, 1, 2)	(1, 2, 1)	(1, 2, 2)	(1, 3, 1)	(1, 3, 2)	(2, 4, 0)
(0, 0, 0)	$p_{(0,0,0)}^\theta$	$p_{(1,0,1)}^\theta$	$p_{(1,0,2)}^\theta$	0	0	0	0	0	0	0

Remark:

- $P(\cdot | s = (0, 0, 0), a = \emptyset) \sim \mathcal{M}(p_{(0,0,0)}^\theta, p_{(1,0,1)}^\theta, p_{(1,0,2)}^\theta)$
- Conjugate distribution : $f(p^\theta | \Theta^\theta) \sim \mathcal{D}(\theta_{(0,0,0)}^\theta, \theta_{(1,0,1)}^\theta, \theta_{(1,0,2)}^\theta)$

Bayes-Adaptive Markov Decision Process (BAMDP)⁶

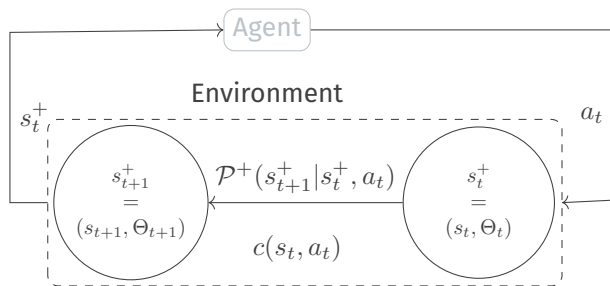


- $s^+ \in \mathcal{S}^+$ the hyper-state space
- \mathcal{P}^+ the transition matrix
- $\Theta_{t+1} = \Theta_t + \Delta_{s_{t+1}}^{a_t}$, with

$$\Delta_{s_{t+1}}^{a_t} = \begin{cases} 1 & \text{if } (s = (\mathbf{0}, \mathbf{0}, \mathbf{0}), a_t, s_{t+1}), \\ 0 & \text{else.} \end{cases}$$

⁶Michael O'Gordon Duff (2002). "Optimal learning: Computational procedures for Bayes -adaptive Markov decision processes". PhD thesis. University of Massachusetts Amherst.

Bayes-Adaptive Markov Decision Process (BAMDP⁶)



- $s^+ \in \mathcal{S}^+$ the hyper-state space
- \mathcal{P}^+ the transition matrix
- $\Theta_{t+1} = \Theta_t + \Delta_{s_{t+1}}^{a_t}$, with

$$\Delta_{s_{t+1}}^{a_t} = \begin{cases} 1 & \text{if } (s = (\mathbf{0}, \mathbf{0}, \mathbf{0}), a_t, s_{t+1}), \\ 0 & \text{else.} \end{cases}$$

Optimization criterion

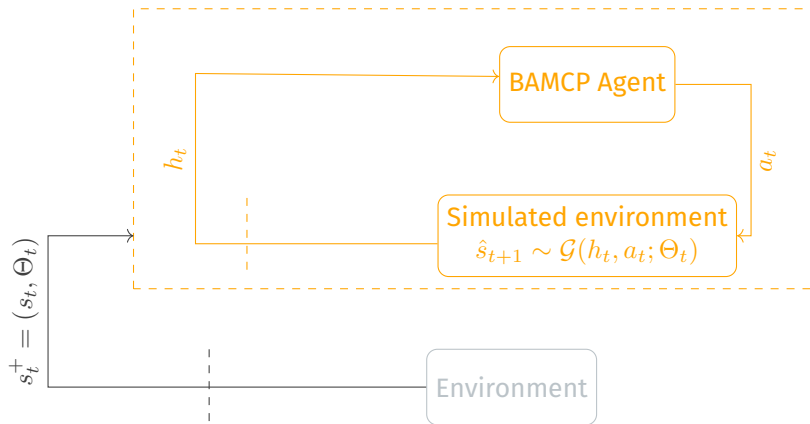
$$V^*(s_t, \Theta_t) = \min_{a_t \in \mathcal{A}} [c(s_t, a_t) + \sum_{s_{t+1}^+ \in \mathcal{S}^+} \mathcal{P}^+(s_{t+1}^+ | s_t^+, a_t) V^*(s_{t+1}, \Theta_{t+1})]$$

⁶Michael O'Gordon Duff (2002). "Optimal learning: Computational procedures for Bayes -adaptive Markov decision processes". PhD thesis. University of Massachusetts Amherst.

A model-based method

Bayes-Adaptive Monte-Carlo Planning (BAMCP⁷)

with $h_t = (s_0, a_0, s_1, \dots, s_{t-1}, a_{t-1}, s_t)$,

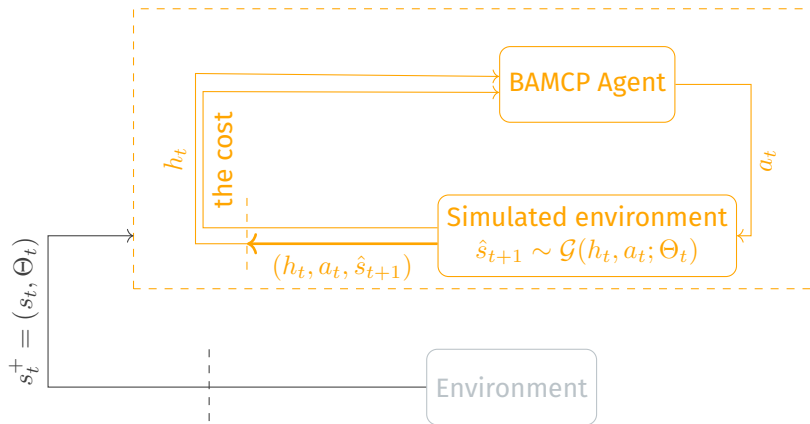


⁷Arthur Guez, David Silver, and Peter Dayan (2012). "Efficient Bayes-Adaptive Reinforcement Learning using Sample-Based Search". In: *Advances in Neural Information Processing Systems*. Ed. by F. Pereira et al. Vol. 25. Curran Associates, Inc.

A model-based method

Bayes-Adaptive Monte-Carlo Planning (BAMCP⁷)

with $h_t = (s_0, a_0, s_1, \dots, s_{t-1}, a_{t-1}, s_t)$,

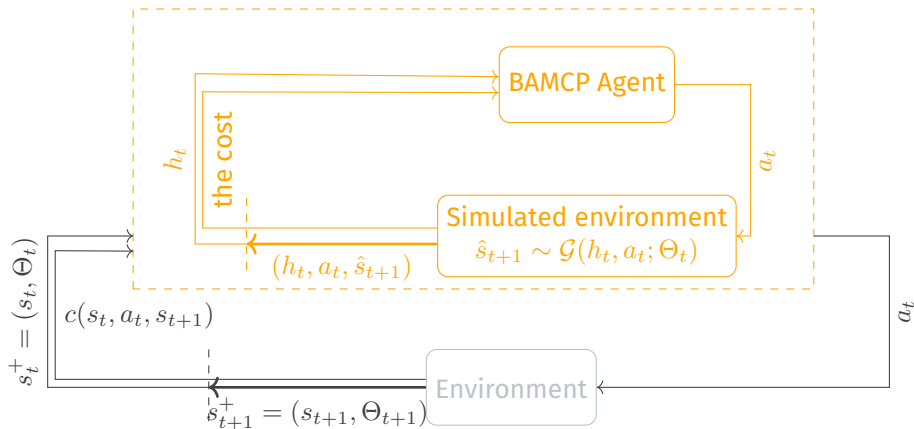


⁷Arthur Guez, David Silver, and Peter Dayan (2012). "Efficient Bayes-Adaptive Reinforcement Learning using Sample-Based Search". In: *Advances in Neural Information Processing Systems*. Ed. by F. Pereira et al. Vol. 25. Curran Associates, Inc.

A model-based method

Bayes-Adaptive Monte-Carlo Planning (BAMCP⁷)

with $h_t = (s_0, a_0, s_1, \dots, s_{t-1}, a_{t-1}, s_t)$,

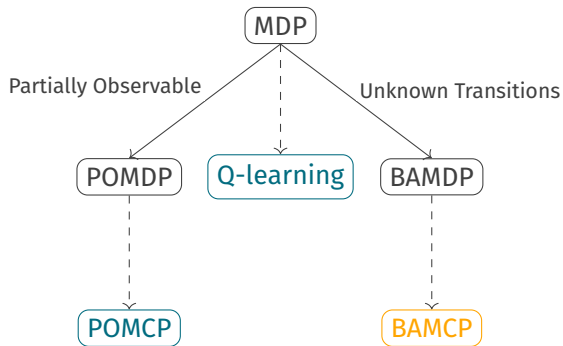


⁷Arthur Guez, David Silver, and Peter Dayan (2012). "Efficient Bayes-Adaptive Reinforcement Learning using Sample-Based Search". In: *Advances in Neural Information Processing Systems*. Ed. by F. Pereira et al. Vol. 25. Curran Associates, Inc.

Contents

- ▶ Introduction
- ▶ Mathematical Model Introduction
- ▶ A Framework for Partial Observability
- ▶ A Framework for Unknown Transitions
- ▶ **Conclusion and Perspectives**

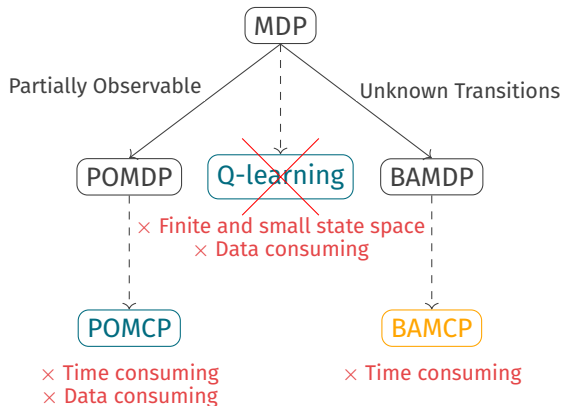
Conclusion



- ◆ Mathematical framework
- ◆ Model-free method
- ◆ Model-based method

Unlike model-free methods and deep reinforcement learning, **model-based approaches** do not require as much interaction with the environment.

Conclusion



- ◆ Mathematical framework
- ◆ Model-free method
- ◆ Model-based method

Unlike model-free methods and deep reinforcement learning, **model-based approaches** do not require as much interaction with the environment.

Perspectives

A real-life problem

Modelling



POMDP

× partially observable

BAMDP

× partially unknown model

BAPOMDP

× partially unknown model
× partially observable

Controlled PDMP

× partially observable
× partially unknown model
× semi-Markov

MDP

× large state space
× continuous state space

Resolution

Exact resolution by DP is no longer possible.
Resolution by **simulations** must be applied.

Perspectives

Modelling

POMDP

BAMDP

Controlled PDMP

BAPOMDP

MDP

- × large state space
- × continuous state space

Resolution

Exact resolution by DP is no longer possible.
Resolution by **simulations** must be applied.

POMCP

BAMCP

BAPOMCP

Deel RL

- ✓ convergence guarantees
 - × applicabilty limited ?
 - × data consuming
 - × on-line time consuming
- ✓ convergence guarantees
 - × applicabilty limited ?
 - × data consuming
 - × on-line time consuming
- ✓ convergence guarantees
 - × applicabilty limited ?
 - × data consuming
 - × on-line time consuming
- × no convergence guarantees
 - ✓ always applicable
 - × data consuming
 - ✓ on-line time

Any questions ?