# CS 730/730W/830: Intro AI

- **Class Outline**

MDPs

Solving MDPs

# Class Outline

1. search: heuristics, CSPs, games
2. knowledge representation: FOL, resolution
3. planning: STRIPS, MDPs
4. learning: supervised, unsupervised
5. uncertainty: particle filters, HMMs

# Markov Decision Processes

# Examples

1. robot navigation

2. driving

3. business

4. war

5. diagnosis

6. life

# Probability

propositional

domain: discrete or continuous

0–1, sum to 1

distribution of continuous $=$ density

$E(X) = \int x \, pdf(x) \, dx$

$P(X = x_1)$ written as $P(x_1)$ or if $X$ is true/false, $P(x)$

conditional ($=$posterior): $P(x|y) = P(x \wedge y)/P(y)$

# Markov Decision Process (MDP)

**initial state:** $s_0$

**transition model:** $T(s, a, s') =$ probability of going from $s$ to $s'$ after doing $a$.

**reward function:** $R(s)$ for landing in state $s$.

**terminal states:** sinks = absorbing states (end the trial).

# Markov Decision Process (MDP)

**initial state:** $s_0$

**transition model:** $T(s, a, s') = $ probability of going from $s$ to $s'$ after doing $a$.

**reward function:** $R(s)$ for landing in state $s$.

**terminal states:** sinks = absorbing states (end the trial).

objective:

**total reward:** reward over (finite) trajectory:
$R(s_0) + R(s_1) + R(s_2)$

**discounted reward:** penalize future by $\gamma$:
$R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) \ldots$

# Markov Decision Process (MDP)

**initial state:** $s_0$

**transition model:** $T(s, a, s') = $ probability of going from $s$ to $s'$ after doing $a$.

**reward function:** $R(s)$ for landing in state $s$.

**terminal states:** sinks $=$ absorbing states (end the trial).

objective:

**total reward:** reward over (finite) trajectory:
$R(s_0) + R(s_1) + R(s_2)$

**discounted reward:** penalize future by $\gamma$:
$R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) \ldots$

find:

**policy:** $\pi(s) = a$

**optimal policy:** $\pi^*$

**proper policy:** reaches terminal state

# What to do?

$$\pi^*(s) =$$

# What to do?

$$\pi^*(s) = \operatorname*{argmax}_a \sum_{s'} T(s, a, s') U^{\pi^*}(s')$$

$$U^\pi(s) =$$

# What to do?

$$\pi^*(s) = \operatorname*{argmax}_a \sum_{s'} T(s, a, s') U^{\pi^*}(s')$$

$$U^\pi(s) = E[\sum_{t=0}^{\infty} \gamma^t R(s_t) | \pi, s_0 = s]$$

# What to do?

$$\pi^*(s) = \operatorname*{argmax}_a \sum_{s'} T(s, a, s') U^{\pi^*}(s')$$

$$U^\pi(s) = E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi, s_0 = s\right]$$

The key:

$$U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$$

(Richard Bellman, 1957)

# Break

- asst 9 due Mon Apr 7
- projects: reading is fine, but probably best to talk with me before starting serious work

# Solving MDPs

# Value Iteration

Repeated Bellman updates:

Repeat until happy
    for each state $s$
        $U'(s) \leftarrow R(s) + \gamma \max_a \sum_{s'} T(s, a, s')U(s')$
    $U \leftarrow U'$

# Value Iteration

Repeated Bellman updates:

Repeat until happy
    for each state $s$
        $U'(s) \leftarrow R(s) + \gamma \max_a \sum_{s'} T(s, a, s')U(s')$
    $U \leftarrow U'$

For infinite updates everywhere, guaranteed to reach equilibrium.

Equilibrium is unique solution to Bellman equations!

asychronous works: converges if every state updated infinitely often (no state permanently ignored)

- What question didn't you get to ask today?
- What's still confusing?
- What would you like to hear more about?

Please write down your most pressing question about AI and put it in the box on your way out.
*Thanks!*