#### **Artificial Intelligence**

Markov decision processes (MDPs)

#### Example: HMM with no actions



# What's missing in HMMs

- HMMs cannot model important aspects of agent interactions:
  - No model for rewards
  - No model for actions which can affect these rewards
- These are actually issues that are faced by many applications:
  - Agents negotiating deals on the web
  - A robot which interacts with its environment

## **Example: No actions**



#### Formal definition of MDPs

- A set of states  $\{s_1 \dots s_n\}$
- A set of rewards {r<sub>1</sub> ... r<sub>n</sub>} \*
- A set of action {a<sub>1</sub> .. a<sub>m</sub>} ←
- Transition probability

One reward for each state

Number of actions could be larger than number of states

$$P_{i,j}^{k} = P(q_{t+1} = s_j \mid q_t = s_i \& h_t = a_k)$$

### Questions

- What is my expected payoff if I am in state *i*
- What is my expected payoff if I am in state *i* and perform action *a*?

# Solving MDPs

- No actions: Value iterations
- With actions: Value iteration, Policy iteration

# Value computation

- An obvious question for such models is what is combined expected value for each state
- What can we expect to earn over our life time if we become Asst. prof.?
- What if we go to industry?

Before we answer this question, we need to define a model for future rewards:

- The value of a current award is higher than the value of future awards
  - Inflation, confidence
  - Example: Lottery

#### **Discounted rewards**

- The discounted rewards model is specified using a parameter  $\boldsymbol{\gamma}$
- Total rewards = current reward +

 $\gamma$  (reward at time t+1) +  $\gamma^2$  (reward at time t+2) +

 $\gamma^{k}$  (reward at time t+k) +

infinite sum

. . . . . . . . . .

#### **Discounted awards**

- The discounted award model is specified using a parameter  $\boldsymbol{\gamma}$
- Total awards = current award +

 $\gamma$  (award at time t+1) +

 $v^2$  (award at time t+2) +

Converges if  $0 < \gamma < 1$  () +

infinite sum

# Determining the total rewards in a state

- Define J\*(s<sub>i</sub>) = expected discounted sum of rewards when starting at state s<sub>i</sub>
- How do we compute  $J^*(s_i)$ ?

$$J^{*}(s_{i}) = r_{i} + \gamma X$$
  
=  $r_{i} + \gamma (p_{i1}J^{*}(s_{1}) + p_{i2}J^{*}(s_{2}) + \cdots + p_{in}J^{*}(s_{n}))$ 

How can we solve this?

#### Computing j\*(s<sub>i</sub>)

$$J^{*}(s_{1}) = r_{1} + \gamma(p_{11}J^{*}(s_{1}) + p_{12}J^{*}(s_{2}) + \cdots + p_{1n}J^{*}(s_{n}))$$

$$J^*(s_2) = r_2 + \gamma(p_{21}J^*(s_1) + p_{22}J^*(s_2) + \cdots + p_{2n}J^*(s_n))$$

$$J^{*}(s_{n}) = r_{n} + \gamma(p_{n}J^{*}(s_{1}) + p_{n}J^{*}(s_{2}) + \cdots + p_{n}J^{*}(s_{n}))$$

- We have n equations with n unknowns
- Can be solved in closed form

#### **Iterative approaches**

- Solving in closed form is possible, but may be time consuming.
- Alternatively, this problem can be solved in an iterative manner
- Let's define J<sup>t</sup>(s<sub>i</sub>) as the expected total discounted rewards after t steps
- How can we compute  $J^t(s_i)$ ?

$$J^{1}(S_{i}) = r_{i}$$
$$J^{2}(S_{i}) = r_{i} + \gamma \left(\sum_{k} p_{i,k} J^{1}(s_{k})\right)$$
$$J^{t+1}(S_{i}) = r_{i} + \gamma \left(\sum_{k} p_{i,k} J^{t}(s_{k})\right)$$

#### **Iterative approaches**

• We know how to solve this!

Lets fill the dynamic programming table

- Lets define  $J^{\kappa}(s_i)$  as the expected discounted awards after k steps
- But wait ...

This is a never ending task!

$$J^{2}(S_{i}) = r_{i} + \gamma \left(\sum_{k} p_{i,k} J^{1}(s_{k})\right)$$
$$J^{t+1}(S_{i}) = r_{i} + \gamma \left(\sum_{k} p_{i,k} J^{t}(s_{k})\right)$$

#### When do we stop?

$$J^{1}(S_{i}) = r_{i}$$
$$J^{2}(S_{i}) = r_{i} + \gamma \left(\sum_{k} p_{i,k} J^{1}(s_{k})\right)$$
$$J^{t+1}(S_{i}) = r_{i} + \gamma \left(\sum_{k} p_{i,k} J^{t}(s_{k})\right)$$

Remember, we have a converging function

We can stop when  $|J^{t-1}(s_i) - J^t(s_i)|_{\infty} < \epsilon$ 

Infinity norm selects maximal element



# Solving MDPs

- No actions: Value iterations  $\sqrt{}$
- With actions: Value iteration, Policy iteration

# Adding actions

A Markov Decision Process:

- A set of states  $\{s_1 \dots s_n\}$
- A set of rewards  $\{r_1 \dots r_n\}$
- A set of action  $\{a_1 \dots a_m\}$
- Transition probability

$$P_{i,j}^{k} = P(q_{t+1} = s_{j} \mid q_{t} = i \& h_{t} = a_{k})$$



#### **Questions for MDPs**

- Now we have actions
- The question changes to the following:

Given our current state and the possible actions, what is the best action for us in terms of long term payment?



# Policy

- A policy maps states to actions
- An optimal policy leads to the highest expected returns
- Note that this does not depend on the start state
- How is this different from the Planning solutions we studies 2 weeks ago?

Gr	В
Go	A
Asst. Pr.	A
Ten. Pr.	В

# Solving MDPs with actions

- It could be shown that for every MDP there exists an optimal policy (we won't discuss the proof).
- Such policy guarantees that there is no other action that is expected to yield a higher payoff

#### Computing the optimal policy: 1. Modified value iteration

- We can compute it by modifying the value iteration method we discussed.
- Define p<sup>k</sup><sub>ij</sub> as the probability of transitioning from state i to state j when using action k
- Then we compute:

$$J^{t+1}(S_i) = \max_k r_i + \gamma \left( \sum_j p_{i,j}^k J^t(s_j) \right)$$
  
Also known as Bellman's equation

#### Computing the optimal policy: 1. Modified value iteration

- We can compute it by modifying the value iteration method we discussed.
- Define p<sup>k</sup><sub>ij</sub> as the probability of transitioning from state i to state j when using action k
- Then we compute:

$$J^{t+1}(S_i) = \max_k r_i + \gamma \left(\sum_j p_{i,j}^k J^t(s_j)\right)$$

Run until convergences

#### Computing the optimal policy: 1. Modified value iteration

- We can compute it by modifying the value iteration method we discussed.
- Define p<sup>k</sup><sub>ij</sub> as the probability of transitioning from state i to state j when using action k
- Then we compute:

$$J^{t+1}(S_i) = \max_k r_i + \gamma \left(\sum_j p_{i,j}^k J^t(s_j)\right)$$

• When the algorithm converges, we have computed the best outcome for each state

• We associate states with the actions that maximize their return

#### Value iteration for $\gamma$ =0.9



#### Computing the optimal policy: 2. Policy iteration

- We can also compute optimal policies by revising an existing policy.
- We initially select a policy at random (mapping from states to actions).
- We re-compute the expected long term reward at each state using the selected policy
- We select a new policy using the expected rewards and iterate until convergences

## Policy iteration: algorithm

- Let  $\pi_t(s_i)$  be the selected policy at time t
- 1. Randomly choose  $\pi_0$ ; set t = 0
- 2. For each state  $s_i$  compute  $J^*(s_i)$ , the long term expected reward using policy  $\pi_t$ .
- expected reward using policy  $\pi_t$ . 3. Set  $\pi_t(s_i) = \max_k r_i + \gamma \left( \sum_j p_{i,j}^k J^*(s_j) \right)$
- 4. Convergence? Yes: output policy. No: t = t + 1, go to 2.

## Policy iteration: algorithm

- Let  $\pi_t(s_i)$  be the selected policy at time t for state i
- 1. Randomly choose  $\pi_0$ ; set t = 0
- 2. For each state  $s_i$  compute  $J^*(s_i)$ , the long term expected reward using policy  $\pi_t$ .

3. set 
$$\pi_{t+1}(s_i) = \max_k r_i + \gamma \left( \sum_j p_{i,j}^k J^*(s_j) \right)$$

4. Convergence? Yes: output policy. No: t = t + 1, go to 2.

Can be computed using  $J^*(s_i)$  for all states

Can be computed using value iteration

#### Value iteration vs. policy iteration

- Depending on the model and the information at hand:
  - If you have a good guess regarding the optimal policy then policy iteration would converge much faster
  - similarly, if there are many possible actions, policy iteration might be faster
  - otherwise value iteration is a safer way

# What you should know

- Models that include rewards and actions
- Value iteration for solving MDPs
- Policy iteration

#### Partially Observed Markov Decision Processes (POMDPs)

- Same model as MDP except: We do not observe the states we are in.
- Thus, we have a distribution over states
- There is an initial distribution for states (initial belief)
- Once we reach a new state and receive a reward we can re-compute a new belief regrading the possible set of states

#### Example

- If we see 1, we can be in any of several locations.
- However, based on past and future observations we can increase a decrease our belief at a given state

1	1	1
3	1	2
1	2	1

POMDPs can be solved by extending the MDP methods to solve for a belief state vector rather than for the original single state MDP