

An example of medical treatment optimization under model uncertainty

Orlane Le Quellenec 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

June 30, 2023



UNIVERSITÉ DE
MONTPELLIER

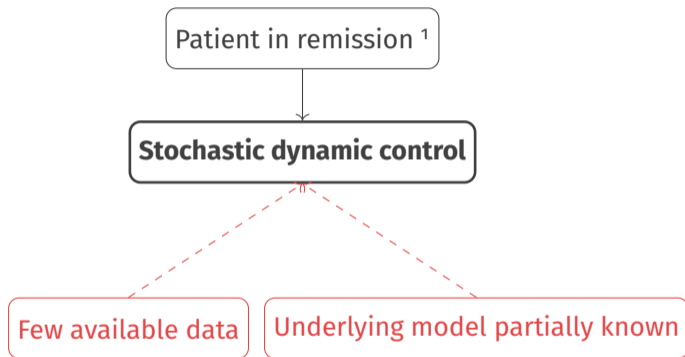
INRAE

IMAG
INSTITUT MONTPELLIERAIN
ALEXANDER GROTHENDIECK



anr[®]

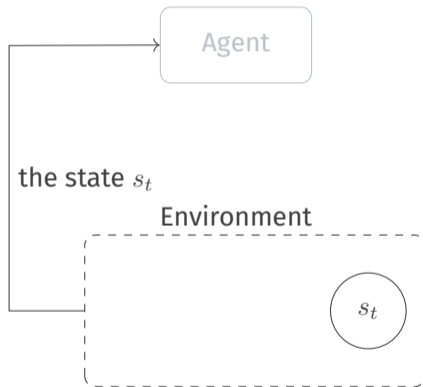
A medical context



How can these issues be addressed in a simplified problem ?

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

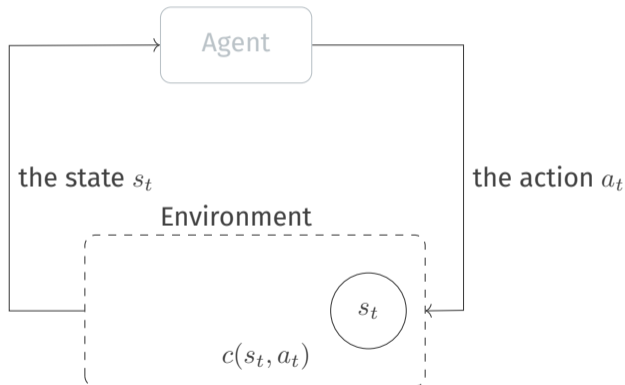
Markov Decision Process (MDP²)



- $s \in \mathcal{S}$ 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.

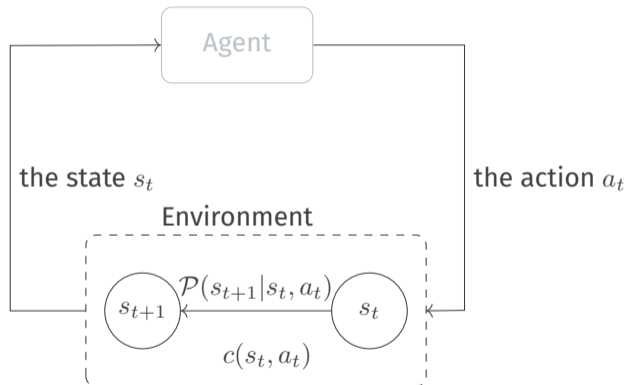
Markov Decision Process (MDP²)



- $s \in \mathcal{S}$ 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.

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.

State transition

The transition matrix is partially known

Table: Transition matrix when patient has no treatment ($a = \emptyset$).

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

State transition

The transition matrix is partially known

Table: Transition matrix when patient has treatment ($a = \rho$).

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

Solving a MDP

Minimizing a cost

The list of costs:

- Treatment: 300
- Disease 1: 200
- Disease 2: 300
- Death: 1000

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.

Policy cost and value function

$$J_H(\pi, s) = \mathbb{E} \left[\sum_{t=0}^{H-1} c(s_t, a_t) \mid \pi, s \right]$$
$$V_H(s) = \inf_{\pi \in \Pi} J_H(\pi, s)$$

Solving a MDP

Minimizing a cost

The list of costs:

- Treatment: 300
- Disease 1: 200
- Disease 2: 300
- Death: 1000

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.

Policy cost and value function

$$J_H(\pi, s) = \mathbb{E} \left[\sum_{t=0}^{H-1} c(s_t, a_t) \mid \pi, s \right]$$
$$V_H(s) = \inf_{\pi \in \Pi} J_H(\pi, s)$$

Optimization criterion

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

A model-free method

Q-learning^{3,4} algorithm

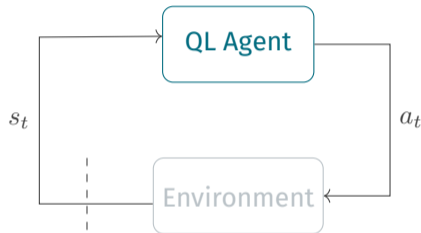


³Christopher J. C. H. Watkins and Peter Dayan (May 1992). “Q-learning”. In: *Mach. Learn.* 8.3, pp. 279–292. ISSN: 1573-0565. DOI: [10.1007/BF00992698](https://doi.org/10.1007/BF00992698).

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



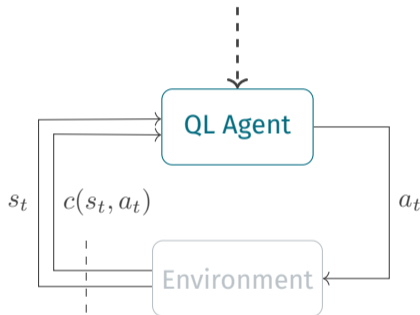
³Christopher J. C. H. Watkins and Peter Dayan (May 1992). "Q-learning". In: *Mach. Learn.* 8.3, pp. 279–292. ISSN: 1573-0565. DOI: [10.1007/BF00992698](https://doi.org/10.1007/BF00992698).

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

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



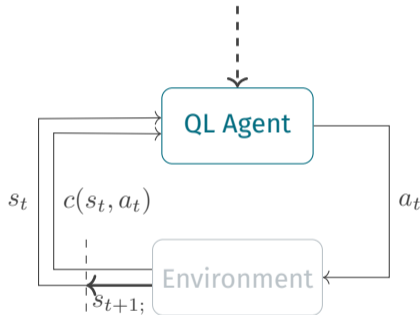
³Christopher J. C. H. Watkins and Peter Dayan (May 1992). “Q-learning”. In: *Mach. Learn.* 8.3, pp. 279–292. ISSN: 1573-0565. DOI: [10.1007/BF00992698](https://doi.org/10.1007/BF00992698).

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).

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).

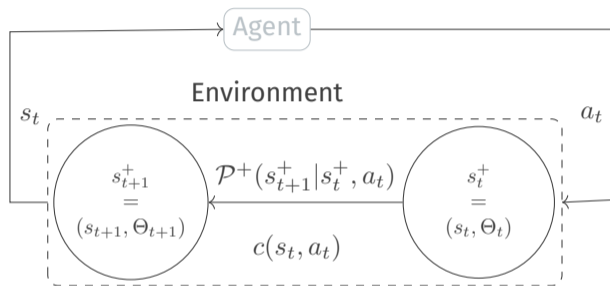
A bayesian approach

$s_t \backslash s_{t+1}$	(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 (BAMDP⁵)

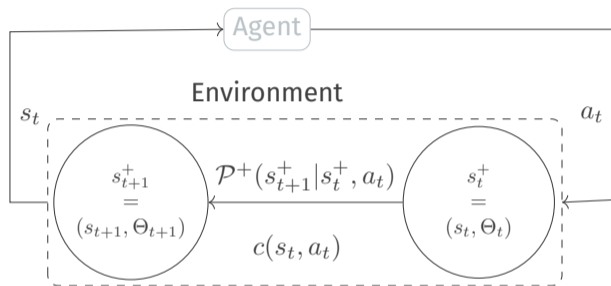


- $s^+ \in \mathcal{S}^+$ the hyper-state space
- \mathcal{P}^+ the transition matrix
- $\Theta_{t+1} = \Theta_t + \Delta_{s_{t+1}}^{a_t}$, with

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

⁵Michael 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 (BAMDP⁵)



- $s^+ \in \mathcal{S}^+$ 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

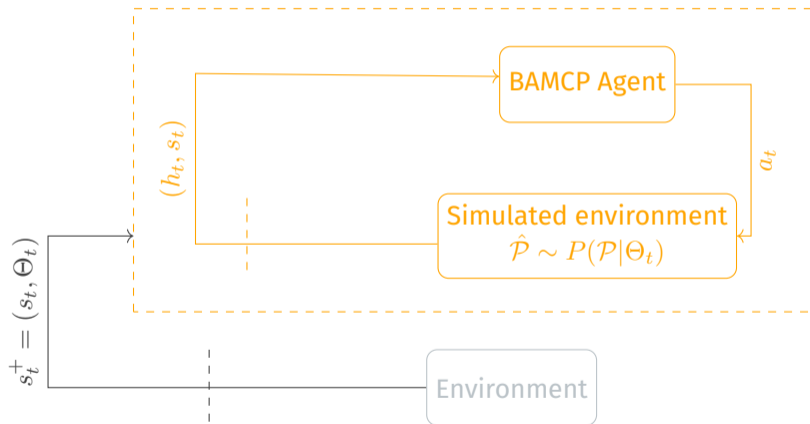
$$V^*(s_t, \Theta_t) = \min_{a \in \mathcal{A}} [c(s_t, a) + \sum_{s_{t+1}^+ \in \mathcal{S}^+} \mathcal{P}^+(s_{t+1}^+ | s_t^+, a) V^*(s_{t+1}, \Theta_{t+1})]$$

⁵Michael 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 (BAMCP)⁶

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

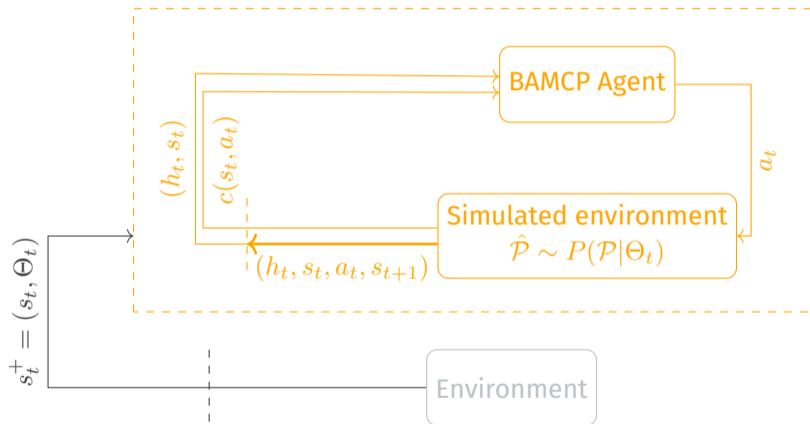


⁶Arthur Guez, David Silver, and Peter Dayan (2012). "Efficient Bayes-Adaptive Reinforcement Learning using Sample-Based Search". In: *Advances in Neural Information Processing Systems*. Ed. by F. Pereira et al. Vol. 25. Curran Associates, Inc.

A model-based method

Bayes-Adaptive Monte-Carlo Planning (BAMCP)⁶

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

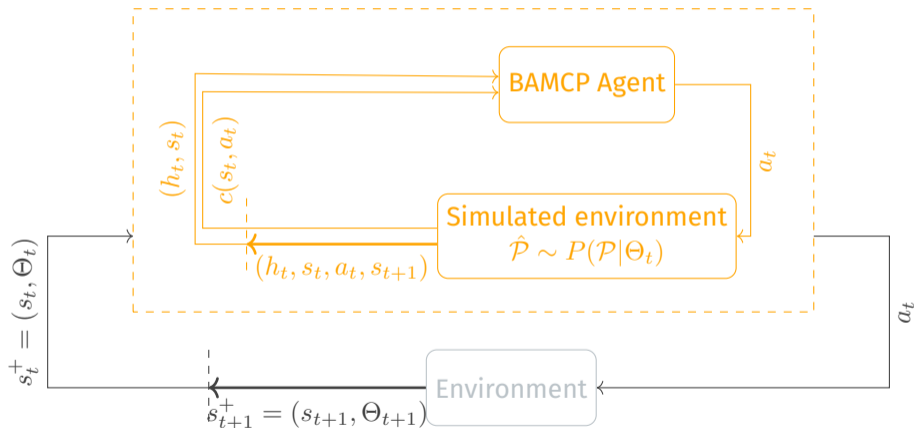


⁶Arthur Guez, David Silver, and Peter Dayan (2012). "Efficient Bayes-Adaptive Reinforcement Learning using Sample-Based Search". In: *Advances in Neural Information Processing Systems*. Ed. by F. Pereira et al. Vol. 25. Curran Associates, Inc.

A model-based method

Bayes-Adaptive Monte-Carlo Planning (BAMCP)⁶

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



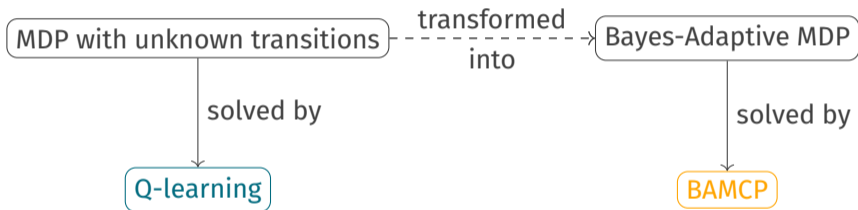
⁶Arthur Guez, David Silver, and Peter Dayan (2012). "Efficient Bayes-Adaptive Reinforcement Learning using Sample-Based Search". In: *Advances in Neural Information Processing Systems*. Ed. by F. Pereira et al. Vol. 25. Curran Associates, Inc.

Results

The optimal policy exact cost: 888.89

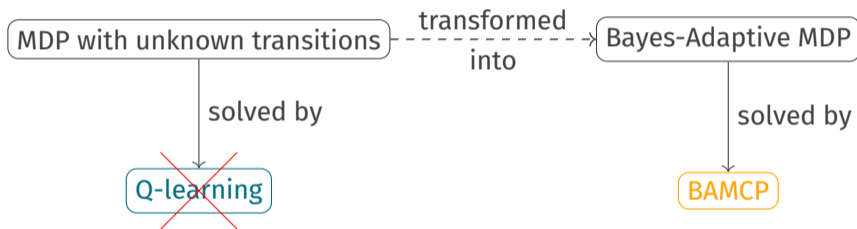
Simulated patients	Q-learning		BAMCP	
	Cost	Time	Cost	Time
10^2	1427.06 ± 1.05	0.07 sec	1302.58 ± 1.32	2.07 hours
10^3	936.96 ± 0.70	2.48 min	1297.64 ± 1.32	2.22 hours
10^4	936.93 ± 0.70	4.17 min	NC	4 days
10^6	891.6 ± 0.68	10.21 min	NC	1.5 years

Conclusion



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

Conclusion



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

Perspectives

MDP model	→	PDMP ⁷ model
Finite state space	→	Continuous state space
Markovian	→	Semi-Markovian
Complete observations	→	Hidden observations

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

⁷Mark H. A. Davis (1984). "Piecewise-Deterministic Markov Processes: A General Class of Non-Diffusion Stochastic Models". In: *Journal of the Royal Statistical Society Series B (Methodological)* 46, pp. 353–376.