

DECISION MAKING UNDER UNCERTAINTY

CMSC 671, Fall 2017

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

SEQUENTIAL DECISION MAKING UNDER UNCERTAINTY

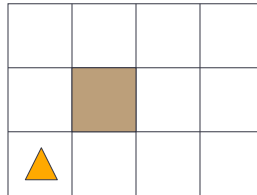
The Big Idea

- “Planning”: Find a sequence of steps to accomplish a goal.
 - Given start state, transition model, goal functions...
- This is a kind of **sequential decision making**.
 - Transitions are deterministic.
- What if they are stochastic (probabilistic)?
 - One time in ten, you drop your sock
- Probabilistic Planning: Make a plan that accounts for probability by **carrying it through the plan**.

Sequential Decision Making

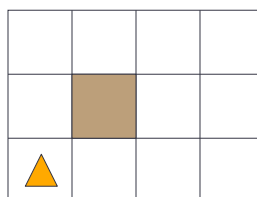
- Finite Horizon
- Infinite Horizon

Simple Robot Navigation Problem



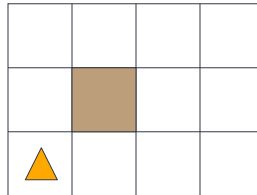
- In each state, the possible actions are **U**, **D**, **R**, and **L**

Probabilistic Transition Model



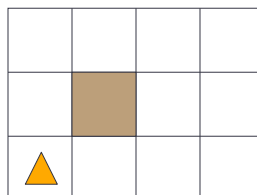
- In each state, the possible actions are **U**, **D**, **R**, and **L**
- The effect of **U** is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)

Probabilistic Transition Model



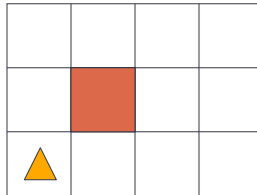
- In each state, the possible actions are **U**, **D**, **R**, and **L**
- The effect of **U** is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)
 - With probability 0.1, the robot moves right one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

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)

Probabilistic Transition Model



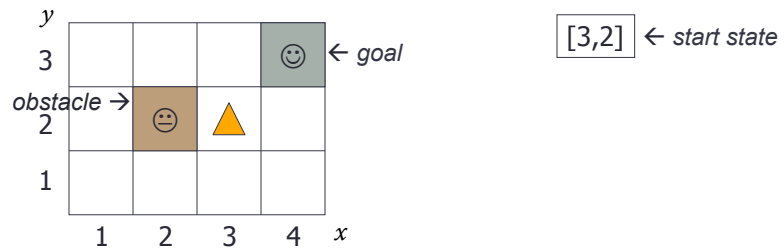
- In each state, the possible actions are U, D, R, and L
- The effect of U is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)
 - With probability 0.1, the robot moves right one square (if the robot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

Markov Property

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

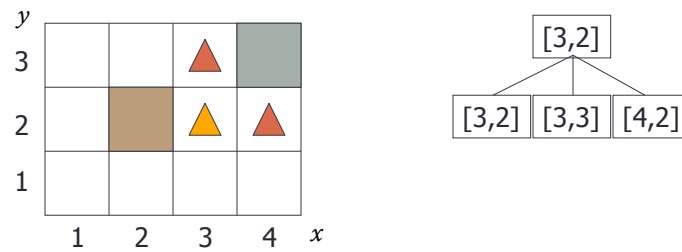
Markov assumption generally: current state only ever depends on previous state (or finite set of previous states).

Sequence of Actions



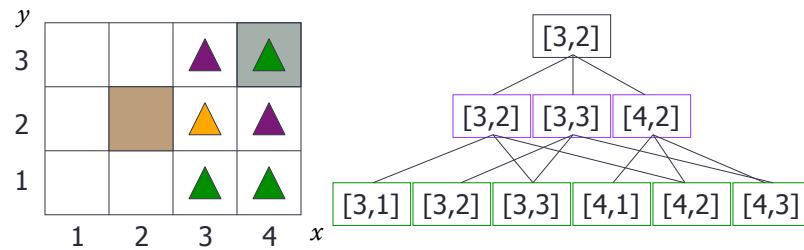
- Planned sequence of actions: (U, R)

Sequence of Actions



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

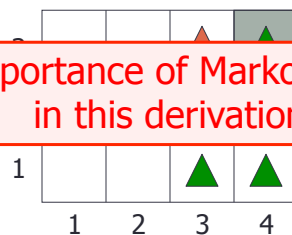
Histories



- 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!

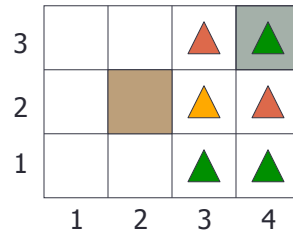
Probability of Reaching the Goal

Note importance of Markov property in this derivation



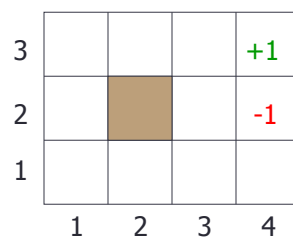
- $P([4,3] \mid (U,R).[3,2]) =$
 $P([4,3] \mid R.[3,3]) \times P([3,3] \mid U.[3,2])$
 $+ P([4,3] \mid R.[4,2]) \times P([4,2] \mid U.[3,2])$
- $P([4,3] \mid R.[3,3]) = 0.8$ • $P([3,3] \mid U.[3,2]) = 0.8$
- $P([4,3] \mid R.[4,2]) = 0.1$ • $P([4,2] \mid U.[3,2]) = 0.1$
- $P([4,3] \mid (U,R).[3,2]) =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

Probability of Reaching the Goal



- Core idea: **multiply backward** probabilities of each step taken from end state reached
- But we still need to consider different ways of reaching a state
 - Going all the way around the obstacle would be “worse”

Utility Function



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

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**

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**

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
- Histories have utility!

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
- Histories have utility!
- 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.

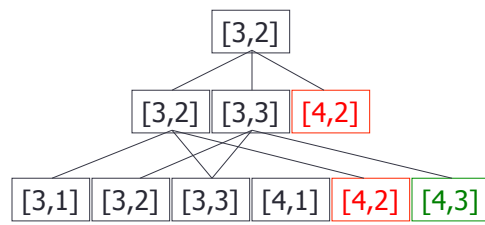
Utility of an Action Sequence

3				+1
2				-1
1				
	1	2	3	4

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

Utility of an Action Sequence

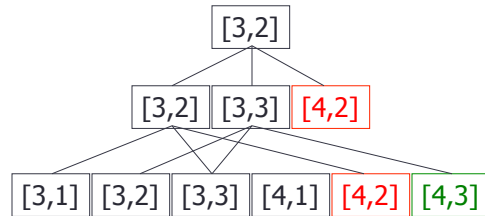
3				+1
2				-1
1				
	1	2	3	4



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

Utility of an Action Sequence

3				+1
2				-1
1				
	1	2	3	4

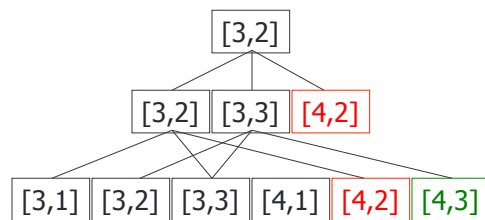


- Consider the action sequence (U,R) from [3,2]
- A run produces one of 7 possible histories, each with some probability
- The **utility of the sequence** is the expected utility of the histories:

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

Optimal Action Sequence

3				+1
2				-1
1				
	1	2	3	4

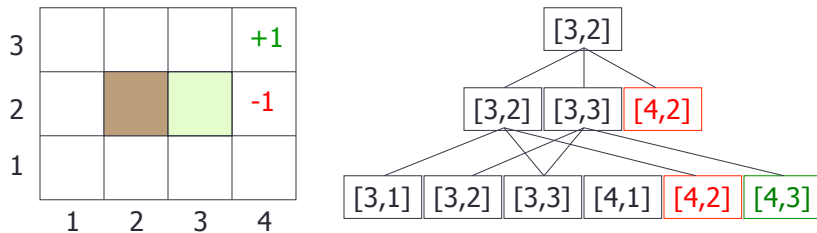


- Consider the action sequence (U,R) from [3,2]
- A run produces one of 7 possible histories, each with some probability
- The **utility of the sequence** is the expected utility of the histories:

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

- The **optimal sequence** is the one with maximal utility

Optimal Action Sequence



- Consider the action sequence (U,R) from [3,2]
- A run produced **only if the sequence is executed blindly!** utility
- 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?**

Reactive Agent Algorithm

Repeat:

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

Accessible or
observable state

Policy (Reactive/Closed-Loop Strategy)

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

- In every state, we need to know what to do
- The **goal** doesn't change
- A **policy (Π)** is a complete mapping from *states* to *actions*
 - "If in [3,2], go up; if in [3,1], go left; if in..."

Reactive Agent Algorithm

Repeat:

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

Optimal Policy

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

- A **policy** Π is a complete strategy for an agent.
- The **optimal policy** Π^* is the policy that maximizes the expected utility over all possible histories (sequence of states and actions).

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

Makes sense because of Markov property

Optimal Policy

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

- A **policy** Π is a complete strategy for an agent.
- 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^* ?

Additive Utility

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

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

- The reward accumulates as you step through states.

Reward

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 $i = 0, \dots, n-1$

Principle of Max Expected Utility

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

→	→	→	+1
↑		↑	-1
↑	←	←	←

*reminder! utility
of a sequence:*

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

First-step analysis →

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

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

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

Value Iteration

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

- For $t = 0, 1, 2, \dots$, do:

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

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

Policy Iteration

- Pick a policy Π at random
- Repeat:
 - Compute the utility of each state for Π

$$\mathbf{U}_{t+1}(i) \leftarrow \mathbf{R}(i) + \sum_k \mathbf{P}(k | \Pi(i), i) \mathbf{U}_t(k)$$

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

$$\Pi'(i) = \arg \max_a \sum_k \mathbf{P}(k | a, i) \mathbf{U}(k)$$

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

$$\Pi'(i) = \arg \max_a \sum_k \mathbf{P}(k | a.i) \mathbf{U}(k)$$

- If $\Pi' = \Pi$ then return Π

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)

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

Value Iteration: Summary

- 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: Summary

- 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!**