

Markov Decision Processes

Wha Sook Jeon

Mobile Computing & Communications Lab.

Contents

- Markov Process with Rewards
 - Discounted rewards
 - Expected discounted sum of rewards
- Markov Decision Process
 - Optimal policy
 - Value iteration
- Partially Observable MDP
 - Observation Process
 - Example

Markov Process

- Stochastic Process (Ransom Process)
 - Collection of random variables, X_t , which represents the system state at time t : $\{X_t, t = 0, 1, \dots\}$
 - Describes the evolution through time of physical process
- Markovian Property
$$\Pr\{X_{t+1} = s_k \mid X_0 = s_0, X_1 = s_1, \dots, X_t = s_i\}$$
$$= \Pr\{X_{t+1} = s_k \mid X_t = s_i\} = p_{ik}$$
- Markov process
 - Stochastic process having Markovian property
 - defined as a tuple (S, P)
 - S : Set of feasible states
 - P : State transition matrix $[p_{ik}]$
 - is used to evaluate the system performance

Markov processes with rewards

- A MP with rewards is a tuple (S, P, R)

- State space: $S = \{s_1, s_2, \dots, s_N\}$
- Transition probability matrix: P

$$P_{ij} = \Pr\{X_{t+1} = s_j \mid X_t = s_i\}$$

- Reward: $R = (r_1, r_2, \dots, r_N)$
 - Each state s_i has a reward r_i

Discounted Rewards

- A reward in the future is not worth as much as a reward now.
- Discounting factor: β
- Discounted sum of future rewards

$$\sum_{t=0}^{\infty} \beta^t R_t$$

- R_t : reward in time t
- R_0 : immediate (now) reward

Expected Reward Sum: $J^*(s_i)$

- $J^*(s_i)$: the expected discounted sum of future rewards, starting in state s_i

$$J^*(s_i) = r_i + \beta \sum_{j=1}^N P_{ij} J^*(s_j)$$

- Matrix Inversion for solving $J^*(s_i)$
 - Using the vector (matrix) notation

$$J^* = \begin{bmatrix} J^*(s_1) \\ J^*(s_2) \\ \vdots \\ J^*(s_N) \end{bmatrix} \quad R = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_N \end{bmatrix} \quad P = \begin{bmatrix} p_{11} & p_{12} & \cdots & p_{1N} \\ p_{21} & p_{22} & \cdots & p_{2N} \\ \vdots & \vdots & \vdots & \vdots \\ p_{N1} & p_{N2} & \cdots & p_{NN} \end{bmatrix}$$

$$J^* = R + \beta P J^*$$

- Then, solve J^* using the matrix inversion

Value Iteration for solving $J^*(s_i)$ (1)

- $J^k(s_i)$: the expected discounted sum of rewards during next k steps, starting at s_i

$$\begin{aligned} J^0(s_i) &\leftarrow r_i && \text{for all } s_i \\ J^1(s_i) &\leftarrow r_i + \beta \sum_{j=1}^N p_{ij} J^0(s_j) && \text{for all } s_i \\ &\vdots \\ J^k(s_i) &\leftarrow r_i + \beta \sum_{j=1}^N p_{ij} J^{k-1}(s_j) && \text{for all } s_i \\ &\vdots \end{aligned}$$

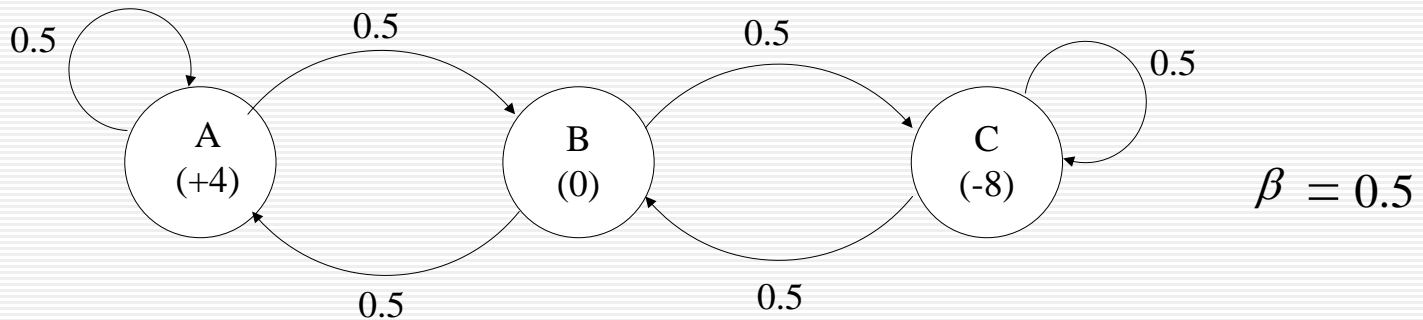
$$\lim_{k \rightarrow \infty} J^k(s_i) = J^*(s_i)$$



$$\text{until } \text{Max}_i |J^{k+1}(s_i) - J^k(s_i)| < \zeta$$

Value Iteration for solving $J^*(s_i)$ (2)

■ Example



k	$J^k(A)$	$J^k(B)$	$J^k(C)$
0	4	0	-8
1	5	-1	-10
2	5	-1.25	-10.75
3	4.94	-1.44	-11

Markov Decision Process

Markov Decision Process (1)

- MDPs provide a mathematical framework for modeling decision-making
 - in situation where outcomes are partly random and partly under the control of the decision maker
- MDPs are useful for studying a wide range of optimization problems via dynamic programming
- A variety of areas including robotics, automated control, economics, etc.

Markov Decision Process (2)

- A discrete time stochastic control process
- Markov chain with rewards and actions
- defined as a tuple (S, A, P, R)
 - State space: $S = \{s_1, s_2, \dots, s_N\}$
 - Action space: A
 - Transition probability matrix: P
$$P_a(i,j) = \Pr\{X_{t+1} = s_j \mid X_t = s_i, a_t = a\}$$
 - Reward: $R = (R(s_1), R(s_2), \dots, R(s_N))$
- A policy is a mapping from states to actions
- What's an **optimal policy**?

Finding the optimal policy: Value Iteration (1)

- Computing the optimal value function using value iteration.
- Optimal policy is the actions for the optimal value function
- **Optimal Value function:** $J^*(s_i)$
 - the expected discounted sum of future rewards, starting at state s_i , when the optimal policy is assumed to be used.
- Computing the optimal value function
 - $J^k(s_i)$: the maximum possible expected discounted sum of rewards we can get, after k time steps starting at s_i
 - $\lim_{k \rightarrow \infty} J^k(s_i) = J^*(s_i)$

Value Iteration (2)

■ Bellman's Equation

$$J^k(s_i) = \max_a \left[r_i + \beta \sum_{j=1}^N p_a(i, j) J^{k-1}(s_j) \right]$$

■ Using the dynamic programming

$k=0$

$J^0(s_i) \leftarrow r_i$ for all s_i

Repeat

$k \leftarrow k+1;$

$J^k(s_i) \leftarrow \max_a [r_i + \gamma \sum_{j=1}^N p_a(i, j) J^{k-1}(s_j)]$ for all s_i

until $(\max_i |J^k(s_i) - J^{k-1}(s_i)| < \zeta)$

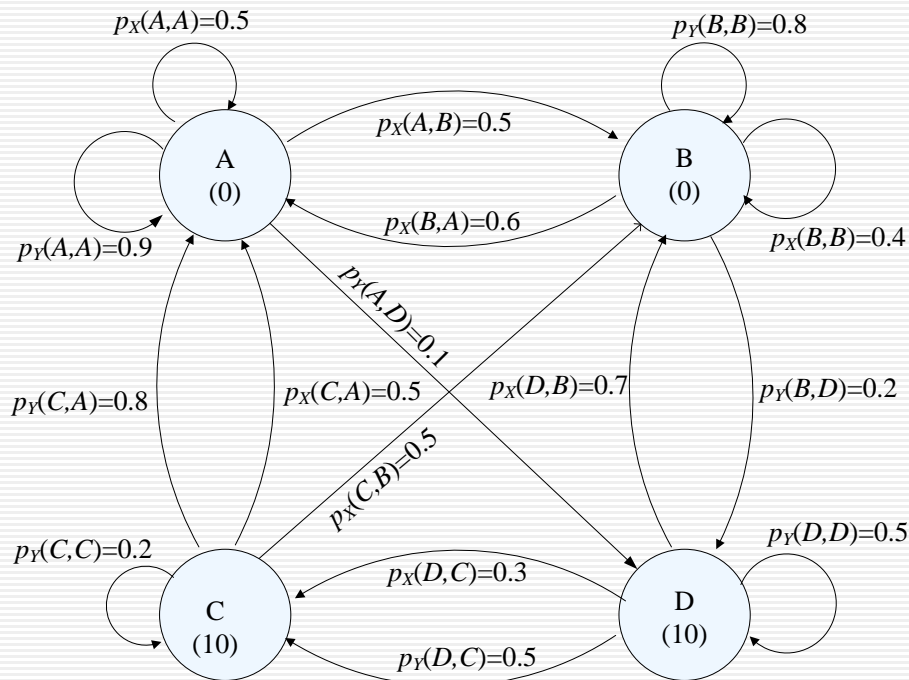
Value Iteration (3)

- Finding the optimal policy
 - Compute $J^*(s_i)$ for all s_i
 - Then, we can obtain the best action in state s_i
 \Rightarrow Optimal policy

$$\arg \max_a \left[r_i + \beta \sum_{j=1}^N p_a(i, j) J^*(s_j) \right]$$

Value Iteration (4)

■ Example



Action set = {X, Y}

$\beta = 0.9$

k	$J^k(A)$	$J^k(B)$	$J^k(C)$	$J^k(D)$
1	0	0	10	10
2	0.9	1.8	11.8	19
3	2.439	4.716	12.772	23.86
4	4.123	7.690	14.055	
5				
6				

$$\begin{aligned}
 J^4(C) &= 10 + 0.9 \times \max\{ \\
 &\quad 0.5 \times J^3(B) + 0.5 \times J^3(A), \\
 &\quad 0.8 \times J^3(A) + 0.2 \times J^3(C) \} \\
 &= 14.055
 \end{aligned}$$

Partially Observable MDP

POMDP (1)

- Defined as a six-tuple (S, A, P, O, Q, R)
- Core process
 - A finite state Markov chain $\{X_t, t \in T\}$, where $T = \{0, 1, \dots\}$.
 - State space: $S = \{1, 2, \dots, N\}$
 - Transition probability matrix: $p_{ij}(a) = \Pr\{X_{t+1} = j \mid X_t = i, a_t = a\}$
 - cannot be directly observable
- Observation process: $\{Y_t, t \in T\}$, where $T = \{0, 1, \dots\}$.
 - By observing Y_t at time t , information regarding the true value of X_t is obtained
- The probabilistic relationship between the core process and observation process when action a is chosen: $q_{ij}(a) = \Pr\{Y_t = j \mid X_t = i, a_{t-1} = a\}$

POMDP (2)

- Random variables

- m_t : the observable value of Y_t

- a_t : the action taken at time t

- d_t : the data available for decision making at time t

$$d_t = (\pi(0), m_1, a_1, m_2, a_2, \dots, a_{t-1}, m_t)$$

- Information vector: $\pi(t) = (\pi_1(t), \pi_2(t), \dots, \pi_N(t))$

- $\pi_i(t) = \Pr\{X_t = i \mid d_t\}$

- Transformation of information vector

- $\pi_i(t+1) = \Pr\{X_{t+1} = i \mid d_{t+1} = (d_t, a_t, m_{t+1} = j)\}$

$$= T_i[\pi(t), a_t, j]$$

$$\begin{aligned} & q_{ij}(a_t) \sum_{k \in S} \pi_k(t) p_{ki}(a_t) \\ &= \frac{\sum_{k \in S} q_{ij}(a_t) \pi_k(t) p_{ki}(a_t)}{\sum_{l \in S} q_{lj}(a_t) \sum_{k \in S} \pi_k(t) p_{kl}(a_t)} \end{aligned}$$

POMDP (3)

- Immediate reward:

- $$r_i(a) = \sum_{j \in S} \sum_{k \in \Theta} R(i, j, k, a) p_{ij}(a) q_{jk}(a)$$

- $R(i, j, k, a)$: immediate reward when action a is taken, the core process is in state i , moves to state j , and observation is k

- Value function

- $$V_{\beta}^n(\pi) = \max_{a \in A} \left\{ \pi \cdot r(a) + \beta \sum_{j \in O} V_{\beta}^{n-1}(T[\pi, j, a]) \eta(j | \pi, a) \right\}$$

- $\eta(j | \pi, a) = \Pr\{Y_{t+1} = j | \pi(t), a_t = a\}$

An Example of POMDP

A POMDP-based Cognitive Radio Sensor Networks

J. A. Han, W. S. Jeon, and D. G. Jeong, "Energy-efficient channel management scheme for cognitive radio sensor networks," *IEEE Transactions on Vehicular Technology*, vol. 60, no. 4, pp. 1905-1910, May 2011

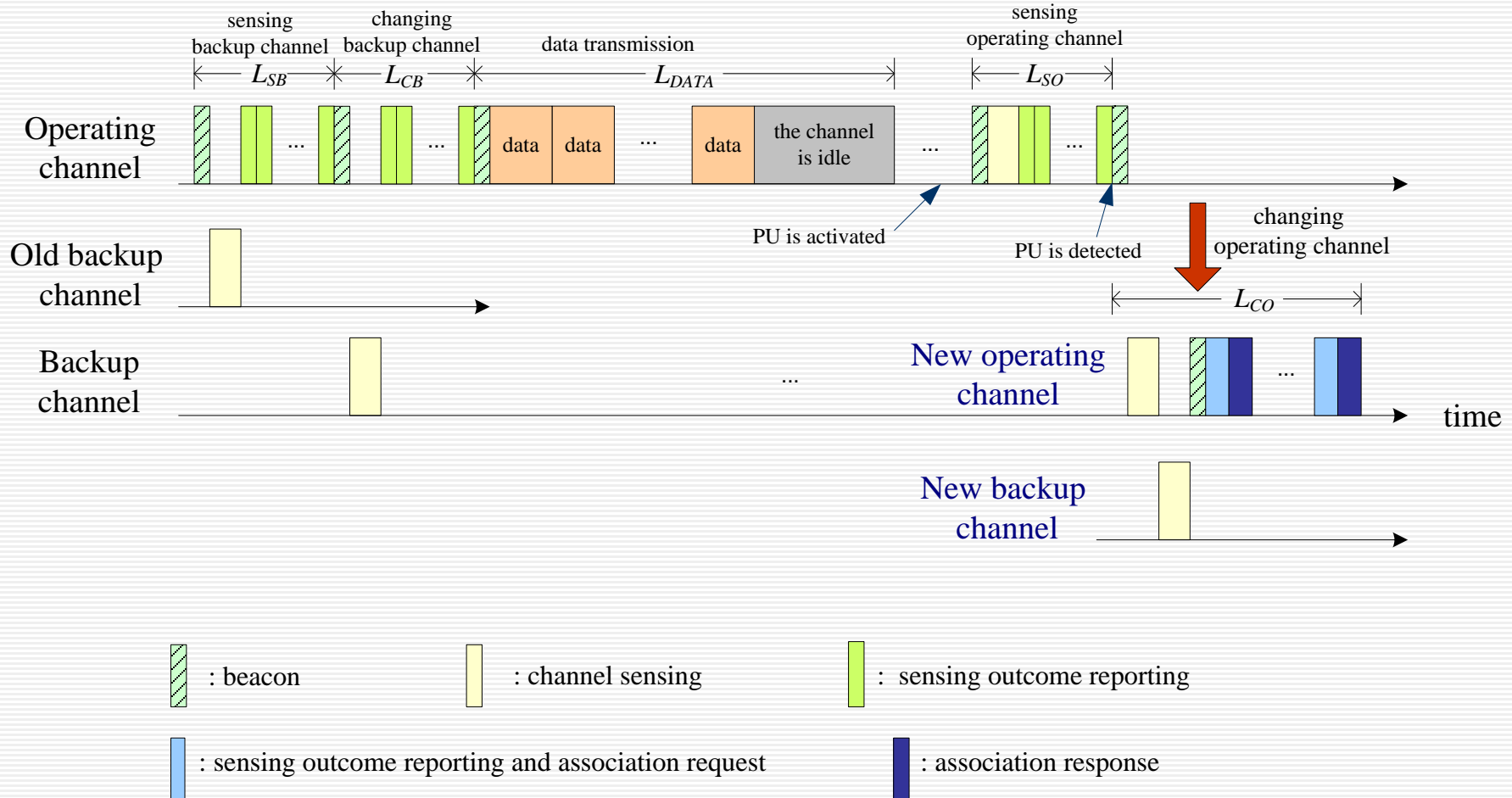
System Description

- Channel pool
 - each channel is licensed to a primary user (PU)
- CR sensor network
 - Cluster with star topology
 - One cluster header (CH) and $(N-1)$ cluster members (CMs)
 - The sensor nodes (CM) opportunistically access to a vacant channel under the control of CH
 - CRSN control of CH: POMDP-based Decision
 - One operating channel and one backup channel

Operation Modes of CRSN

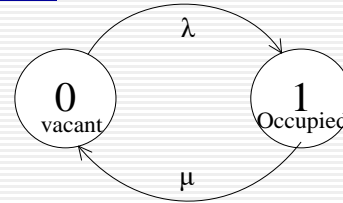
- *DATA* mode
 - The sensor nodes transmit data to CH
 - according to the transmission schedule given by CH
- *SO* mode
 - Sense the operating channel and report the sensing result to CH
- *SB* mode
 - Sense the backup channel and report the sensing result to CH
- *CO* mode
 - Switch to the backup channel (new operating channel)
 - Sense the new operating channel and the randomly selected backup channel
 - Report the sensing result and send new association message
- *CB* mode
 - CH randomly chooses new backup channel
 - All sensor nodes sense the backup channel and report their results

An Example Scenario



Sensing Model

- PU activation model on a channel:
- Operating channel: channel 1
Backup channel: channel 2
- Sensing model: **Energy detection**
 - Each cluster member reports the received energy to CH
 - $s_t^{(m)}$: Sum of the sensing results on the channel m at decision epoch t
 - *Chi-square* distribution
 - Quantize $s_t^{(m)}$ into K levels with thresholds $\gamma_0, \gamma_1, \dots, \gamma_K$
 - Probability that the quantized value (observation value) is k
 - H_0 : Channel m is empty
 - $v_0(k) \equiv \Pr\{\gamma_{k-1} < s_t^{(m)} < \gamma_k \mid H_0\}$
 - H_1 : A PU exists on channel m
 - $v_1(k) \equiv \Pr\{\gamma_{k-1} < s_t^{(m)} < \gamma_k \mid H_1\}$



POMDP Model (1)

- A six-tuple (X, A, O, P, Q, R)
 - X : State space of the core process
 - $X_t = (x_t^{(1)}, x_t^{(2)})$
 - $x_t^{(i)} = 0$ (vacant) or 1 (occupied)
 - A : Action space
 - $\{DATA, SO, SB, CO, CB\}$
 - O : State space of the observation process
 - $O_t = (o_t^{(1)}, o_t^{(2)})$
 - If the channel m is sensed, $o_t^{(m)} \in \{1, \dots, K\}$
 - Otherwise, $o_t^{(m)} = 0$
 - P : Transition probability matrix
$$p_{(i,j)(i',j')}(a) = \Pr\{X_{t+1} = (i', j') \mid X_t = (i, j), A_t = a\}$$

POMDP Model (2)

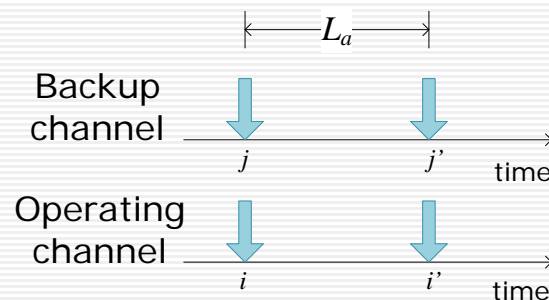
- Transition probability of core process

$u_{i,i'}(a)$: prob. that a channel transits from state i to i'

$w_i(a)$: prob. that a randomly selected channel is in state i

- $DATA, SO, SB$

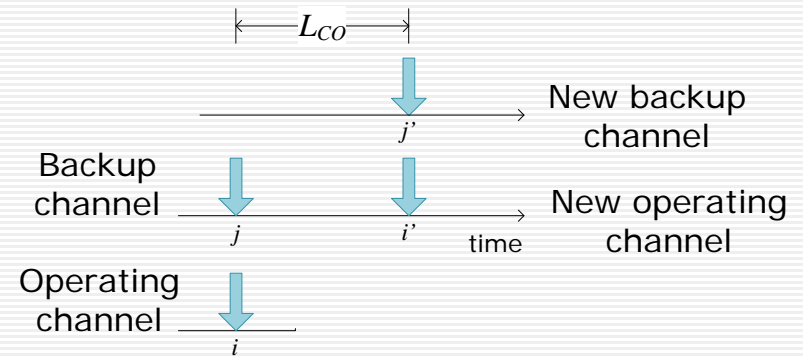
- $$p_{(i,j)(i',j')}(a) = u_{i,i'}(a) \times u_{j,j'}(a)$$



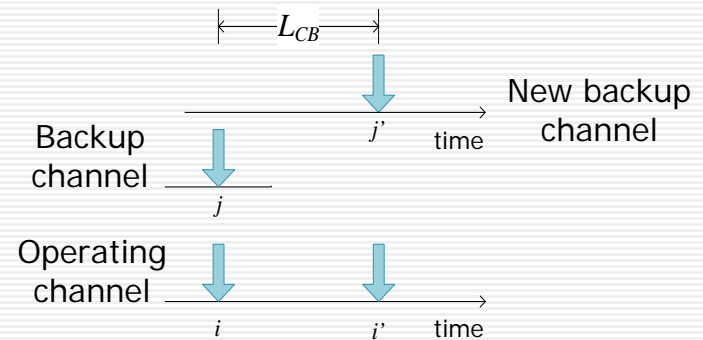
POMDP Model (3)

- P : Transition probability of core process

- CO : $p_{(i,j)(i',j')}(a) = u_{j,i'}(a) \times w_{j'}(a)$



- CB : $p_{(i,j)(i',j')}(a) = u_{i,i'}(a) \times w_{j'}(a)$



POMDP Model (4)

- Q : Probabilistic relation between X and O

- $q_{(i,j)(k,l)}(a) = \Pr\{O_t = (k,l) \mid X_t = (i,j), A_{t-1} = a\}$

Action a	$q_{(i,j)(k,l)}(a)$
DATA	$n(k) \times n(l)$
SO	$v_i(k) \times n(l)$
SB	$n(k) \times v_j(l)$
CO	$v_i(k) \times v_j(l)$
CB	$n(k) \times v_j(l)$

- $v_i(k)$: Prob. that $o_t^{(m)} = k$ when the channel m in state i is sensed
- $n(k)$: Prob. that $o_t^{(m)} = k$ when the channel m is not sensed

POMDP Model (5)

- ***R*: Rewards**

- Control parameter for getting the required performance
 - Penalty on unnecessary energy consumption
 - ex) CO/CB by false alarm
 - Positive reward on protecting PU
- $r_{(i,j)(i',j')}(a)$: reward by taking the action a in state (i,j) which results in the transition (i',j')

- Example

$$r_{(1,\cdot)(\cdot,\cdot)}(DATA) = -10, \quad r_{(0,\cdot)(\cdot,\cdot)}(DATA) = 10$$

$$r_{(0,\cdot)(\cdot,\cdot)}(SO) = r_{(\cdot,0)(\cdot,\cdot)}(SB) = -1$$

- $$R_{(i,j)}(a) = \sum_{i'=0}^1 \sum_{j'=0}^1 r_{(i,j)(i',j')}(a) p_{(i,j)(i',j')}(a)$$

Decision Making (1)

- Information vector: $\Pi(t) = (\pi_{(0,0)}(t), \pi_{(0,1)}(t), \pi_{(1,0)}(t), \pi_{(1,1)}(t))$
 - $\pi_{(i,j)}(t)$: Prob. that the core process is in state (i,j) at decision epoch t .
 - summarizes all information required for the decision-making
 - Update of information vector

- $$\begin{aligned}\pi_{(i,j)}(t+1) &\equiv T_{(i,j)}(\Pi(t), a, (k, l)) \\ &= \Pr\{X_{t+1} = (i, j) \mid \Pi(t), A_t = a, O_{t+1} = (k, l)\} \\ &= \frac{q_{(i,j)(k,l)}(a) \sum_{i'=0}^1 \sum_{j'=0}^1 p_{(i',j')(i,j)}(a) \pi_{(i',j')}(t)}{\sum_{\tilde{i}=0}^1 \sum_{\tilde{j}=0}^1 q_{(\tilde{i},\tilde{j})(k,l)}(a) \sum_{i'=0}^1 \sum_{j'=0}^1 p_{(i',j')(\tilde{i},\tilde{j})}(a) \pi_{(i',j')}(t)}\end{aligned}$$

Decision Making (2)

- Optimal value function

- $$V^*(\Pi) = \max_{a \in \mathbf{A}} \left(\sum_{i=0}^1 \sum_{j=0}^1 \pi_{(i,j)} R_{(i,j)}(a) + \beta \sum_{k=0}^K \sum_{l=0}^K V^*(T(\Pi, a, (k, l))) \times \Pr\{(k, l) | \Pi, a\} \right)$$

- Optimal policy

- $$\delta^*(\Pi) = \arg \max_{a \in \mathbf{A}} \left(\sum_{i=0}^1 \sum_{j=0}^1 \pi_{(i,j)} R_{(i,j)}(a) + \beta \sum_{k=0}^K \sum_{l=0}^K V^*(T(\Pi, a, (k, l))) \times \Pr\{(k, l) | \Pi, a\} \right)$$

Conclusion

- The solution to a MDP is an optimal policy, which gives the action to take for a given state
- When the action is fixed to each state, the resulting MDP behaves like a Markov process
- A POMDP is a generation of a MDP which permits uncertainty