

# Cours 3 et 4 Processus décisionnels Markoviens

Patrice PERNY



Université Paris 6

Patrice.Perny@lip6.fr

1

Patrice Perny

MADI – Cours 3

## 1. PROCESSUS DÉCISIONNELS MARKOVIENS

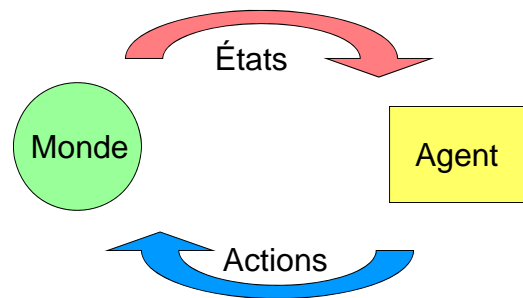
2

Patrice Perny

MADI – Cours 3

## PROCESSUS DÉCISIONNELS MARKOVIENS

Représentation d'une interaction synchrone  
entre un agent et le monde



Planification des actions d'un agent dans l'incertain

3

Patrice Perny

MADI – Cours 3

## MULTIPLES APPLICATIONS

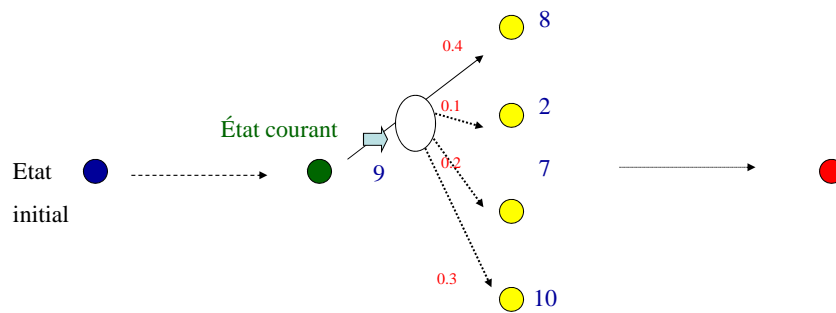


*Robotique* : sélection de l'action, planification  
*Gestion de Stocks* : planification des commandes  
*Jeux* : calcul de stratégies  
...

ROADEF'05 – Tours

MADI – Cours 3

## REPRÉSENTATION DANS UN GRAPHE D'ÉTAT



## FORMALISATION D'UN PDM

$$PDM = \langle S, A, T, R \rangle$$

$S$  : un ensemble fini d'états du monde

$A$  : un ensemble fini d'actions

$T : S \times A \rightarrow L(S)$  fonction de transition  
 $(s, a) \mapsto T(s, a)$  loi de proba sur les états de la nature

Abus de notation :  $T(s, a, s') = T(s, a)(s')$ ,  $s' \in S$

$R : S \times A \rightarrow \mathbb{R}$   
 $(s, a) \mapsto R(s, a)$  récompense immédiate

## DÉCISIONS ET POLITIQUES

**Règles de décision :** « si l'état est  $s$  alors exécuter l'action  $a$  »

Représentation par une fonction de décision

$$d : S \rightarrow A$$

$$s \mapsto a = d(s) \text{ action choisie dans l'état } s$$

A l'instant  $n + 1 - t$  : états possibles  $S_t$ , actions possibles  $A_t$

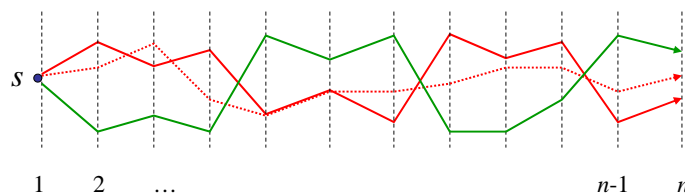
$$d_t : S_t \rightarrow A_t$$

$$s \mapsto a = d_t(s) \text{ action choisie à l'instant } n + 1 - t$$

**Hypothèse :** *observabilité totale* (on connaît l'état courant)

## DÉCISION DYNAMIQUE

$n$  étapes de décision ( $n$  = horizon, fini ou infini)



**politique** = séquence de fonctions de décisions  $(d_n, \dots, d_1)$

stratégie stationnaire :  $(d, d, \dots, d)$

## CRITÈRES À OPTIMISER

Horizon fini :

$$E \left( \sum_{t=1}^N R(s, d_t(s)) \right)$$

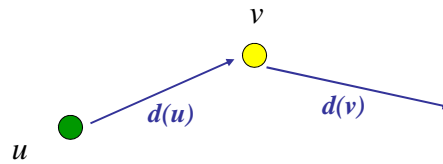
Horizon infini :

$$E \left( \sum_{t=1}^N \gamma^t R(s, d_t(s)) \right)$$

## 2. RÉOLUTION DE PDMs

## ÉVALUATION LOCALE D'UNE POLITIQUE

## Le cas déterministe



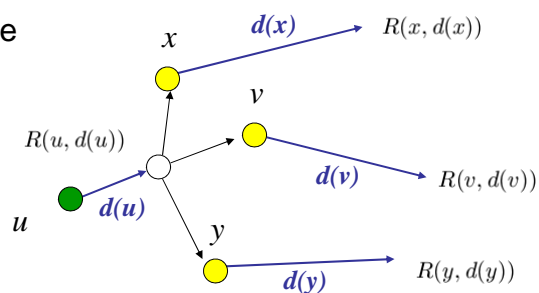
à horizon 1

à horizon 2

récompense :  $R(u, d(u)) + R(v, d(v))$ 

## ÉVALUATION LOCALE D'UNE POLITIQUE

## Le cas non-déterministe



à horizon 1

 $V_d(u) = R(u, d(u))$ 

à horizon 2

 $V_d(u) = R(u, d(u))$ 

$$+ \gamma [T(u, d(u), x)R(x, d(x)) + T(u, d(u), v)R(v, d(v)) + T(u, d(u), y)R(y, d(y))]$$
facteur de réduction :  $0 < \gamma \leq 1$  ( $\gamma < 1$  si horizon infini)

## EVALUATION LOCALE D'UNE POLITIQUE

## Dernière décision

$$V_{d,1}(s) = R(s, d_1(s))$$

## Décision à t étapes de la fin

$$V_{d,t}(s) = R(s, d_t(s)) + \gamma \sum_{s' \in S} T(s, d_t(s), s') V_{d,t-1}(s')$$

## CALCUL DES DÉCISIONS OPTIMALES (HORIZON FINI)

## Dernière décision

$$V_1^*(s) = \max_a R(s, a) \quad d_1^*(s) = \arg \max_a R(s, a)$$

## Décision à t étapes de la fin

$$V_t^*(s) = \max_a \left[ R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V_{t-1}^*(s') \right]$$

$$d_t^*(s) = \arg \max_a \left[ R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V_{t-1}^*(s') \right]$$

$\gamma = 1$

## ALGORITHME 1 : INDUCTION ARRIÈRE

Résolution d'un PDM à horizon fini  $N$ **Algorithm 1:** Backward Induction

---

```


foreach  $s \in S$  do  $V_0^*(s) = 0$ 
for  $t \leftarrow 1$  to  $N$  do
  for  $s \in S$  do
    for  $a \in A$  do
       $Q_t^a(s) = R(s, a) + \sum_{s' \in S} T(s, a, s') V_{t-1}^*(s')$ 
    end
     $a^* \leftarrow \text{choice}[\text{argmax}_{a \in A} Q_t^a(s)]$ 
     $d_t(s) \leftarrow a^*$ 
     $V_t^*(s) \leftarrow Q_t^{a^*}(s)$ 
  end
end
return  $d_t(s), V_t^*(s)$  for all  $t = 1, \dots, N, s \in S$ 

```

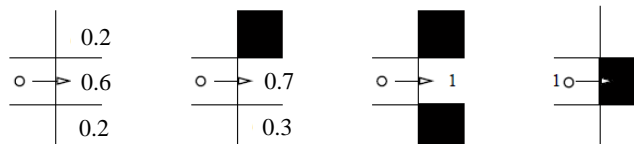
---

choice (X) : choix d'un élément dans X

## EXEMPLE

	1	2	3
1			
2			0.5 G
3		0.5 G	1 G

Possibilités de transition







## RÉSOLUTION D'UN PDM À HORIZON INFINI

Évaluation d'une stratégie stationnaire :  $(d, d, \dots, d)$

$$V_d(s) = R(s, d(s)) + \gamma \sum_{s' \in S} T(s, d(s), s') V_d(s'), \quad s \in S$$

**Théorème** (Howard, 1960)

Il existe une stratégie stationnaire optimale pour tout état initial.

La valeur de cette stratégie est caractérisée par le système:

$$V^*(s) = \max_a \left[ R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V^*(s') \right], \quad s \in S$$

qui admet une solution unique.

Equations de Bellman

## FORMULATION VECTORIELLE ET APPROCHE ITÉRATIVE

$$\forall s \in S, LV(s) = \max_{a \in A} \left\{ R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V(s') \right\}$$

$$LV = \max_{\pi} \{ R_{\pi} + \gamma T_{\pi} V \}$$

Equations de Bellman :  $V = LV$

Proposition :  $\|LV - LV'\| \leq \gamma \|V - V'\|$

Théorème du point fixe de Banach  $\Rightarrow$   
la suite  $\{V_n\}$  définie par :  $V_0 = 0$ ,  $V_t = LV_{t-1}$  for all  $t \geq 1$   
converge vers la solution  $V^*$  de l'équation  $V = LV$

$\rightarrow$  calculer  $V_0, V_1, \dots, V_t$  jusqu'à  $\|V_t - V_{t-1}\| < \varepsilon$

## ALGORITHME 2 : ITÉRATION DE LA VALEUR

## Algorithme de l'itération de la valeur (Bellman, 57)

## Algorithm 2: Value Iteration

---

```

foreach  $s \in S$  do  $V_0(s) = 0$ 
 $t \leftarrow 0$ 
repeat
   $t \leftarrow t + 1$ 
  for  $s \in S$  do
    for  $a \in A$  do
       $Q_t^a(s) = R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V_{t-1}(s')$ 
    end
     $V_t(s) \leftarrow \max_{a \in A} Q_t^a(s)$ 
  end
until  $\max_{s \in S} \{|V_t(s) - V_{t-1}(s)|\} < \varepsilon$ ;
foreach  $s \in S$  do  $d(s) \leftarrow \text{choice}[\arg\max_{a \in A} Q_t^a(s)]$ 
return  $d(s), V_t(s)$  for all  $s \in S$ 

```

---

*Amélioration* : Gauss-Seidel (utiliser  $V_t(s')$  à la place de  $V_{t-1}(s')$  si déjà calculé)

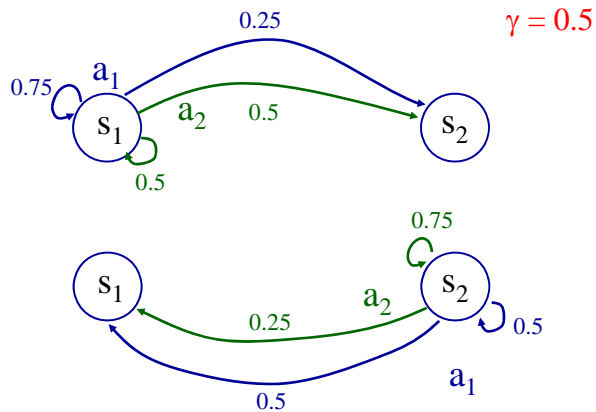
## GARANTIE DE PERFORMANCE

## Borne de l'erreur (Williams and Baird, 1993)

Si  $|V_t(s) - V_{t-1}(s)| < \epsilon$  pour tout  $s \in S$  alors :

$$\max_{s \in S} |V_{d_{V_t}}(s) - V^*(s)| < 2\epsilon \frac{\gamma}{1 - \gamma}$$

## EXEMPLE



Rewards	s1	s2
a1	8	11
a2	12	9

## ITÉRATION DE LA VALEUR

$$v_1 = \max\{8 + \frac{3}{8}v_1 + \frac{1}{8}v_2, 12 + \frac{1}{4}v_1 + \frac{1}{4}v_2\}$$

$$v_2 = \max\{11 + \frac{1}{4}v_1 + \frac{1}{4}v_2, 9 + \frac{1}{8}v_1 + \frac{3}{8}v_2\}$$

v1	v2	Q11	Q12	Q21	Q22
0	0				
12,00	11,00	8,00	12,00	11,00	9,00
17,75	16,75	13,88	17,75	16,75	14,63
20,63	19,63	16,75	20,63	19,63	17,50
22,06	21,06	18,19	22,06	21,06	18,94
22,78	21,78	18,91	22,78	21,78	19,66
23,14	22,14	19,27	23,14	22,14	20,02
23,32	22,32	19,45	23,32	22,32	20,20
23,41	22,41	19,54	23,41	22,41	20,29
23,46	22,46	19,58	23,46	22,46	20,33
23,48	22,48	19,60	23,48	22,48	20,35
23,49	22,49	19,61	23,49	22,49	20,36
23,49	22,49	19,62	23,49	22,49	20,37
23,50	22,50	19,62	23,50	22,50	20,37
23,50	22,50	19,62	23,50	22,50	20,37

## L'EXEMPLE DU ROBOT AVEC HORIZON INFINI

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
1	Probas							Bas				Droite				Altitude	
2																	
3	0,2	1	0,00		0,00		-1,00		-1,00		0,00		0,00		1		3
4	0,6		0,00	0,00	0,00		0,00	33,00	14,00		-1,00	9,00	0,00		2	3	4
5	0,2			0,00	100,00			0,00	100,00			15,00	100,00			5	3
6																	
7		2	0,00		0,00		10,55		12,65		0,00		0,00				
8	0,3		0,00	33,00	14,00		0,00	35,25	19,25		9,20	13,20	7,00				
9	0,7			15,00	100,00			7,50	100,00			19,90	100,00				
10																	
11	Gamma	3	10,55		12,65		12,72		14,23		5,28		6,33				
12	0,5		9,20	35,25	19,25		4,60	35,99	20,97		11,25	16,04	9,63				
13				19,90	100,00			9,95	100,00			21,74	100,00				
14																	
15		4	12,72		14,23		13,28		14,74		6,36		7,11				
16			11,25	35,99	20,97		5,63	36,26	21,61		12,01	16,71	10,48				
17				21,74	100,00			10,87	100,00			22,34	100,00			Politique optimale	
18																	
19		5	13,28		14,74		13,49		14,93		6,64		7,37		v		v
20			12,01	36,26	21,61		6,00	36,35	21,82		12,26	16,96	10,80		>	v	v
21				22,34	100,00			11,17	100,00			22,56	100,00			>	>
22																	

ROADEF'05 – Tours

MADI – Cours 3

## L'EXEMPLE DU ROBOT AVEC HORIZON INFINI

5	13,28		14,74		13,49		14,93		6,64		7,37		v		v
	12,01	36,26	21,61		6,00	36,35	21,82		12,26	16,96	10,80		>	v	v
		22,34	100,00			11,17	100,00			22,56	100,00			>	>
6	13,49		14,93		13,56		15,00		6,75		7,47				
	12,26	36,35	21,82		6,13	36,38	21,90		12,35	17,04	10,91				
		22,56	100,00			11,28	100,00			22,64	100,00				
7	13,56		15,00		13,59		15,02		6,78		7,50				
	12,35	36,38	21,90		6,17	36,40	21,92		12,38	17,07	10,95				
		22,64	100,00			11,32	100,00			22,66	100,00				
8	13,59		15,02		13,60		15,03		6,79		7,51				
	12,38	36,40	21,92		6,19	36,40	21,93		12,39	17,08	10,96				
		22,66	100,00			11,33	100,00			22,67	100,00				

ROADEF'05 – Tours

MADI – Cours 3

## AUTRE APPROCHE : ITÉRATION DE LA POLITIQUE

*Proposition* : soit  $\pi$  une politique stationnaire et  $V_\gamma^\pi$  sa valuation en chaque état. Alors toute politique  $\pi'$  choisie dans  $\operatorname{argmax}_\delta \{R_\delta + \gamma T_\delta V_\gamma^\pi\}$  vérifie l'inégalité  $V_\gamma^{\pi'} \geq V_\gamma^\pi$

$\Rightarrow$  idée d'un autre algo itératif

Partir avec une politique arbitraire  $\pi_0$

résoudre le système  $V_0 = L_{\pi_0} V_0$

choisir  $\pi_1$  dans  $\operatorname{argmax}_\delta \{R_\delta + \gamma T_\delta V_0\}$

etc...

$$V_0 \leq V_1 \leq \dots \leq V_t$$

## ALGORITHME 3 : ITÉRATION DE LA POLITIQUE

---

**Algorithm 3:** Policy Iteration
 

---

choose an initial decision rule  $d_0$

$t \leftarrow 0$

**repeat**

résoudre

$$\{V_t(s) = R(s, d_t(s)) + \gamma \sum_{s' \in S} T(s, d_t(s), s') V_t(s'), \quad \forall s \in S\}$$

**for**  $s \in S$  **do**

$$d_{t+1}(s) \leftarrow \operatorname{choice}[\operatorname{argmax}_{a \in A} \{R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V_t(s')\}]$$

**end**

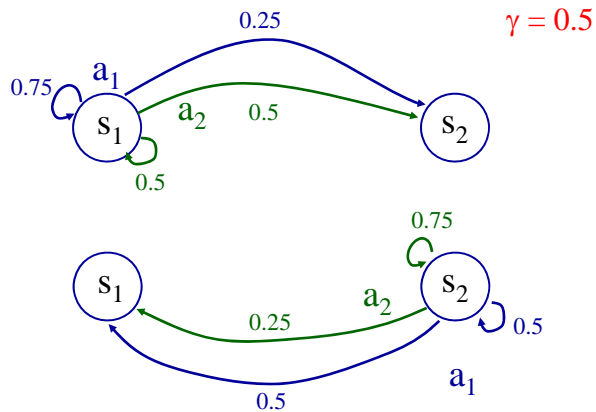
$t \leftarrow t + 1$

**until**  $d_t(s) = d_{t+1}(s), \forall s \in S;$

**return**  $d_{t+1}(s), V_t(s)$  for all  $s \in S$

---

## EXEMPLE



Rewards	s1	s2
a1	8	11
a2	12	9

## EXEMPLE

Evaluation de la politique :  $a_1$  si  $s_1$  et  $a_2$  si  $s_2$

Résoudre le système :

$$v_1 = 8 + \frac{3}{8}v_1 + \frac{1}{8}v_2$$

$$v_2 = 9 + \frac{1}{8}v_1 + \frac{3}{8}v_2$$

Utiliser la solution pour mettre à jour la politique courante

Itérer jusqu'à stabilité...

## FORMULATION PAR PL

*Proposition :*  $V \geq LV \Rightarrow V \geq V^*$

Donc on peut résoudre le problème d'optimisation

$$\min \sum_{s \in S} V(s)$$

$$\text{s.c. } V \geq LV$$

## RÉSOLUTION PAR PL

---

**Algorithm 4:** Linear Programming Formulation
 

---

résoudre  $\min_V \sum_{s \in S} V(s)$

s.c.

$$V(s) \geq R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V^*(s'), \quad \forall s \in S, \forall a \in A$$

**foreach**  $s \in S$  **do**

$$d(s) \leftarrow \text{choice}[\text{argmax}_{a \in A} \{R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V(s')\}]$$

**end**

**return**  $d(s), V(s)$  for all  $s \in S$

---

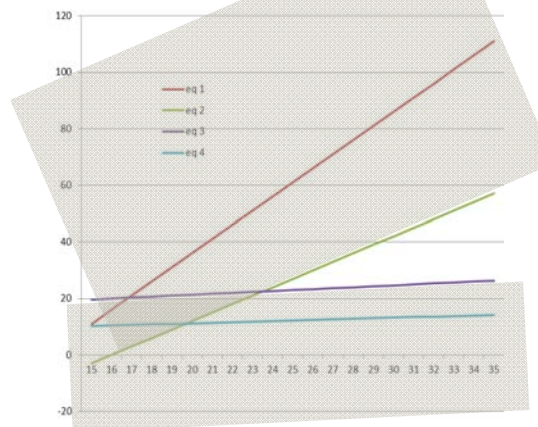


## EXEMPLE

$$v_1 = \max\{8 + \frac{3}{8}v_1 + \frac{1}{8}v_2, 12 + \frac{1}{4}v_1 + \frac{1}{4}v_2\}$$

$$v_2 = \max\{11 + \frac{1}{4}v_1 + \frac{1}{4}v_2, 9 + \frac{1}{8}v_1 + \frac{3}{8}v_2\}$$

$$\begin{aligned} & \min v_1 + v_2 \\ \text{s.t.} \quad & \begin{cases} v_1 \geq 8 + \frac{3}{8}v_1 + \frac{1}{8}v_2 & (1) \\ v_1 \geq 12 + \frac{1}{4}v_1 + \frac{1}{4}v_2 & (2) \\ v_2 \geq 11 + \frac{1}{4}v_1 + \frac{1}{4}v_2 & (3) \\ v_2 \geq 9 + \frac{1}{8}v_1 + \frac{3}{8}v_2 & (4) \end{cases} \end{aligned}$$



33

Patrice Perny

MADI – Cours 3

PDMs classiques

## RÉSOLUTION PAR PL



Deutsch | Support | Contact Us

Search website or documentation

PRODUCTS SUPPORT RESOURCES DOWNLOAD COMPANY

### AN EASIER WAY TO BETTER DECISIONS

#### Gurobi Optimizer 5.6 - State-of-the-Art Mathematical Programming Solver

- Superior optimization algorithms for faster times to feasibility and optimality
- Flexible interfaces and modeling language support for maximum productivity
- Great support from easy-to-reach optimization experts
- Transparent pricing and flexible licensing so no surprises when it's time to deploy

#### ACCOUNT LOGIN

email password

Register for Free | Forgot Password

#### HOW MAY WE HELP YOU

We are happy to set you up with a free trial, provide pricing information, or discuss your other needs.

CONTACT US

FREE EVALUATION

GUROBI OPTIMIZER 5.6  
High-end Optimization Libraries

ABOUT GUROBI  
and why so many are choosing us

SWITCHING TO GUROBI  
We've made migrating models easy

34

Patrice Perny

MADI – Cours 3

## RÉSOLUTION PAR PL AVEC GUROBI

Optimize a model with 4 rows, 2 columns and 8 nonzeros

Presolve time: 0.00s

Presolved: 4 rows, 2 columns, 8 nonzeros

Iteration	Objective	Primal Inf.	Dual Inf.	Time
0	0.000000e+00	4.000000e+01	0.000000e+00	0s
2	4.600000e+01	0.000000e+00	0.000000e+00	0s

Solved in 2 iterations and 0.00 seconds

Optimal objective 4.60000000e+01

Valeurs de la politique optimale:

<gurobi.Var v1 (value 23.5)>

<gurobi.Var v2 (value 22.5)>

Valeur de la fonction objectif : 46.0

35

Patrice Perny

MADI – Cours 3

### 3. APPRENTISSAGE PAR RENFORCEMENT DANS LES PDMS

ROADEF'05 – Tours

MADI – Cours 3

## LE Q-LEARNING

### CADRE DES MDPs :

Itération de la valeur, de la politique, PL  $\rightarrow$   
 $v(s), Q(a, s), \pi$  politique optimale

Imaginons que  $T$  ou  $R$  n'est pas connu, deux approches :

- approche indirecte (model based) : estimer  $T$  et  $R$  puis résoudre le MDP
- approche directe (model free) : on va directement essayer d'estimer  $Q$ , c'est le Q-learning

### CADRE DU Q-LEARNING :

Essayer des actions dans des états, observer, pour apprendre  
 $Q(a, s), v(s), \pi$

Episodes d'observation :  $(s, a, r, s', a, r', s'', a'', r'', s''', \dots)$

- partir avec des estimations initiales  $\hat{Q}(a, s)$  des  $Q(a, s)$
- mettre à jour les valeurs  $\hat{Q}(a, s)$  en fonction des observations faites

37

Patrice Perny

MADI – Cours 3

## LE Q-LEARNING

Dans le modèle MDP :

$$Q(a, s) = R(a, s) + \gamma \sum_{s' \in S} T(s, a, s') \max_{a'} Q(s, a')$$

Mais ici on ne connaît plus  $T$  et/ou  $R$

Imaginons que l'on fasse l'observation  $(s, a, s', r)$

$\rightarrow$  nouvelle estimation de  $\hat{Q}(s, a) : r + \gamma \max_{a'} \hat{Q}(s', a')$

Incorporation de cette nouvelle estimation :

$$\hat{Q}(s, a) \leftarrow (1 - \alpha_t) \hat{Q}(s, a) + \alpha_t [r + \gamma \max_{a'} \hat{Q}(s', a')]$$

avec  $\alpha_t = 1/t$  où une fonction de ce type qui décroît avec  $t$   
( $\alpha_t$  : taux d'apprentissage)

$$\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \alpha_t [r + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)]$$

38

Patrice Perny

MADI – Cours 3

## L'ALGORITHME DE Q-LEARNING

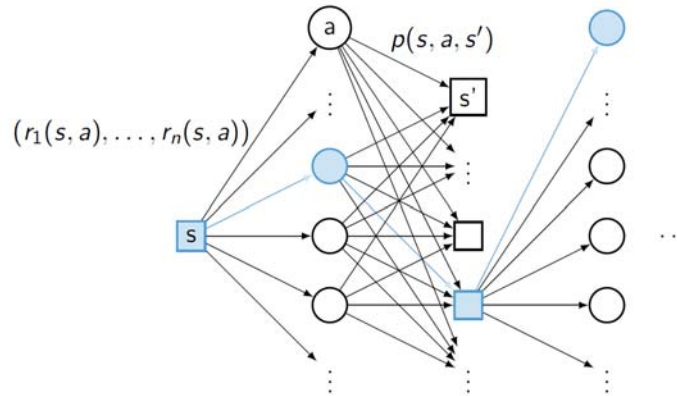
(Watkins, 1989)

```
Initialize  $Q(s, a)$  arbitrarily
Repeat (for each episode):
  Initialize  $s$ 
  Repeat (for each step of episode):
    Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Take action  $a$ , observe  $r, s'$ 
     $Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$ 
     $s \leftarrow s'$ 
  until  $s$  is terminal
```

## 4. PROCESSUS DÉCISIONNELS MARKOVIENS MULTIOBJECTIFS

## PROCESSUS MARKOVIENS MULTIOBJECTIFS

Minimiser distance, temps, consommation d'énergie, risque encouru...



41

Patrice Perny

MADI – Cours 3

## FORMULATIONS PL POUR UN SEUL OBJECTIF

### Programme primal

$$(\mathcal{P}) \begin{cases} \min \sum_{s \in S} \mu(s) v(s) \\ \text{s.t. } v(s) - \gamma \sum_{s' \in S} T(s, a, s') v(s') \geq R(s, a) \quad \forall s \in S, \forall a \in A \end{cases}$$

### Programme dual

$$(\mathcal{D}) \begin{cases} \max \sum_{s \in S} \sum_{a \in A} R(s, a) x_{sa} \\ \text{s.t. } \sum_{a \in A} x_{sa} - \gamma \sum_{s' \in S} \sum_{a \in A} T(s', a, s) x_{s'a} = \mu(s) \quad \forall s \in S \\ x_{sa} \geq 0 \quad \forall s \in S, \forall a \in A \end{cases}$$

42

Patrice Perny

MADI – Cours 3

## CHARACTÉRISATION DES POLITIQUES MIXTES

$$\begin{aligned} \text{s.t. } & \sum_{a \in A} x_{sa} - \gamma \sum_{s' \in S} \sum_{a \in A} T(s', a, s) x_{s'a} = \mu(s) \quad \forall s \in S \\ & x_{sa} \geq 0 \quad \forall s \in S, \forall a \in A \end{aligned}$$

**Proposition 1** For a policy  $\pi$ , if  $x^\pi$  is defined as  $x^\pi(s, a) = \sum_{t=0}^{\infty} \gamma^t p_t^\pi(s, a)$ ,  $\forall s \in S, \forall a \in A$  where  $p_t^\pi(s, a)$  is the probability of reaching state  $s$  and choosing  $a$  at step  $t$ , then  $x^\pi$  is a feasible solution of  $\mathcal{D}$ .

**Proposition 2** If  $x_{sa}$  is a solution of  $\mathcal{D}$ , then the stationary randomized policy  $\delta^\infty$ , defined by  $\delta(s, a) = x_{sa} / \sum_{a' \in A} x_{sa'}$ ,  $\forall s \in S, \forall a \in A$  defines  $x^{\delta^\infty}(s, a)$  as in Proposition 1, that are equal to  $x_{sa}$ .

Les solutions du polyèdre des contraintes caractérisées par les variables  $x_{sa}$  sont les politiques mixtes du MDP

43

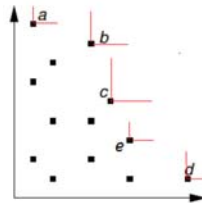
Patrice Perny

MADI – Cours 3

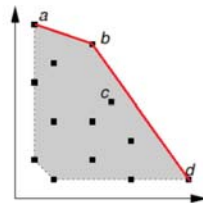
## PL ET DUALITÉ POUR LES MOMDPS

$$\max f_i(x) = \sum_{s \in S} \sum_{a \in A} R_i(s, a) x_{sa} \quad \forall i = 1, \dots, n$$

$$\begin{aligned} \text{s.t. } & \sum_{a \in A} x_{sa} - \gamma \sum_{s' \in S} \sum_{a \in A} T(s', a, s) x_{s'a} = \mu(s) \quad \forall s \in S \\ & x_{sa} \geq 0 \quad \forall s \in S, \forall a \in A \end{aligned}$$



politiques pures



politiques mixtes

Pour forcer des politiques pures

$$\begin{aligned} \sum_{a \in A} d_{sa} &\leq 1 \quad \forall s \in S \\ (1 - \gamma)x_{sa} &\leq d_{sa} \quad \forall s \in S, \forall a \in A \\ d_{sa} &\in \{0, 1\} \quad \forall s \in S, \forall a \in A \end{aligned}$$

$$x_{sa} = \sum_{t=0}^{\infty} \gamma^t p_t^\pi(s, a) \leq \frac{1}{1-\gamma}$$

44

Patrice Perny

MADI – Cours 3

## POLITIQUE MAX-MIN OPTIMALE

$$\begin{aligned}
 & \max z \\
 & z \leq f_i(x), i = 1, \dots, n \\
 & f_i(x) = \sum_{s \in S} \sum_{a \in A} R_i(s, a) x_{sa} \quad \forall i = 1, \dots, n \\
 & \text{s.t.} \quad \sum_{a \in A} x_{sa} - \gamma \sum_{s' \in S} \sum_{a \in A} T(s', a, s) x_{s'a} = \mu(s) \quad \forall s \in S \\
 & x_{sa} \geq 0 \quad \forall s \in S, \forall a \in A
 \end{aligned}$$

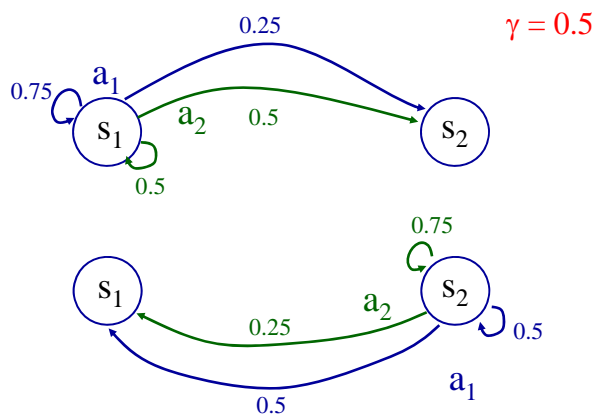
45

Patrice Perny

MADI – Cours 3

## PDMs classiques

### EXEMPLE



R1	s1	s2
a1	8	11
a2	12	9

R2	s1	s2
a1	13	7
a2	6	15

46

Patrice Perny

MADI – Cours 3

## LES CONTRAINTES CHANGENT ENTRE R1 ET R2

$$\begin{array}{ll} \min \frac{1}{2}v_1 + \frac{1}{2}v_2 & \min \frac{1}{2}v_1 + \frac{1}{2}v_2 \\ \left\{ \begin{array}{l} \frac{5}{8}v_1 - \frac{1}{8}v_2 \geq 8 \\ \frac{3}{4}v_1 - \frac{1}{4}v_2 \geq 12 \\ -\frac{1}{4}v_1 + \frac{3}{4}v_2 \geq 11 \\ -\frac{1}{8}v_1 + \frac{5}{8}v_2 \geq 9 \end{array} \right. & \left\{ \begin{array}{l} \frac{5}{8}v_1 - \frac{1}{8}v_2 \geq 13 \\ \frac{3}{4}v_1 - \frac{1}{4}v_2 \geq 6 \\ -\frac{1}{4}v_1 + \frac{3}{4}v_2 \geq 7 \\ -\frac{1}{8}v_1 + \frac{5}{8}v_2 \geq 15 \end{array} \right. \end{array}$$

47

Patrice Perny

MADI – Cours 3

## FORMES DUALES, LES OBJECTIFS CHANGENT

$$\begin{array}{l} \max 8x_{11} + 12x_{12} + 11x_{21} + 9x_{22} \\ \left\{ \begin{array}{l} \frac{5}{8}x_{11} + \frac{3}{4}x_{12} - \frac{1}{4}x_{21} - \frac{1}{8}x_{22} = \frac{1}{2} \\ -\frac{1}{8}x_{11} - \frac{1}{4}x_{12} + \frac{3}{4}x_{21} + \frac{5}{8}x_{22} = \frac{1}{2} \end{array} \right. \\ x_{11} \geq 0, x_{12} \geq 0, x_{21} \geq 0, x_{22} \geq 0, \end{array}$$


---


$$\begin{array}{l} \max 13x_{11} + 6x_{12} + 7x_{21} + 15x_{22} \\ \left\{ \begin{array}{l} \frac{5}{8}x_{11} + \frac{3}{4}x_{12} - \frac{1}{4}x_{21} - \frac{1}{8}x_{22} = \frac{1}{2} \\ -\frac{1}{8}x_{11} - \frac{1}{4}x_{12} + \frac{3}{4}x_{21} + \frac{5}{8}x_{22} = \frac{1}{2} \end{array} \right. \\ x_{11} \geq 0, x_{12} \geq 0, x_{21} \geq 0, x_{22} \geq 0, \end{array}$$

48

Patrice Perny

MADI – Cours 3



## OPTIMISATION MAX-MIN

max z

$$\begin{cases} z \leq 8x_{11} + 12x_{12} + 11x_{21} + 9x_{22} \\ z \leq 13x_{11} + 6x_{12} + 7x_{21} + 15x_{22} \\ \frac{5}{8}x_{11} + \frac{3}{4}x_{12} - \frac{1}{4}x_{21} - \frac{1}{8}x_{22} = \frac{1}{2} \\ -\frac{1}{8}x_{11} - \frac{1}{4}x_{12} + \frac{3}{4}x_{21} + \frac{5}{8}x_{22} = \frac{1}{2} \end{cases}$$

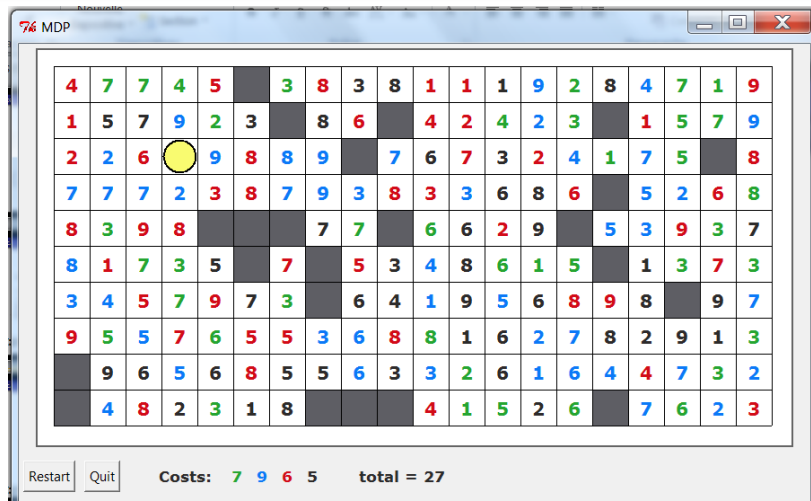
$$x_{11} \geq 0, x_{12} \geq 0, x_{21} \geq 0, x_{22} \geq 0,$$

49

Patrice Perny

MADI – Cours 3

## APPLICATION :



50

Patrice Perny

MADI – Cours 3