#### **Ecommerce:**

#### Agents and Rational Behavior Lecture 8: Decision-Making under Uncertainty Complex Decisions

#### Ralf Möller Hamburg University of Technology

#### Literature



Stuart Russell • Peter Norvig Prentice Hall Series in Artificial Intelligence

#### • Chapter 17

Material from Lise Getoor, Jean-Claude Latombe, Daphne Koller, and Russell

# **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)

#### **Markov Property**

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







#### **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
- There are 9 possible sequences of states

   called histories and 6 possible final states for the robot!



## **Utility Function**



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

## **Utility Function**



- [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**



- [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

# **Utility of a History**



- [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

#### **Utility of an Action Sequence**



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

### **Utility of an Action Sequence**



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

# **Utility of an Action Sequence**



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

$$\mathbf{U} = \Sigma_{h} \mathbf{U}_{h} \mathbf{P}(h)$$

## **Optimal Action Sequence**



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

## **Optimal Action Sequence**



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



#### **Policy** (Reactive/Closed-Loop Strategy)



• A policy  $\Pi$  is a complete mapping from states to actions

#### **Reactive Agent Algorithm**

Repeat:

- s ← sensed state
- If s is terminal then exit
- a ← Π(s)
- Perform a

# **Optimal Policy**



- A policy  $\Pi$  is a complet Note that [3,2] is a "dangerous"
- The optimal policy Π\* i history (ending at a teres to avoid
   The optimal policy Π\* i tries to avoid

expected utility

Makes sense because of Markov property

## **Optimal Policy**



- A policy Π is a comp • The entire leading This problem is called a ns • The entire leading This problem (MDP)
- The optimal policy T Markov Decision Problem (MDP) history with maximal expected utility

How to compute  $\Pi^*$ ?

# **Additive Utility**

- History  $H = (s_0, s_1, ..., s_n)$
- The utility of H is additive iff:  $\mathbf{U}(s_0, s_1, \dots, s_n) = \mathbf{R}(0) + \mathbf{U}(s_1, \dots, s_n) = \sum \mathbf{R}(i)$

Reward

# **Additive Utility**

- History  $H = (S_0, S_1, ..., S_n)$
- The utility of H is additive iff:  $\mathbf{U}(s_0, s_1, \dots, s_n) = \mathbf{R}(0) + \mathbf{U}(s_1, \dots, s_n) = \Sigma \mathbf{R}(0)$
- Robot navigation example:

• 
$$\mathbf{R}(n) = +1 \text{ if } \mathbf{S}_n = [4,3]$$

• 
$$\mathbf{R}(n) = -1 \text{ if } \mathbf{S}_n = [4,2]$$

• 
$$\mathbf{R}(i) = -1/25$$
 if  $i = 0, ..., n-1$ 

#### **Principle of Max Expected Utility**

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



First-step analysis  $\rightarrow$ 

- $\mathbf{U}(i) = \mathbf{R}(i) + \max_{a} \sum_{k} \mathbf{P}(k \mid a.i) \mathbf{U}(k)$
- $\Pi^*(i) = \arg \max_a \sum_k \mathbf{P}(k \mid a.i) \mathbf{U}(k)$



- Initialize the utility of each non-terminal state  $s_i$  to  $U_o(i) = 0$
- For t = 0, 1, 2, ..., do:  $U_{t+1}(i) \in \mathbf{R}(i) + \max_{a} \sum_{k} \mathbf{P}(k \mid a.i) U_{t}(k)$



## **Value Iteration**

• Initialize the utility of each  $\mathbf{U}_{0}(\mathbf{i}) = 0$ Note the importance of terminal states and connectivity of the state-transition graph

For t = 0, 1, 2, ..., do:  

$$\mathbf{U}_{t+1}(i) \in \mathbf{R}(i) + \max_{a} \sum_{k} \mathbf{P}(k \mid a.i) \mathbf{U}_{t}(k)$$



• Pick a policy  $\Pi$  at random

- Pick a policy Π at random
- Repeat:
  - Compute the utility of each state for  $\Pi$  $U_{t+1}(i) \in \mathbf{R}(i) + \sum_{k} \mathbf{P}(k \mid \Pi(i).i) U_{t}(k)$

- Pick a policy Π at random
- Repeat:
  - Compute the utility of each state for  $\Pi$  $U_{t+1}(i) \in \mathbf{R}(i) + \sum_{k} \mathbf{P}(k \mid \Pi(i).i) U_{t}(k)$
  - Compute the policy Π' given these utilities

 $\Pi'(i) = \arg \max_{a} \sum_{k} \mathbf{P}(k \mid a.i) \mathbf{U}(k)$ 

- Pick a policy Π at random
- Repeat:
  - Compute the utility of each state for  $\Pi$  $\mathbf{U}_{t+1}(i) \in \mathbf{R}(i) + \sum_{k} \mathbf{P}(k \mid \Pi(i).i) \mathbf{U}_{t}(k)$
  - Compute the policy  $\Pi$ ' given these utilities  $\Pi$ '(i) = arg max<sub>a</sub>  $\Sigma_k P($  $U(i) = R(i) + \Sigma_k P(k \mid \Pi(i).i) U(k)$ (often a sparse system)
  - If  $\Pi' = \Pi$  then return  $\Pi$

#### n-Step decision process



Assume that:

- Each state reached after n steps is terminal, hence has known utility
- There is a single initial state
- Any two states reached after i and j steps are different

#### n-Step Decision Process



$$\Pi^{*}(i) = \arg \max_{a} \sum_{k} \mathbf{P}(k \mid a.i) \mathbf{U}(k)$$
$$\mathbf{U}(i) = \mathbf{R}(i) + \max_{a} \sum_{k} \mathbf{P}(k \mid a.i) \mathbf{U}(k)$$

For j = n-1, n-2, ..., 0 do:

For every state S<sub>i</sub> attained after step j

- Compute the utility of S<sub>i</sub>
- Label that state with the corresponding action

#### What is the Difference?





 $\Pi^{*}(i) = \arg \max_{a} \Sigma_{k} \mathbf{P}(k \mid a.i) \mathbf{U}(k)$  $\mathbf{U}(i) = \mathbf{R}(i) + \max_{a} \Sigma_{k} \mathbf{P}(k \mid a.i) \mathbf{U}(k)$ 

#### **Infinite Horizon**

In many problems, e.g., the robot navigat What if the robot lives forever? potentially unbounded and the same state can be reach One trick:



Use discounting to make infinite Horizon problem mathematically tractable

## **Example: Tracking a Target**



#### **POMDP** (Partially Observable Markov Decision Problem)

- A sensing operation returns multiple states, with a probability distribution
- Choosing the action that maximizes the expected utility of this state distribution assuming "state utilities" computed as above is not good enough, and actually does not make sense (is not rational)

#### **Example: Target Tracking**

There is uncertainty in the robot's and target's positions; this uncertainty grows with further motion

> There is a risk that the target may escape behind the corner, requiring the robot to move appropriately

But there is a positioning landmark nearby. Should the robot try to reduce its position uncertainty?



## Summary

- Decision making under uncertainty
- Utility function
- Optimal policy
- Maximal expected utility
- Value iteration
- Policy iteration