

Decision Making Under Uncertainty

CMSC 671, Fall 2016
Class #23

material from Marie desJardin, Lise Getoor,
Jean-Claude Latombe, and Daphne Koller

1

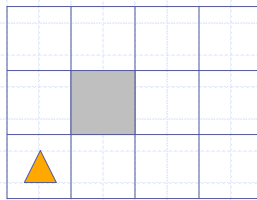
SEQUENTIAL DECISION MAKING UNDER UNCERTAINTY

Sequential Decision Making

- ◆ Finite Horizon
- ◆ Infinite Horizon

3

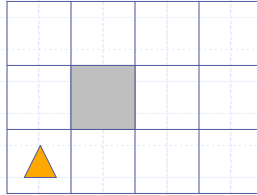
Simple Robot Navigation Problem



- In each state, the possible actions are **U**, **D**, **R**, and **L**

4

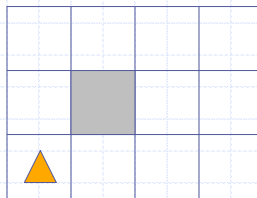
Probabilistic Transition Model



- In each state, the possible actions are **U**, **D**, **R**, and **L**
- The effect of **U** is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)

5

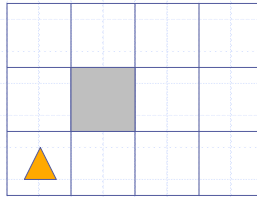
Probabilistic Transition Model



- In each state, the possible actions are **U**, **D**, **R**, and **L**
- The effect of **U** is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)
 - With probability 0.1, the robot moves right one square (if the robot is already in the rightmost row, then it does not move)

6

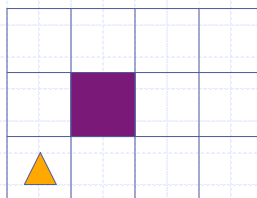
Probabilistic Transition Model



- In each state, the possible actions are **U**, **D**, **R**, and **L**
- The effect of **U** is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)
 - With probability 0.1, the robot moves right one square (if the robot is already in the rightmost row, then it does not move)
 - With probability 0.1, the robot moves left one square (if the robot is already in the leftmost row, then it does not move)

7

Probabilistic Transition Model



- In each state, the possible actions are **U**, **D**, **R**, and **L**
- The effect of **U** is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)
 - With probability 0.1, the robot moves right one square (if the robot is already in the rightmost row, then it does not move)
 - With probability 0.1, the robot moves left one square (if the robot is already in the leftmost row, then it does not move)
- **D**, **R**, and **L** have similar probabilistic effects

8

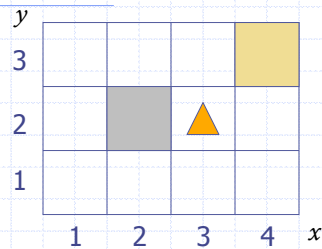
Markov Property

The transition properties depend only on the current state, not on the previous history (how that state was reached)

Markov assumption generally: current state only ever depends on previous state (or finite set of previous states).

9

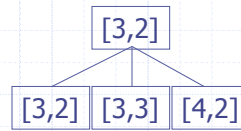
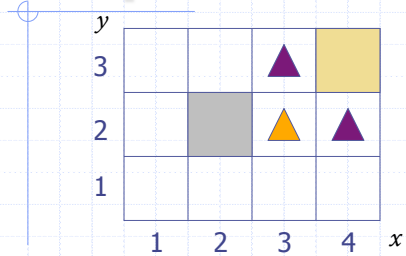
Sequence of Actions



- Planned sequence of actions: (U, R)

10

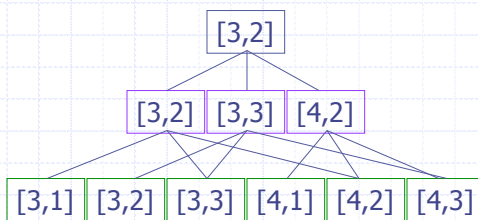
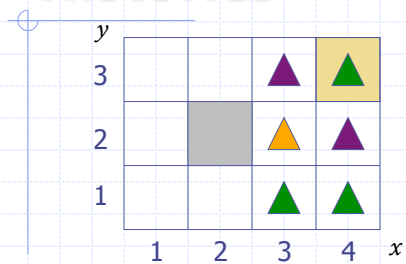
Sequence of Actions



- Planned sequence of actions: (U, R)
- U is executed

11

Histories

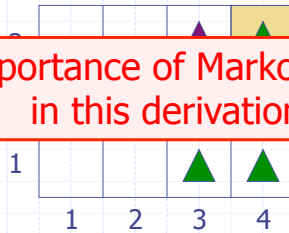


- Planned sequence of actions: (U, R)
- U has been executed
- R is executed
- 9 possible sequences of states – called histories
- 6 possible final states for the robot!

12

Probability of Reaching the Goal

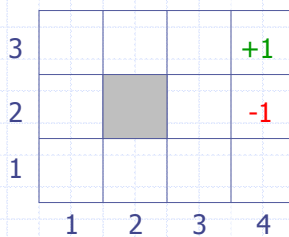
Note importance of Markov property in this derivation



- $P([4,3] | (U,R).[3,2]) =$
 $P([4,3] | R.[3,3]) \times P([3,3] | U.[3,2])$
 $+ P([4,3] | R.[4,2]) \times P([4,2] | U.[3,2])$
- $P([4,3] | R.[3,3]) = 0.8$ • $P([3,3] | U.[3,2]) = 0.8$
- $P([4,3] | R.[4,2]) = 0.1$ • $P([4,2] | U.[3,2]) = 0.1$
- $P([4,3] | (U,R).[3,2]) = 0.65$

13

Utility Function



- [4,3] provides power supply
- [4,2] is a sand area from which the robot cannot escape

14

Utility Function

3			+1	
2			-1	
1				
	1	2	3	4

- [4,3] provides power supply
- [4,2] is a sand area from which the robot cannot escape
- The robot needs to recharge its batteries

15

Utility Function

3			+1	
2			-1	
1				
	1	2	3	4

- [4,3] provides power supply
- [4,2] is a sand area from which the robot cannot escape
- The robot needs to recharge its batteries
- [4,3] and [4,2] are terminal states

16

Utility of a History

3				+1
2				-1
1				
	1	2	3	4

- [4,3] provides power supply
- [4,2] is a sand area from which the robot cannot escape
- The robot needs to recharge its batteries
- [4,3] or [4,2] are terminal states
- The utility of a history is defined by the utility of the last state (+1 or -1) minus $n/25$, where n is the number of moves
 - Many utility functions possible, for many kinds of problems.

17

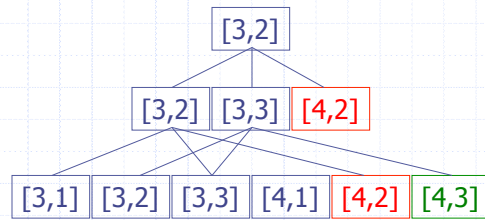
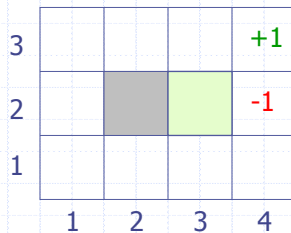
Utility of an Action Sequence

3				+1
2				-1
1				
	1	2	3	4

- Consider the action sequence (U,R) from [3,2]

18

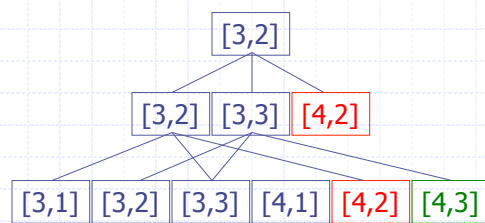
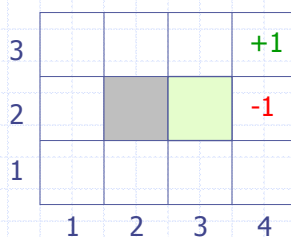
Utility of an Action Sequence



- Consider the action sequence (U,R) from [3,2]
- A run produces one of 7 possible histories, each with some probability

19

Utility of an Action Sequence

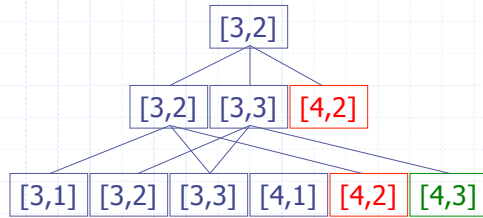
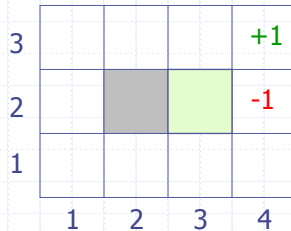


- Consider the action sequence (U,R) from [3,2]
- A run produces one of 7 possible histories, each with some probability
- The **utility of the sequence** is the expected utility of the histories:

$$U = \sum_h U_h P(h)$$

20

Optimal Action Sequence



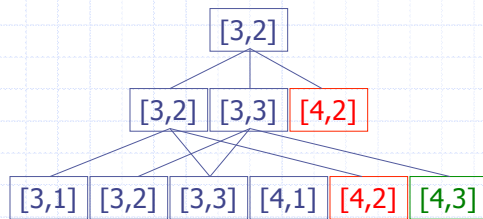
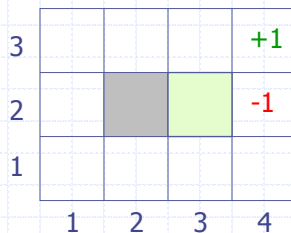
- Consider the action sequence (U,R) from [3,2]
- A run produces one of 7 possible histories, each with some probability
- The utility of the sequence is the expected utility of the histories:

$$U = \sum_h U_h P(h)$$

- The optimal sequence is the one with maximal utility

21

Optimal Action Sequence



- Consider the action sequence (U,R) from [3,2]
- A run produces only if the sequence is executed blindly! ability
- The utility of the sequence is the expected utility of the histories
- The optimal sequence is the one with maximal utility
- **But is the optimal action sequence what we want to compute?**

22

Reactive Agent Algorithm

Repeat:

- ◆ $s \leftarrow$ sensed state
- ◆ If s is terminal then exit
- ◆ $a \leftarrow$ choose action (given s)
- ◆ Perform a

Accessible or
observable state

23

Policy (Reactive/Closed-Loop Strategy)

3	→	→	→	+1
2	↑		↑	-1
1	↑	←	←	←
	1	2	3	4

- A **policy Π** is a complete mapping from states to actions

24

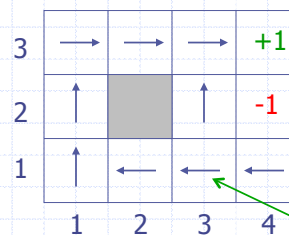
Reactive Agent Algorithm

Repeat:

- ◆ $s \leftarrow$ sensed state
- ◆ If s is terminal then exit
- ◆ $a \leftarrow \Pi(s)$
- ◆ Perform a

25

Optimal Policy



- A **policy** Π is a complete
- The **optimal policy** Π^* is the policy that maximizes the **expected utility** over all possible histories (ending at a terminal state)

Note that [3,2] is a “dangerous” state that the optimal policy tries to avoid

Makes sense because of Markov property

26

Optimal Policy

3	→	→	→	+1
2	↑		↑	-1
1	↑	←	←	←
	1	2	3	4

- A **policy** Π is a complete strategy that specifies an action for every state.
- The **optimal policy** Π^* is the policy that maximizes the expected utility over all possible histories with maximal expected utility.

This problem is called a Markov Decision Problem (MDP)

How to compute Π^* ?

27

Additive Utility

- ◆ History $H = (s_0, s_1, \dots, s_n)$
- ◆ The utility of H is **additive** iff:

$$U(s_0, s_1, \dots, s_n) = R(0) + U(s_1, \dots, s_n) = \sum R(i)$$

Reward

28

Additive Utility

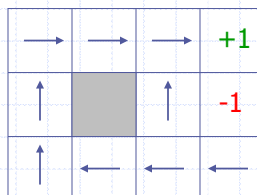
- ◆ History $H = (s_0, s_1, \dots, s_n)$
- ◆ The utility of H is **additive** iff:

$$U_{(s_0, s_1, \dots, s_n)} = R(0) + U_{(s_1, \dots, s_n)} = \sum R(i)$$
- ◆ Robot navigation example:
 - $R(n) = +1$ if $s_n = [4, 3]$
 - $R(n) = -1$ if $s_n = [4, 2]$
 - $R(i) = -1/25$ if $i = 0, \dots, n-1$

29

Principle of Max Expected Utility

- ◆ History $H = (s_0, s_1, \dots, s_n)$
- ◆ Utility of H : $U_{(s_0, s_1, \dots, s_n)} = \sum R(i)$



*reminder! utility
of a sequence:*
 $U = \sum_h U_h P(h)$

First-step analysis →

- ◆ $U(i) = R(i) + \max_a \sum_k P(k | a, i) U(k)$
- ◆ $\Pi^*(i) = \arg \max_a \sum_k P(k | a, i) U(k)$

30

Defining State Utility

◆ Problem:

- When making a decision, we only know the reward so far, and the possible actions
- We've defined utility retroactively (i.e., the utility of a history is known *once we finish it*)
- What is the utility of a particular **state** in the middle of decision making?
- Need to compute **expected utility** of possible future histories

31

Value Iteration

- ◆ Initialize the utility of each non-terminal state s_i to $U_0(i) = 0$ } or some uniform or uniformly distributed value

- ◆ For $t = 0, 1, 2, \dots$, do:

$$U_{t+1}(i) \leftarrow R(i) + \max_a \sum_k P(k | a, i) U_t(k)$$

3				+1
2				-1
1				
	1	2	3	4

32

Value Iteration

- ◆ Initialize the utility of each state s_i to $U_0(i) = 0$

- ◆ For $t = 0, 1, 2, \dots$, do:

$$U_{t+1}(i) \leftarrow R(i) + \max_a \sum_k P(k | a, i) U_t(k)$$

Note the importance of terminal states and connectivity of the state-transition graph

3	→ 0.812 → 0.868 → ??? → +1
2	↑ 0.762 ↑ 0.660 ↑ -1
1	↑ 0.705 ← 0.655 ← 0.611 ← 0.388
	1 2 3 4

EXERCISE: What is $U^*([3,3])$ (assuming that the other U^* are as shown)?

Value Iteration

- ◆ Initialize the utility of each non-terminal state s_i to $U_0(i) = 0$

- ◆ For $t = 0, 1, 2, \dots$, do:

$$U_{t+1}(i) \leftarrow R(i) + \max_a \sum_k P(k | a, i) U_t(k)$$

3	→ 0.812 → 0.868 → 0.918 → +1
2	↑ 0.762 ↑ 0.660 ↑ -1
1	↑ 0.705 ← 0.655 ← 0.611 ← 0.388
	1 2 3 4

$$U_{3,3}^* = R_{3,3} + [P_{3,2} U_{3,2}^* + P_{3,3} U_{3,3}^* + P_{4,3} U_{4,3}^*]$$

Policy Iteration

- ◆ Pick a policy Π at random

37

Policy Iteration

- ◆ Pick a policy Π at random
- ◆ Repeat:
 - Compute the utility of each state for Π
$$\mathbf{U}_{t+1}(i) \leftarrow \mathbf{R}(i) + \sum_k \mathbf{P}(k | \Pi(i), i) \mathbf{U}_t(k)$$

38

Policy Iteration

- ◆ Pick a policy Π at random
- ◆ Repeat:
 - Compute the utility of each state for Π
$$\mathbf{U}_{t+1}(i) \leftarrow \mathbf{R}(i) + \sum_k \mathbf{P}(k | \Pi(i).i) \mathbf{U}_t(k)$$
 - Compute the policy Π' given these utilities
$$\Pi'(i) = \arg \max_a \sum_k \mathbf{P}(k | a.i) \mathbf{U}(k)$$

39

Policy Iteration

- ◆ Pick a policy Π at random
- ◆ Repeat:
 - Compute the utility of each state for Π
$$\mathbf{U}_{t+1}(i) \leftarrow \mathbf{R}(i) + \sum_k \mathbf{P}(k | \Pi(i).i) \mathbf{U}_t(k)$$
 - Compute the policy Π' given these utilities
$$\Pi'(i) = \arg \max_a \sum_k \mathbf{P}(k | a.i) \mathbf{U}(k)$$

Or solve the set of linear equations:
$$\mathbf{U}(i) = \mathbf{R}(i) + \sum_k \mathbf{P}(k | \Pi(i).i) \mathbf{U}(k)$$

(often a sparse system)
 - If $\Pi' = \Pi$ then return Π

40

Infinite Horizon

In many problems, e.g., the robot navigation example, histories are potentially unbounded and the same state can be reached many times

3				+1
2				-1
1				
	1	2	3	4

What if the robot lives forever?

One trick:
Use discounting to make an infinite horizon problem mathematically tractable

41

Value Iteration

◆ Value iteration:

- Initialize state values (expected utilities) randomly
- Repeatedly update state values using best action, according to current approximation of state values
- Terminate when state values stabilize
- Resulting policy will be the best policy because it's based on accurate state value estimation

◆ Policy iteration:

- Initialize policy randomly
- Repeatedly update state values using best action, according to current approximation of state values
- Then update policy based on new state values
- Terminate when policy stabilizes
- Resulting policy is the best policy, but state values may not be accurate (may not have converged yet)
- Policy iteration is often faster (because we don't have to get the state values right)

◆ **Both methods have a major weakness: They require us to know the transition function exactly in advance!**

42