

# Processus de Décision Markoviens

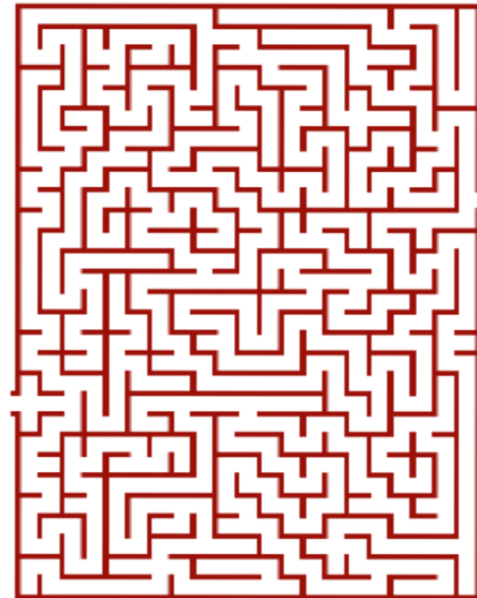
**Hélène Fargier** / (Philippe Leray)

`fargier@irit.fr`



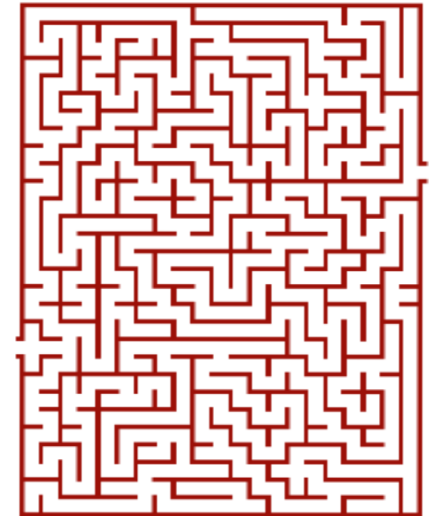
# Introduction

- Un problème de décision est souvent dynamique...
- Comment prendre en compte l'aspect séquentiel ?
- Ex : robot dans un labyrinthe



# Notations

- A l'instant  $t$ ,
  - l'environnement est dans l'état  $s_t$
  - l'agent décide d'une action  $a_t$
  - pour cette action, l'agent recevra une récompense immédiate  $r_{t+1}$   
(*cout de l'action + reward ou cout à etre dans l'etat suivant*)
  - l'action fait passer l'environnement dans l'état  $s_{t+1}$



- A l'instant  $t$ ,**
- $s_t$  = position du robot dans le labyrinthe
  - $a_t$  = aller à l'Ouest / Est / Nord / Sud
  - $r_t$  = -1 si la case n'est pas la sortie, +1 sinon

# Processus de Décision Markovien (MDP)

- Soient

- un ensemble fini d'états  $S$
- un ensemble fini d'actions  $A$

- A chaque instant  $t$  (discret), l'agent

- observe un état  $s \in S$
- choisit une action  $a \in A$

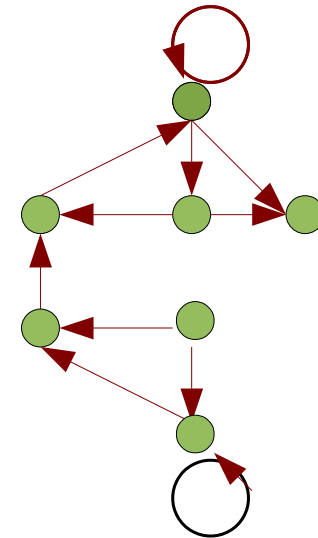
- Répercussions :

- l'environnement passe dans l'état (cas déterministe)
- l'agent recoit la récompense

$$s' = T(s; a)$$

$$r' = r(s; a) (= r(s') - c(s, a))$$

- **Hypothèse de Markov :  $T(s; a)$  et  $r(s; a)$  ne dépendent que de l'état et de l'action choisie**



# Objectif (cas déterministe)

**Cas déterministe** : appliquée dans un état donné, une action ne peut conduire qu'à un seul nouvel état

- Choisir une séquence d'actions  $\tau = (s_0, a_0, s_1, a_1, s_2, a_2, \dots)$  - un plan, ou « trajectoire »- qui va maximiser une fonction fixée
- A l'instant  $t$ , essayer de recevoir de futures récompenses  $r_{t+i}$  les meilleures possibles
- Récompense cumulée (à horizon fini):  $R(t) = \sum_{i=0, h} r_{t+i+1}$
- Récompense cumulée (à horizon infini) avec un terme d'oubli

$$R(t) = \sum_{i=0, \infty} \gamma^i \cdot r_{t+i+1}$$

*$\gamma^i$  représente la valeur  $t=i$  d'une unité de récompense reçue  $t=i+1$*

# Représenter les actions

- Actions déterministes:

$$T: S \times A \rightarrow S ;$$

Pour tout état et toute action T spécifieait un nouvel état

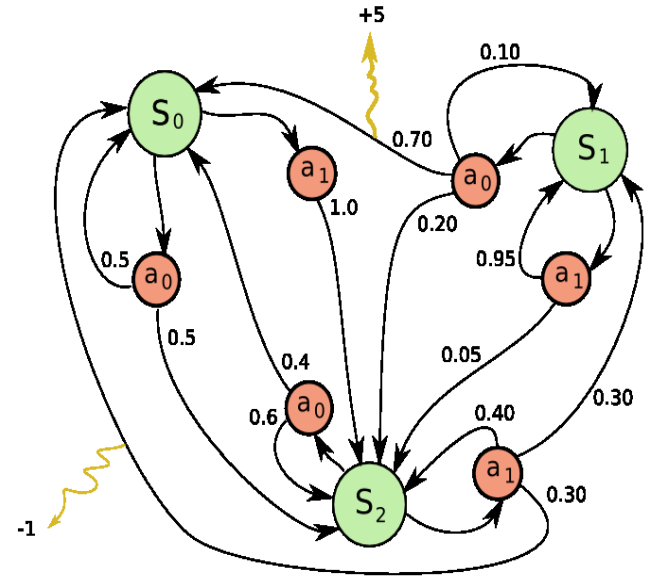
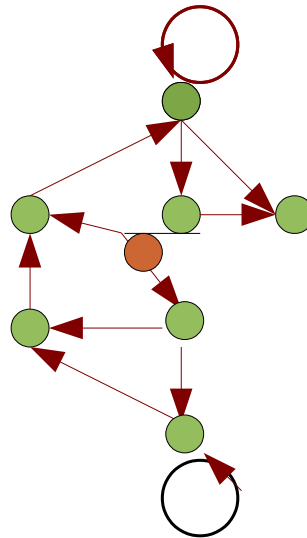
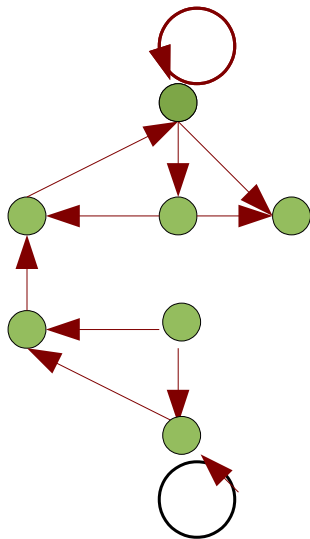
- Actions stochastiques :

$$T: S \times A \rightarrow P(S)$$

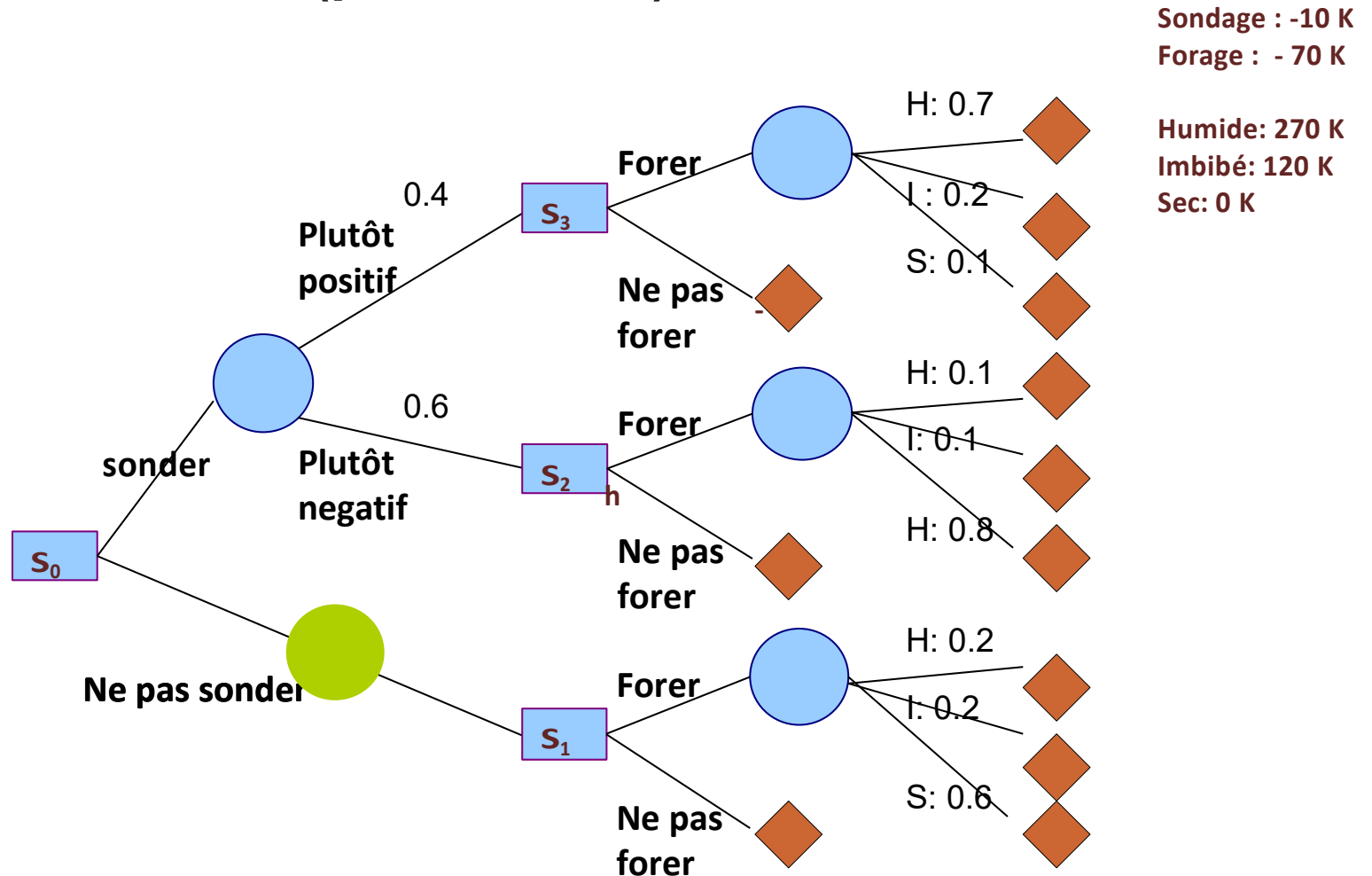
Pour tout état et tout action, T spécifie une distribution de probabilité sur les états suivants  $P(s'|s,a)$ .

**Hypothèse de Markov :  $T(s; a)$  et  $r(s; a)$  ne dépendent que de l'état et de l'action choisie**

# Automates déterministes, automates stochastiques



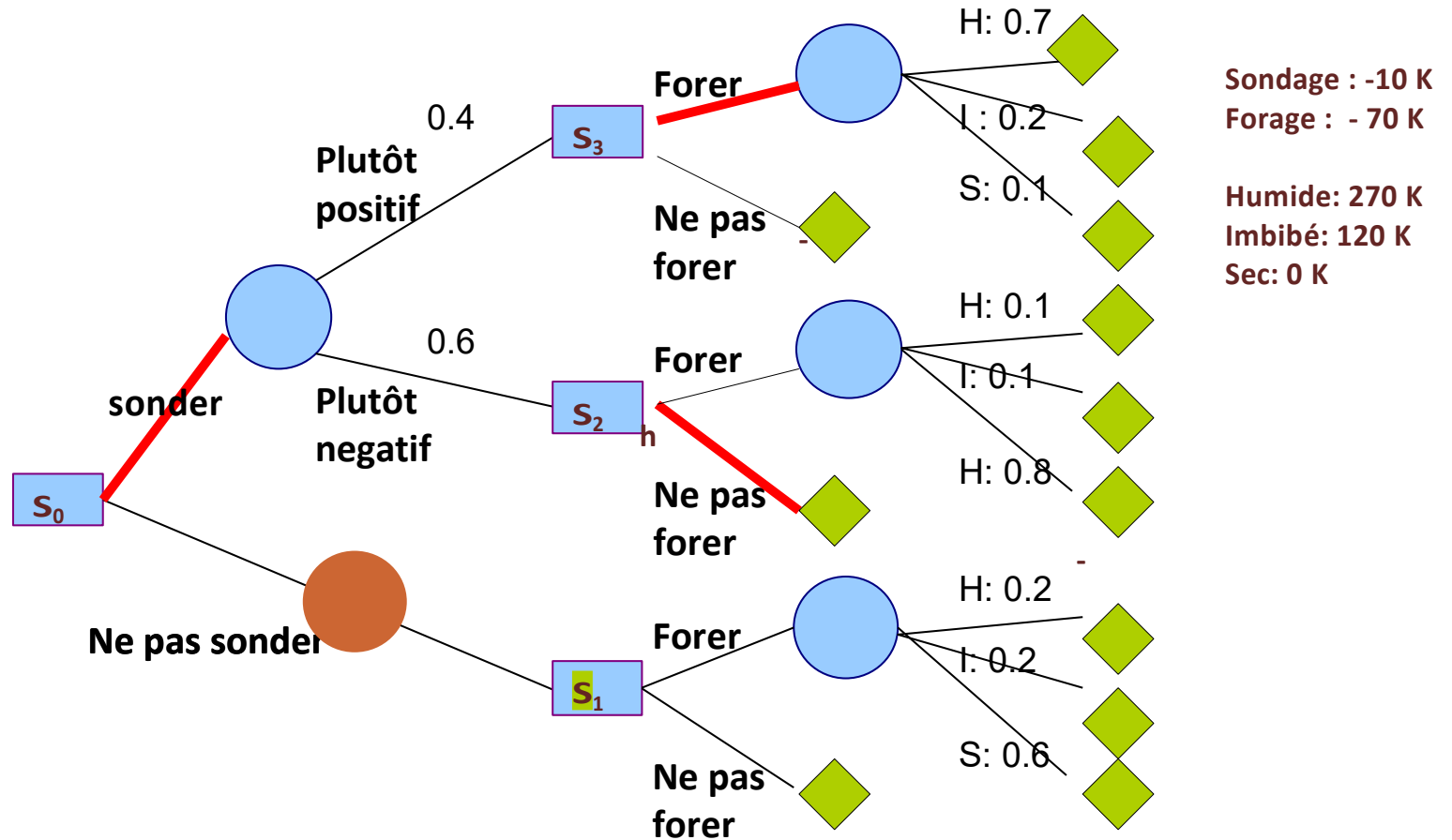
# Arbre de décision (probabiliste)



Sur chaque arc (Ci,N) issu d'un noeud chance :  $P(N|Ci, PRED(Ci))$

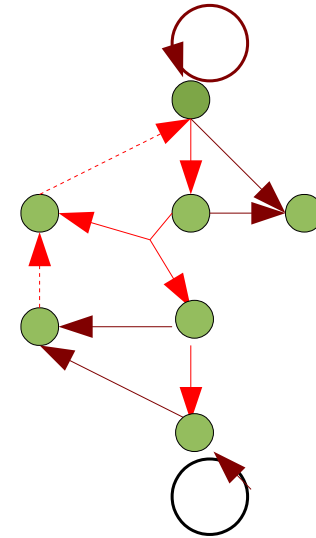


# Politique



# Représenter une solution

- Une politique  $\pi$  est une fonction de  $S$  dans  $A$ .
- Suivre une politique
  1. Déterminer l'état courant  $s$
  2. Exécuter l'action  $\pi(s)$
  3. Retourner en 1
- Suppose l'observabilité totale du système (on sait toujours où on est)



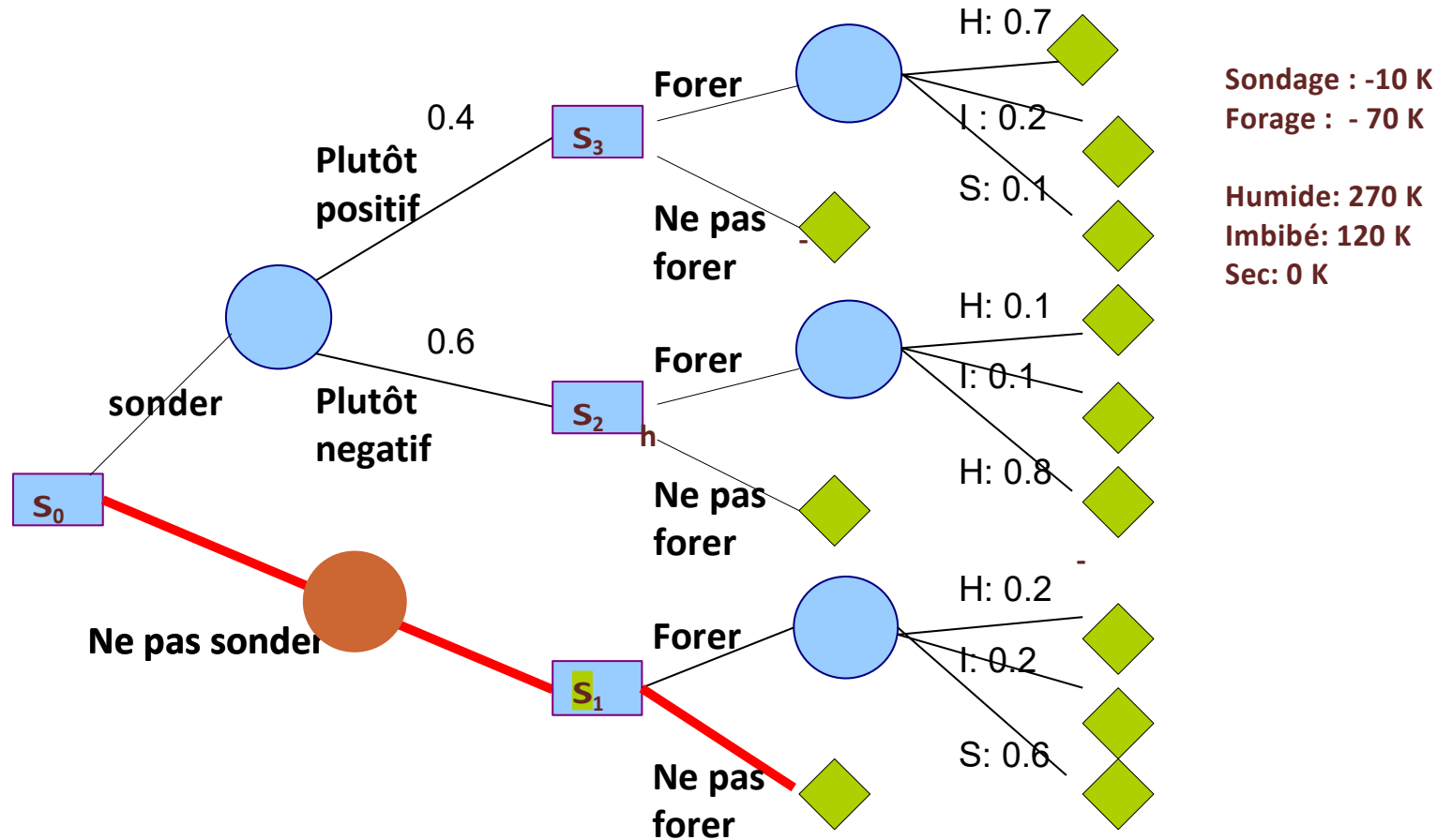
# Evaluer la qualité d'une politique (cas stochastique)

- Une politique d'action = Une fonction  $\pi : S \rightarrow A$ .
- Cas déterministe,  $\pi$  définit une trajectoire pour tout état initial  $s_0$  :

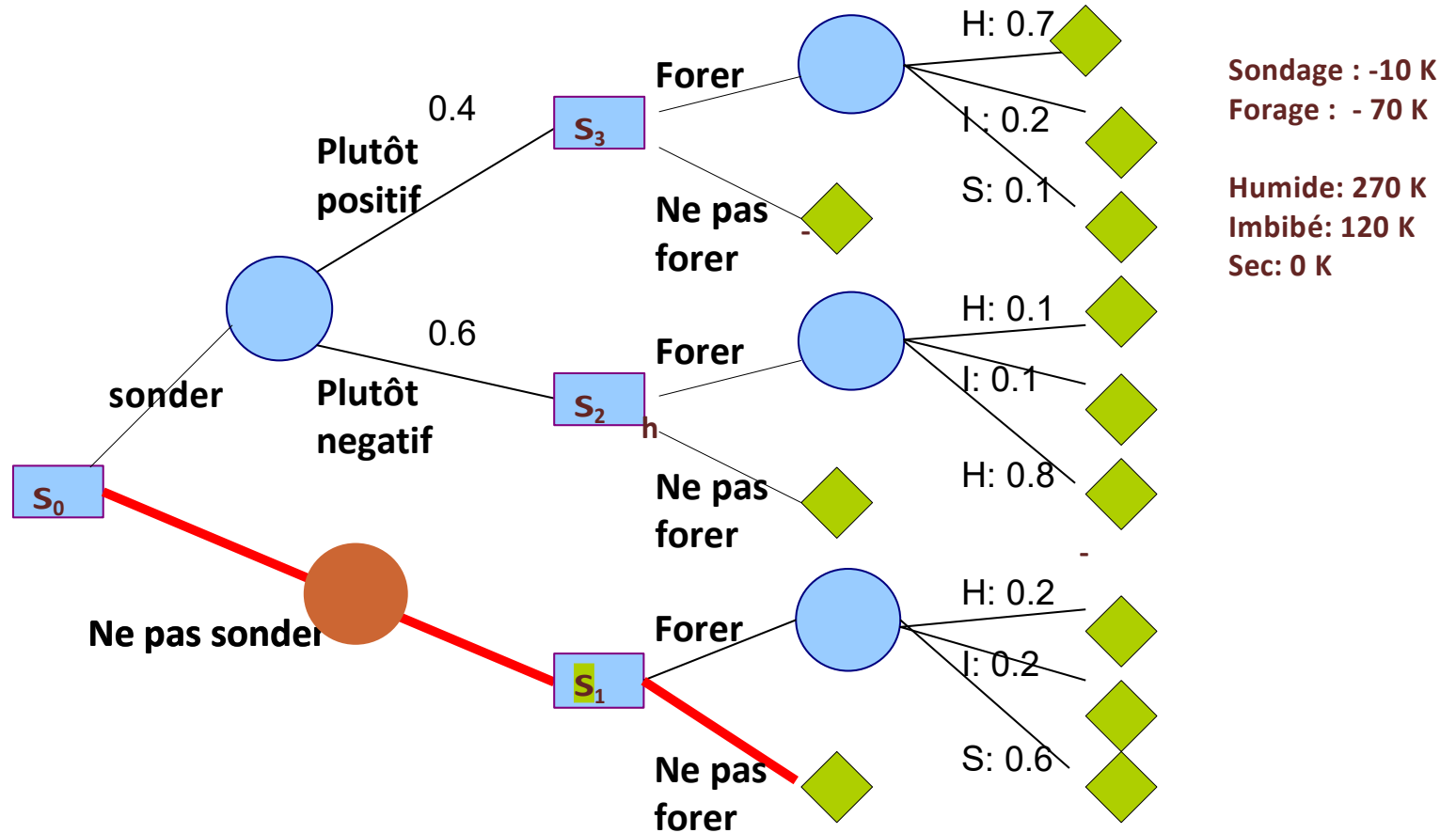
$$\tau(\pi, s_0) = (s_0, \pi(s_0), s_1, \pi(s_1), s_2, \pi(s_2), \dots)$$

On peut choisir la politique maximisant l'utilité  $R(t) = \sum_{i=0, h} r_{t+i+1}$  pour  $t=0$

# Politique



# Politique



# Evaluer la qualité d'une politique (cas stochastique)

- Une politique d'action = Une fonction  $\pi : S \rightarrow A$ .
- Cas déterministe,  $\pi$  définit une trajectoire pour tout état initial  $s_0$  :

$$\tau(\pi, s_0) = (s_0, \pi(s_0), s_1, \pi(s_1), s_2, \pi(s_2), \dots)$$

On peut choisir la politique maximisant l'utilité  $R(t) = \sum_{i=0, h} r_{t+i+1}$  pour  $t=0$

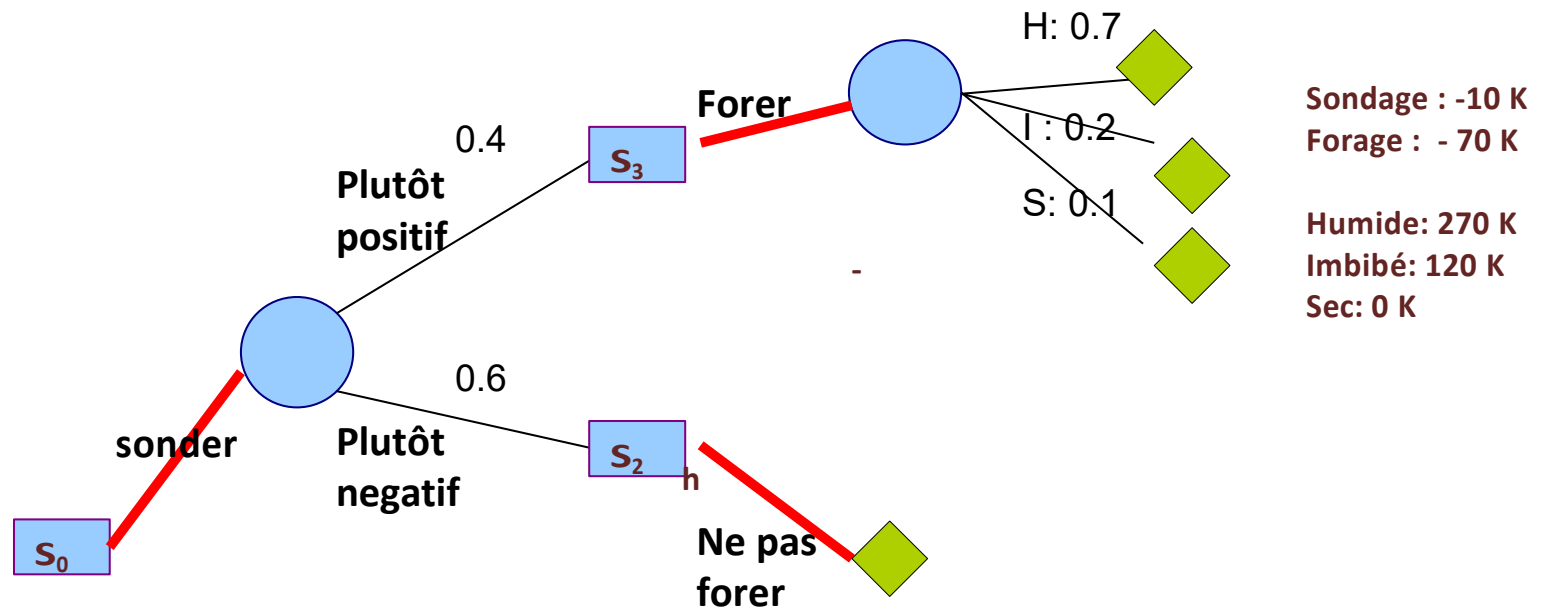
- Maintenant,  $\pi$  génère une distribution de probabilité sur les trajectoires

$$p(\tau|\pi, s_0) = \prod_{t=0, h-1} p(s_{i+1} \mid s_i, \pi(t, s_i))$$

A chaque trajectoire une utilité

**Politique optimale :  $\pi$  maximisant l'espérance mathématique de l'utilité**

# Politique → Distribution de Probabilité sur les utilités



# Equation de Bellman

- Equations de Bellman: EU pour chaque etat selon la politique

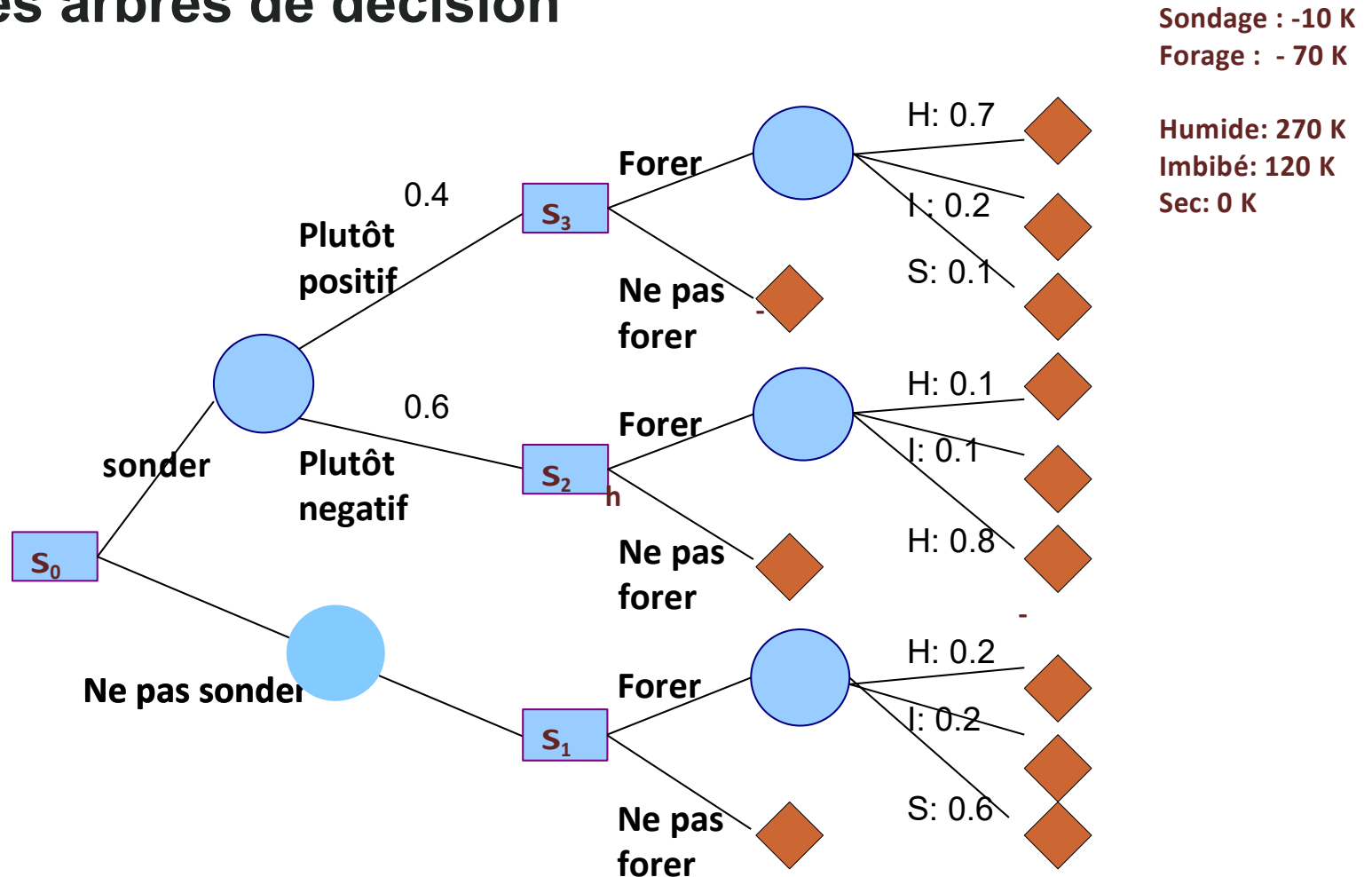
$$V_{\pi}(s) = r(s, \pi(s)) + \sum_{s' \in \mathcal{S}} P(s'|s, a) \cdot (\gamma) \cdot V_{\pi}(s')$$

Maximiser  $U(\pi) = V_{\pi}(s_0)$  ; en fait, maximiser  $V_{\pi}(s)$  pour chaque  $s$ ;

- Comment résoudre l'équation de Bellman ?
  - Pour chaque  $s$ , 1 équation linéaire à  $|\mathcal{S}|$  inconnues :  $|\mathcal{S}|$  à  $|\mathcal{S}|$  inconnues: Programmation linéaire ( $\mathcal{S}$  petit)
  - Faire des simulations : methode de Monte Carlo
  - Programmation dynamique



# Calcul d'une politique par programmation dynamique dans les arbres de décision



**En partant de la fin, maximiser l'utilité esperée**

# Calcul d'une politique par programmation dynamique dans les arbres de décision

- Partir de la fin : pour chaque  $s$ , choisir l'action qui maximise EU

Pour  $t = N-1$  à 0 faire

  Pour tout  $s$  de  $S$  faire

    Pour tout  $a$  de  $A(s)$  faire

$$U(s, a) = r(s, a) + \sum_{s' \in S} P(s' | s, a) \cdot U(s')$$

    Fin pour

$$\pi(s) = \operatorname{argmax}_a U(s, a)$$

$$U(s) = U(s, \pi(s))$$

  Fin pour

Fin pour

Algo fondé sur principe de monotonie + decomposition de l'EU

# Calcul d'une politique par programmation dynamique dans les arbres de décision

- **Monotonie :**

**toute sous politique d'une politique optimale est optimale**

$$EU(a) > EU(a') \Leftrightarrow EU(p \cdot a + (1 - p) a'') > EU(p \cdot a' + (1 - p) a'')$$

- **Decomposition :**

**Les utilités espérées se calculent à partir des utilités espérées**

$$EU(p \cdot a + (1 - p) a'') = p \cdot EU(a) + (1 - p) EU(a'')$$

# Calcul d'une politique par programmation dynamique dans MDP à Horizon fini

- On peut faire comme pour les arbres de décision :-)

Pour  $t = N-1$  à 0 faire

  Pour tout  $s$  de  $S_t$  faire

    Pour tout  $a$  de  $A(s)$  faire

$$U(s, a) = r(s, a) + \sum_{s' \in S} P(s' | s, a) \cdot U(s')$$

    Fin pour

$$\pi(s) = \operatorname{argmax}_a U(s, a)$$

$$U(s) = U(s, \pi(s))$$

  Fin pour

Fin pour

# Horizon Infini : Iterative *policy evaluation*

Pour une *politique fixée*  $\pi$  et un coéf d'oubli  $\gamma$  donné

Initialisation : pour tout  $s$ ,  $U(s) = 0$

Repete:

$\Delta = 0$ ;

Pour tout  $s$  de  $S$ :

$u = U(s)$ ;

$U(s) = \sum_{s' \in S} (P(s' | s, \pi(s)) [r(s, \pi(s), s') + \gamma \cdot U(s')])$ ;

$\Delta = \max(\Delta, u - U(s))$ ;

Jusqu'à  $\Delta < \epsilon$

Résultat :  $U \approx U_\pi$

# Iterative policy evaluation (exemple)

→	→	→	○ Fin (10)
↑	Obstacle (-10)	↑	↑
↑	←	↑	←

Cout du déplacement : -1 si nord, 0 pour les autres, discount 0.75

Politique:

→	→	→	○ Fin (10)
↑	Obstacle (-10)	↑	↑
↑	←	↑	←

Cas stochastique : Est : Est (0.5), reste sur place (0.5)

# Policy iteration

→	→	→	○ Fin (1)
↑	obstacle	↑	↑
↑	←	↑	←

■ Et si on allait à l'Est en (1,2)?

■ Idée:

- on part d'une politique quelconque  $\pi_0$
- on calcule (policy evaluation)  $U_{\pi_0}$
- grace à  $U_{\pi_0}$  on trouve une politique améliorée  $\pi_1$
- et on continue jusqu'à stabilité ...

$$\pi_0 \xrightarrow{\text{red}} U_{\pi_0} \xrightarrow{\text{green}} \pi_1 \xrightarrow{\text{red}} U_{\pi_1} \xrightarrow{\text{green}} \pi_2 \xrightarrow{\text{red}} U_{\pi_2} \xrightarrow{\text{green}} \dots \xrightarrow{\text{green}} \pi^* \xrightarrow{\text{red}} U_{\pi^*}$$

**Et ça converge !**

# Policy iteration

Pour coéf d'oubli donné ....

Initialisation : pour tout  $s$ ,  $U(s) = 0$  et  $\pi$  quelconque

Repeter:

Policy Evaluation (c.f. algo prédédent)

(Policy improvement)

$\text{stable} \leftarrow \text{true}$

Pour tout  $s$  de  $S$ :

$b \leftarrow U(s);$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s' \in S} P(s' | s, T(a, s)) [r(s, T(a, s)) + \gamma \cdot U(s')];$

$\text{stable} \leftarrow (b == \pi(s)) \ \&\& \ \text{stable} ;$

Tant que  $\text{stable} = \text{false}$



# Iteration de la valeur (Value Iteration)

- **Chacune des itérations de Policy Iteration nécessite l'exécution de toute l'évaluation ...**
- **Principe de l'algorithme Value Iteration = effectuer seulement la première itération de l'algorithme Policy Evaluation**
- **... et ça converge aussi !**

# Value Iteration

Pour coéf d'oubli donné ....

Initialisation : pour tout  $s$ ,  $U(s) = 0$  et  $\pi$  quelconque

$U \leftarrow \text{Policy Evaluation } (\pi)$

Repeter:

$\Delta = 0;$

Pour tout  $s$  de  $S$ :

$u \leftarrow U(s);$

$U(s) = \max_a \sum_{s' \in S} P(s' | s, T(a, s)) [r(s, T(a, s)) + \gamma \cdot U(s')];$

$\Delta = \max(\Delta, |u - U(s)|);$

Jusqu'à  $\Delta < \epsilon$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s' \in S} P(s' | s, T(a, s)) [r(s, T(a, s)) + \gamma \cdot U(s')];$

# Pour aller plus loin

- **POMDP** : l'état ( $s'$ ) resultat de l'action ( $a$ ) dans ( $s$ ) est partiellement observable : de l'observation  $o$  on deduit une distribution de probabilité sur l'etat  $s'$  à partir de
  - $p(o \mid s', a)$
  - $T(s' \mid a, s)$
  - $p(s \mid \dots)$ , mon etat de croyance sur  $s$
- **Diagrammes d'influence** :  
Réseau Bayesien + variables de décision + utilités
- **D'autres critères que l'utilité esperée**
  - Utilité possibiliste (optimiste, pessimiste) : MDP possibilistes
  - MDP algébriques
  - Choquet (RDU, Probabilities imprecises) ? Aie ! → Complexité accrue !

# Références

- *Decision Theory - An introduction to Dynamic Programming and Sequential Decisions*, J. Bather, Wiley, 2000.
- *Reinforcement Learning - an introduction*, R. Sutton & A. Barto, MIT Press, 1998.
- *Learning to solve Markovian Decision Processes* (<http://www.cs.colorado.edu/baveja/>), S.P. Singh, PhD Dissertation, Department on Computer Science, MIT, 1994.
- *Algorithms for Sequential Decision Making* (<ftp://ftp.cs.brown.edu/pub/techreports/96/cs96-09.pdf>), M. Littman, Department of Computer Science, Brown University, 1996.
- Algebraic Markov Decision Processes. Perny, Patrice and Spanjaard, Olivier and Weng, Paul (2005). In 19th International Joint Conference on Artificial Intelligence, pp. 1372--1377.