

RN, Chapter  
17.4 — 17.5



# Partially-Observable MDPs

---



# Decision Theoretic Agents

---

- Introduction to Probability [Ch13]
- Belief networks [Ch14]
- Dynamic Belief Networks [Ch15]
- Single Decision [Ch16]
- Sequential Decisions [Ch17]
  - MDPs [Ch17.1 – 17.3]  
(Value Iteration, Policy Iteration,  $TD(\lambda)$ )
  - POMDPs [Ch17.4 – 17.5]  
Dynamic Decision Networks
- Game Theory [Ch17.6 – 17.7]

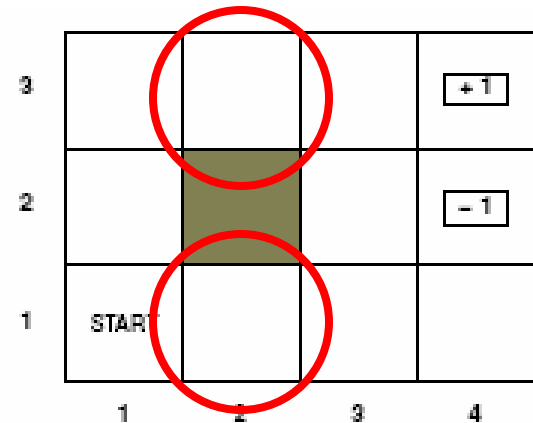
# Partially Accessible Environment

- In inaccessible environment  
percept NOT enough to determine state

**Partially Observable**

**Markov Decision Problem "POMDP"**

⇒ Need to base decision on  
DISTRIBUTION over possible states,  
based all previous percepts, . . . (E)



Eg: Given only distance to walls in 4 directions,

"[ 2, 1 ] ≡ [ 2, 3 ]"

but DIFFERENT actions for each!

If  $P(\text{Loc}[ 2,1 ] \mid E) = 0.8$ ,  $P(\text{Loc}[ 2,3 ] \mid E) = 0.2$

then utility of action  $a$  is

$$0.8 \times U(a \mid \text{Loc}[ 2,1 ]) + 0.2 \times U(a \mid \text{Loc}[ 2,3 ])$$



# Dealing with POMDPs

- Why not view “percept == state”... and just apply MDP alg to “percept”??
  1. Markov property does NOT hold for percepts (percept  $\neq$  states)
    - MDP means  
***next state depends only on current state***
    - But in POMDP:  
next percept does NOT depend only on current percept
  2. May need to take action to *reduce uncertainty* . . . not needed in MDP, as always KNOW state  $\Rightarrow$  utility should include ValueOfInfo. . .

# Extreme Case: Senseless Agent

- What if **NO** observations?
- Perhaps
  - act to reduce uncertainty
  - then go to goal
- (a) Initially: could be ANYWHERE
- (b) After "Left" 5 times
- (c) ... then "Up" 5 times
- (d) ... then "Right" 5 times
- Prob of reaching [4,3]: 77.5%  
but slow: Utility  $\approx 0.08$

0.111	0.111	0.111	<u>0.000</u>
0.111		0.111	<u>0.000</u>
0.111	0.111	0.111	0.111

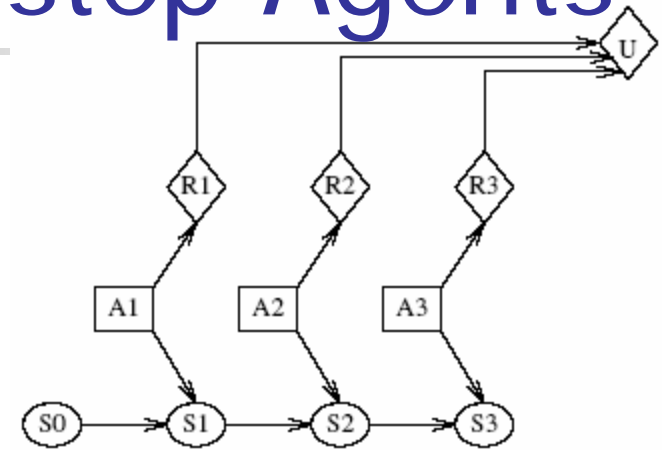
0.300	0.010	0.008	<u>0.000</u>
0.221		0.059	<u>0.012</u>
0.371	0.012	0.008	0.000

0.622	0.221	0.071	<u>0.024</u>
0.005		0.003	<u>0.022</u>
0.003	0.024	0.003	0.000

0.005	0.007	0.019	<u>0.775</u>
0.034		0.007	<u>0.105</u>
0.005	0.006	0.008	0.030



# "Senseless" Multi-step Agents



- Want **sequence of actions**  $[a_1, \dots, a_n]$  that maximizes the expected utility:

$$\operatorname{argmax}_{[a_1, \dots, a_n]} \sum_{[s_0, \dots, s_n]} P(s_0, \dots, s_n \mid a_1, \dots, a_n) \times U([s_0, a_1, \dots, a_n, s_n])$$

- If deterministic, use problem solving techniques to “solve”
  - (finding optimal sequence)
- Stochastic  $\Rightarrow$  don't know state. . .  
but deal w/ DISTRIBUTION OVER STATES

# Unobservable Environments

## ■ View Action-Sequence as BIG action

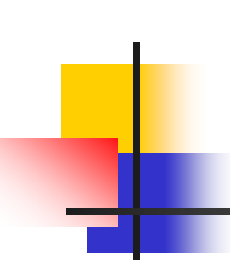
For each possible actions-sequence  $[a_1, \dots, a_n]$   
compute  $P(S_0, S_1, \dots, S_n | a_1, \dots, a_n)$   
compute  $U([s_0, a_1, \dots, a_n, s_n])$   
compute  $score = \sum_{[s_0, \dots, s_n]} P(s_0, \dots, s_n | a_1, \dots, a_n) \cdot U([s_0, a_1, \dots, a_n, s_n])$   
Return action-sequence that gave maximum *score*

## ■ As Markovian:

- $P(S_0, S_1, \dots, S_n | a_1, \dots, a_n) =$   
 $P(S_0) P(S_1 | S_0, a_1) \times P(S_2 | S_1, a_2) \times \dots \times P(S_n | S_{n-1}, a_n)$
- $U([s_0, a_1, \dots, a_n, s_n]) = \sum_t R(s_t)$

⇒ For each action sequence,  
requires searching over all possible sequences of resulting states.

- If  $P(S_{t+1} | S_t, A_{t+1})$  deterministic, can be solved using search...

- 
- Next action must depend on Complete Sequence of Percepts,  $\mathbf{o}$ 
    - (That is all available to agent!)
  - Compress  $\mathbf{o}$  into “distribution over states”
    - $\mathbf{p} = [p_1, \dots, p_n]$  where  $p_i = P(\text{state} = i \mid \mathbf{o})$
  - Given new percept  $\mathbf{o}_t$ ,  
 $\mathbf{p}' = [ P(\text{state} = i \mid \mathbf{o}, \mathbf{o}_t) ]$



# POMDPs

0.111	0.111	0.111	0.000
0.111		0.111	0.000
0.111	0.111	0.111	0.111

$\mathbf{b}_{init}$

- Partially Observable Markov Decision Problem

- $\mathbf{M}_{s,s'}^a \equiv P(s' | s, a)$  : transition
    - $\mathbf{R}(s)$  : reward function
    - $\mathbf{O}(s, \mathbf{o}) \equiv P(\mathbf{o} | s)$  : observation model
- [If senseless:  $O(s, \{\}) = 1.0$ ]

- Belief state  $\mathbf{b}(\cdot) \equiv$  distribution over states

- $\mathbf{b}(s) \equiv P(s | \dots)$  is prob  $\mathbf{b}$  assigns to  $s$
  - Eg:  $\mathbf{b}_{init} = \langle 1/9, 1/9, \dots, 1/9, 0, 0 \rangle$

- Given  $\mathbf{b}(\cdot)$ , after action  $\mathbf{a}$ , observation  $\mathbf{o}$

- $\mathbf{b}'(s') = O(s', \mathbf{o}) \sum_s P(s | a, s') \mathbf{b}(s)$
  - $\mathbf{b}' = \text{Forward}(\mathbf{b}, \mathbf{a}, \mathbf{o})$

## Filtering!

- Optimal action depends only on current belief state!
- ... not on actual state

# What to do, in POMDP?

$$b'(s') = \alpha O(s', o) \sum_s P(s' | a, s) b(s)$$

- Policy  $\pi$  maps BELIEF STATE  $b$  to ACTION  $a$   
 $\pi(b) = a \quad \pi: [0, 1]^n \mapsto \{ \text{North, East, South, West} \}$
- Given optimal policy  $\pi^*$ 
  - 1. Given  $b_i$  compute/execute action  $a_i = \pi(b_i)$
  - 2. Receive observation  $o_i$
  - 3. Compute  $b_{i+1} = \text{Forward}(b_i, a_i, o_i)$
- With MDPs, can just "reach" new state ... no observations...  
With POMDPs, need to know observation  $o_i$  to determine  $b'$
- Some POMDP actions may be
  - to reduce uncertainty
  - to gather information
- How to compute optimal  $\pi^*$  ?  
. . . perhaps make POMDP look like MDP?

# Transform POMDP into MDP ?

- Every MDP needs

- Transition  $M$ : State Action  $\mapsto$  Distribution over State
- Reward  $R$ : State  $\mapsto \mathfrak{R}$

$\Rightarrow$  Given “belief state”  $b$ , need

- $\rho(b) =$  (expected) reward for being in  $b$   
 $= \sum_s b(s) R(s)$

- $\mu(b, a, b') = P( b' \mid b, a )$

... prob of reaching  $b'$  if take action  $a$  in  $b$ . . .

Depends on observation  $o$ :

- $P( b' \mid a, b ) = \sum_o P( b' \mid o, a, b ) P( o \mid a, b )$   
 $= \sum_o \delta[ b' = \text{Forward}(b, a, o) ] P( o \mid a, b )$

- where  $\delta[ b' = \text{Forward}(b, a, o) ] = 1$  iff  $b' = \text{Forward}(b, a, o)$

- Need DISTRIBUTION over observations . . .

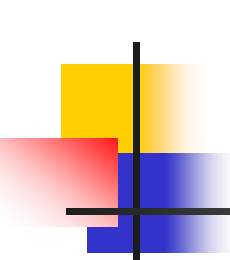


# Distribution over Observations

- $P(o | a, b)$ 
$$= \sum_{s'} P(o | a, s', b) P(s' | a, b)$$
$$= \sum_{s'} O(s', o) P(s' | a, b)$$
$$= \sum_{s'} O(s', o) \sum_s P(s' | a, s) b(s)$$

- So...

$$\begin{aligned} \mu_{b,b'}^a &= P(b' | a, b) \\ &= \sum_o P(b' | o, a, b) P(o | a, b) \\ &= \sum_o \mathbb{1}[b' = \text{Forward}(b, a, o)] \sum_{s'} O(s', o) \sum_s P(s' | a, s) b(s) \end{aligned}$$



# POMDP $\Rightarrow$ ? MDP ??

- $\mu_{b,b'}^a = P(b' | b, a)$   
 $\rho(b)$  = (expected) reward  
... define OBSERVABLE MDP!  
(Agent can always observe its beliefs!)
- Optimal policy for this MDP  $\pi^*(b)$   
is optimal for POMDP

**Solving POMDP on physical state space**

≡

**solving MDP on corresponding BELIEF STATE SPACE!**

- But. . . this MDP has continuous  
(and usually HIGH-Dimension)  
state space!
- Fortunately . . .

# Transform POMDP into MDP

- Fortunately,  $\exists$  versions of

- value iteration
- policy iteration

that apply to such continuous-space MDPs  
(Represent  $\pi(\mathbf{b})$  as set of REGIONS of belief space  
each with specific optimal action)

$U \equiv$  LINEAR FUNCTION of  $\mathbf{b}$  w/in each region

Each iteration refines boundaries of regions . . .

- Solution:

[Left, Up, Up, Right, Up, Up, Right, Up, Up, ...]

(Left ONCE to ensure NOT at [4,1],

then go Right and Up until reaching [4, 3].)

Succeeds 86.6%, quickly. . .

Utility = 0.38

- In general: finding optimal policies is PSPACE-Hard!

3				<span style="border: 1px solid black; padding: 2px;">+1</span>
2				<span style="border: 1px solid black; padding: 2px;">-1</span>
1	START			
	1	2	3	4

# Solving POMDP, in General

**function** DECISION-THEORETIC-AGENT( *percept*) **returns** *action*  
calculate updated probabilities for current state based on  
available evidence including current percept and previous action  
calculate outcome probabilities for actions  
given action descriptions and probabilities of current states  
select *action* with highest expected utility  
given probabilities of outcomes and utility information  
**return** *action*

- To determine current state  $S_t$ :
  - Deterministic: previous action  $a_{t-1}$  from  $S_{t-1}$  determines  $S_t$
  - Accessible: current percepts identify  $S_t$
  - Partially accessible: use BOTH action and percepts
- Computing outcome probabilities:  
... as above
- Computing *expected utilities*:  
At time  $t$ , need to think about making decision  $D_{t+i}$   
At that time  $t+i$ , agent will THEN have percepts  $E_{t+1}, \dots, E_{t+i}$   
But not known now (at time  $t$ ). . .

# Challenges

- To decide about  $A_t$  (action at time t), need distribution of current state based on
  - all evidence ( $E_i$  is evidence at time i)
  - all actions ( $A_i$  is action at time i)

$$Bel(S_t) \equiv P(S_t | E_1, \dots, E_t, A_1, \dots, A_{t-1})$$

⇒ very hard to compute, in general

- But. . . some simplifications:
  - $P(S_t | S_1, \dots, S_{t-1}, A_1, \dots, A_{t-1}) = P(S_t | S_{t-1}, A_{t-1})$   
Markov
  - $P(E_t | S_1, \dots, S_t, E_1, \dots, E_t, A_1, \dots, A_{t-1}) = P(E_t | S_t)$   
Evidence depends only on current world
  - $P(A_{t-1} | A_1, \dots, A_{t-2}, E_1, \dots, E_{t-1}) = P(A_{t-1} | E_1, \dots, E_{t-1})$   
Agent acts based only input. . . and knows what it did

**RECURSIVE form of  $Bel()$**  updated with each evidence:

- **Prediction Phase:**

Predict distribution over state, before evidence

$$\underline{Bel}(S_t) = \sum_{s_{t-1}} P(S_t | S_{t-1} = s_{t-1}, A_{t-1}) Bel(S_{t-1} = s_{t-1})$$

- **Estimation Phase:** ... Incorporate  $E_t$

$$Bel(S_t) = \alpha P(E_t | S_t) \underline{Bel}(S_t)$$



# Decision-Theoretic Agent

function DECISION-THEORETIC-AGENT( $E_t$ ) returns an action

inputs:  $E_t$ , the percept at time  $t$

static:  $BN$ , a belief network with nodes  $X$

$Bel(X)$ , a vector of probabilities, updated over time

$$\widehat{Bel}(X_t) \leftarrow \sum_{X_{t-1}} P(X_t | X_{t-1}=x_{t-1}, A_{t-1}) Bel(X_{t-1}=x_{t-1})$$

$$Bel(X_t) \leftarrow \alpha P(E_t | P X_t) \widehat{Bel}(X_t)$$

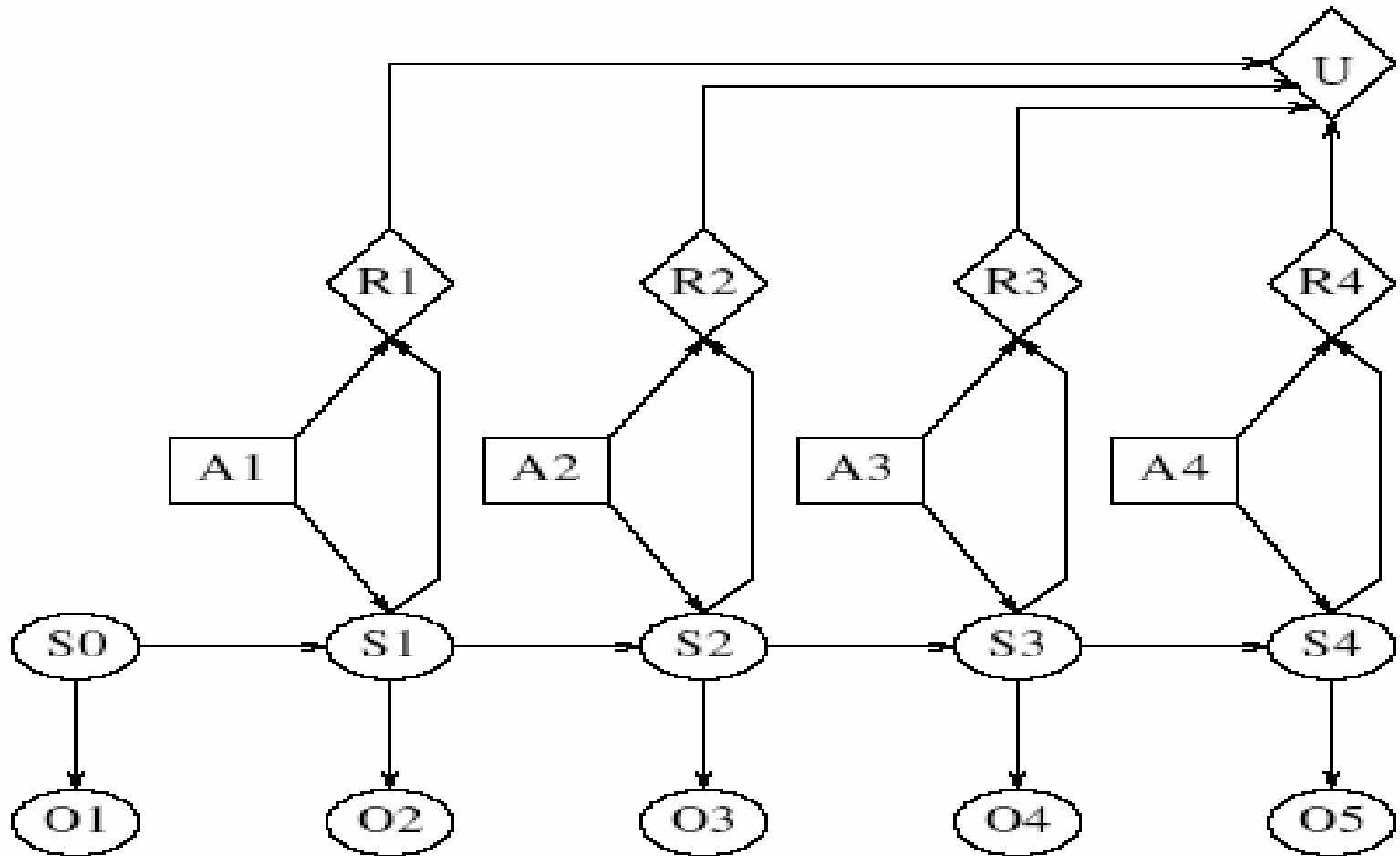
$$action \leftarrow \arg \max_{A_t} \sum_{X_t} \left[ Bel(X_t=x_t) \sum_{X_{t+1}} P(X_{t+1}=x_{t+1} | X_t=x_t, A_t) U(x_{t+1}) \right]$$

return  $action$

- Dependencies are reasonable:
  - action model:  $P(S_t | S_{t-1}, A_{t-1})$
  - sensor model:  $P(E_t | S_t)$

# Partially Observable MDPs

## Dynamic Decision Networks





# Approximate Method for Solving POMDP's

## Two Key Ideas:

- Compute optimal value function  $U(S)$  assuming complete observability  
(Whatever will be needed later, will be available)
- Maintain  $\text{Bel}(S_t) = P(S_t | E_t, A_t, S_{t-1}, \dots, S_0, E_0)$

## At each time $t$

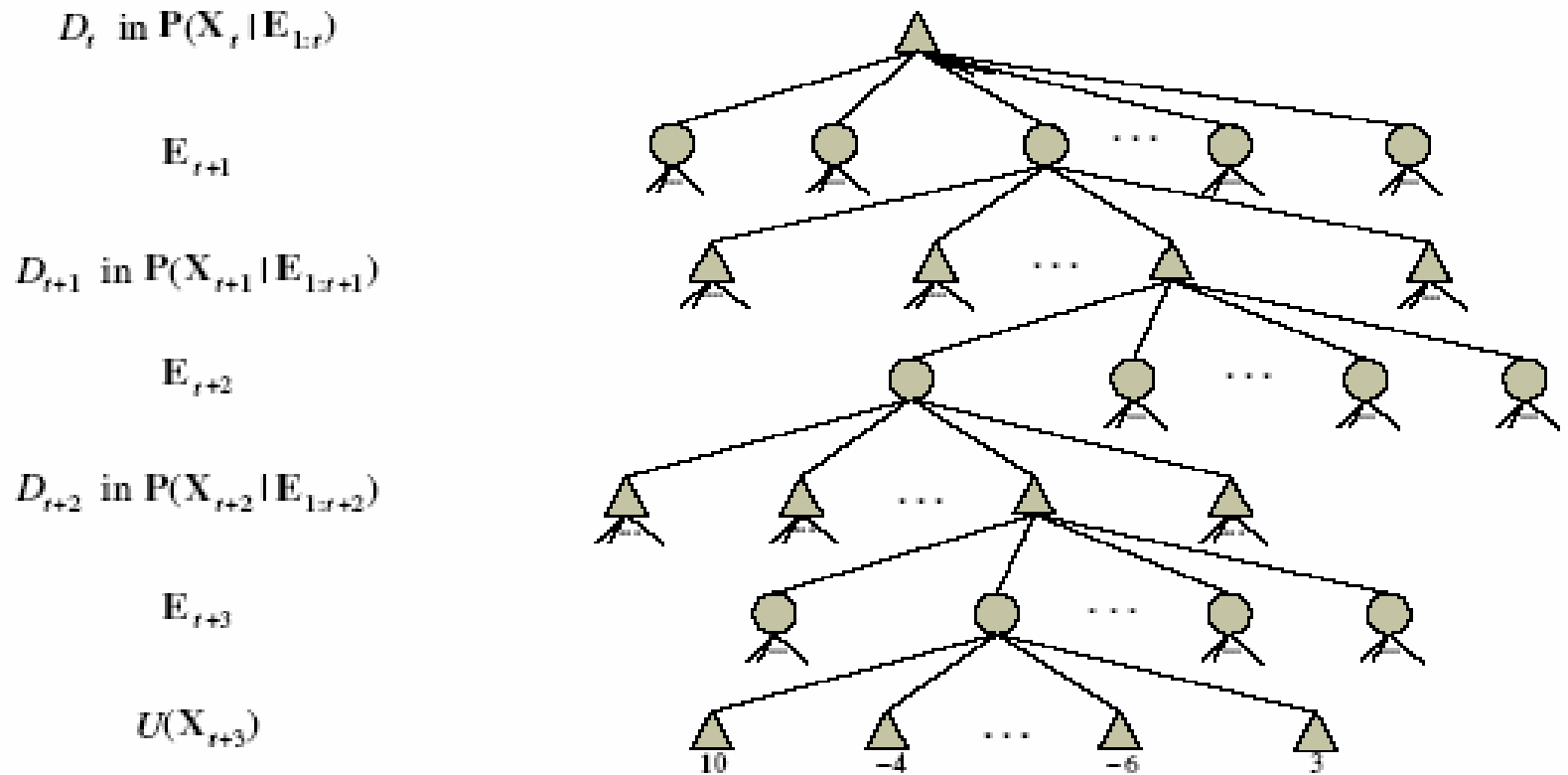
- Observe current percept  $E_t$
- Update  $\text{Bel}(S_t)$
- Choose next  $k$  optimal actions  $[a_{t+1}, \dots, a_{t+k}]$  to maximize

$$\sum_{S_{t+1}, \dots, S_{t+k}} \sum_{E_{t+1}, \dots, E_{t+k}} P(S_{t+1} | S_t, a_{t+1}) P(E_{t+1} | S_{t+1}) \dots P(S_{t+k} | S_{t+k-1}, a_{t+k})$$
$$[\sum_{i=1}^k R(S_{t+i} | S_{t+i-1}, a_{t+i}) + U(S_{t+k})]$$

- Perform action  $a_{t+1}$

# Look-ahead Search

## ExpectiMiniMax





# Wrt Dynamic Decision Networks

---

- Handle uncertainty correctly... sometimes efficiently...
- Deal with streams of sensor input
- Handle unexpected events (as have no fixed “plan”)
- Handle noisy sensors, sensor failure
- Act in order to obtain information as well as to receive rewards
- Handle relatively large state spaces as they decompose state into set of state var's with sparse connections
- Exhibit graceful degradation under time pressure and in complex environments using various approximation techniques



# Open Problems

## wrt Probabilistic Agents

- First-order probabilistic representations
  - If any car hits lamp post going over 30mph, occupants of car injured with probability 0.60.*
- Methods for scaling up MDP's
- More efficient algorithms for POMDP's
- Learning environment
  - $M_{ij}^a, P(E | S), \dots$



# Probabilistic Agents Summary

- **Three key components:**

- $P(S' | S, A)$  (action model)
- $P(E | S)$  (sensor model)
- $R(S' | S, A)$  (reward function)

- In accessible environments,

{ Value iteration, Policy iteration } work well.

Each computes local (state) utility function, optimal policy.

- In { unobservable, partially-observable } environments,

lookahead search gives approx solutions

- Updating current beliefs in a DDN is easy.
- Look-ahead search is hard.