## Slide 1

# CS 188: Artificial Intelligence

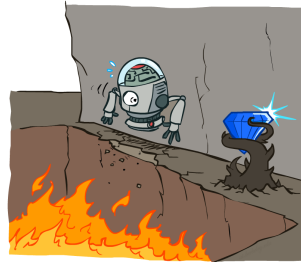## Markov Decision Processes I



Instructor: Marco Alvarez

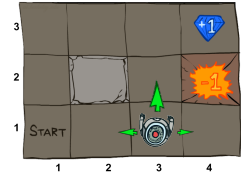University of Rhode Island

1

## Slide 2

# Non-Deterministic Search



2

## Slide 3

# Example: Grid World

- A maze-like problem
  - The agent lives in a grid
  - Walls block the agent's path
- Noisy movement: actions do not always go as planned
  - 80% of the time, the action North takes the agent North
    (if there is no wall there)
  - 10% of the time, North takes the agent West; 10% East
  - If there is a wall in the direction the agent would have been taken, the agent stays put
- The agent receives rewards each time step
  - Small "living" reward each step (can be negative)
  - Big rewards come at the end (good or bad)
- Goal: maximize sum of rewards



3

## Slide 4

# Grid World Actions
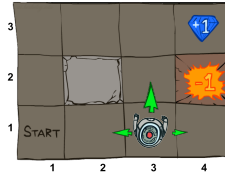
Deterministic Grid World

Stochastic Grid World



4

## Slide 5

# Markov Decision Processes

- An MDP is defined by:
  - A set of states $s \in S$
  - A set of actions $a \in A$
  - A transition function $T(s, a, s')$
    - Probability that a from s leads to s', i.e., $P(s' | s, a)$
    - Also called the model or the dynamics
  - A reward function $R(s, a, s')$
    - Sometimes just $R(s)$ or $R(s')$
  - A start state
  - Maybe a terminal state

- MDPs are non-deterministic search problems
  - One way to solve them is with expectimax search
  - We'll have a new tool soon



[Demo – gridworld manual intro (L8D1)]

5

## Slide 6

# Video of Demo Gridworld Manual Intro



6

## Slide 7

# What is Markov about MDPs?

- "Markov" generally means that given the present state, the future and the past are independent

- For Markov decision processes, "Markov" means action outcomes depend only on the current state

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1}, \ldots S_0 = s_0)$$
$$=$$
$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

- This is just like search, where the successor function could only depend on the current state (not the history)
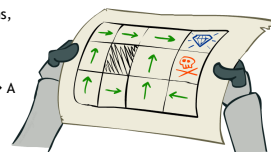
Andrey Markov
(1856-1922)
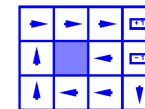
7

## Slide 8

# Policies

- In deterministic single-agent search problems, we wanted an optimal plan, or sequence of actions, from start to a goal

- For MDPs, we want an optimal policy $\pi^*$: S → A
  - A policy $\pi$ gives an action for each state
  - An optimal policy is one that maximizes expected utility if followed
  - An explicit policy defines a reflex agent

- Expectimax didn't compute entire policies
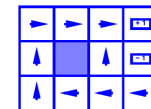  - It computed the action for a single state only



Optimal policy when R(s, a, s') = -0.03 for all non-terminals s
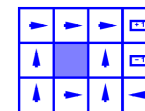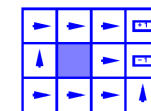
8

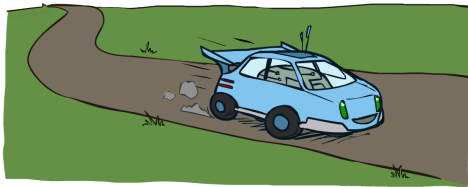## Slide 9

# Optimal Policies
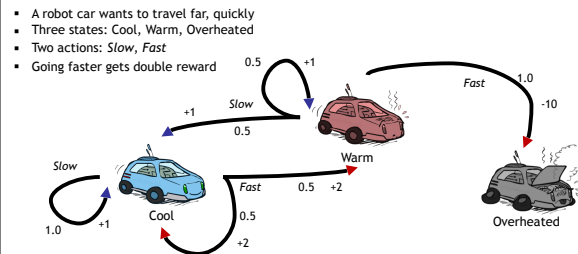


R(s) = -0.01

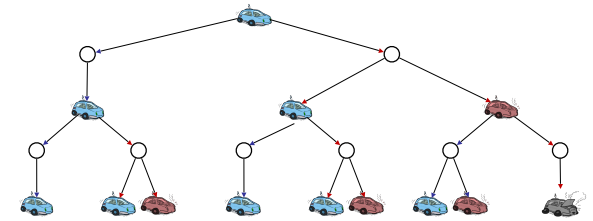R(s) = -0.03

R(s) = -0.4

R(s) = -2.0

9

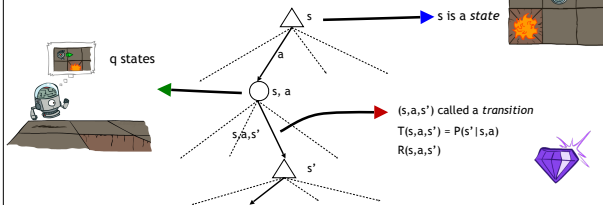## Example: Racing



10

## Example: Racing

- A robot car wants to travel far, quickly
- Three states: Cool, Warm, Overheated
- Two actions: *Slow*, *Fast*
- Going faster gets double reward



11

## Racing Search Tree



12

## MDP Search Trees

- Each MDP state projects an expectimax-like search tree



s is a *state*

q states

(s,a,s') called a *transition*

$T(s,a,s') = P(s'|s,a)$

$R(s,a,s')$

13

## Utilities of Sequences



14

## Utilities of Sequences

- What preferences should an agent have over reward sequences?

- More or less?   [1, 2, 2]   or   [2, 3, 4]

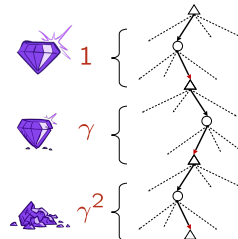- Now or later?   [0, 0, 1]   or   [1, 0, 0]



15

## Discounting

- It's reasonable to maximize the sum of rewards
- It's also reasonable to prefer rewards now to rewards later
- One solution: values of rewards decay exponentially



$1$     $\gamma$     $\gamma^2$

Worth Now    Worth Next Step    Worth In Two Steps

16

## Discounting

- How to discount?
    - Each time we descend a level, we multiply in the discount once

- Why discount?
    - Sooner rewards probably do have higher utility than later rewards
    - Also helps our algorithms converge

- Example: discount of 0.5
    - U([1,2,3]) = 1*1 + 0.5*2 + 0.25*3
    - U([1,2,3]) < U([3,2,1])



$1$

$\gamma$

$\gamma^2$

17

## Stationary Preferences

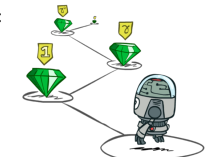- Theorem: if we assume stationary preferences:

$$[a_1, a_2, \ldots] \succ [b_1, b_2, \ldots]$$
$$\Updownarrow$$
$$[r, a_1, a_2, \ldots] \succ [r, b_1, b_2, \ldots]$$



- Then: there are only two ways to define utilities
    - Additive utility: $U([r_0, r_1, r_2, \ldots]) = r_0 + r_1 + r_2 + \cdots$
    - Discounted utility: $U([r_0, r_1, r_2, \ldots]) = r_0 + \gamma r_1 + \gamma^2 r_2 \cdots$

18

## Quiz: Discounting

- Given:

| 10 | | | ♡ | 1 |
|---|---|---|---|---|
| a | b | c | d | e |

  - Actions: East, West, and Exit (only available in exit states a, e)
  - Transitions: deterministic

- Quiz 1: For $\gamma = 1$, what is the optimal policy?

| 10 | | | ♡ | 1 |
|---|---|---|---|---|

- Quiz 2: For $\gamma = 0.1$, what is the optimal policy?

| 10 | | | ♡ | 1 |
|---|---|---|---|---|

- Quiz 3: For which $\gamma$ are West and East equally good when in state d?

## Infinite Utilities?!

- Problem: What if the game lasts forever? Do we get infinite rewards?

- Solutions:
  - Finite horizon: (similar to depth-limited search)
    - Terminate episodes after a fixed T steps (e.g. life)
    - Gives nonstationary policies ($\pi$ depends on time left)
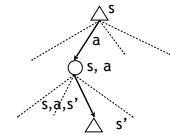  - Discounting: use $0 < \gamma < 1$
    $$U([r_0, \dots r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{\max}/(1 - \gamma)$$
    - Smaller $\gamma$ means smaller "horizon" – shorter term focus

  - Absorbing state: guarantee that for every policy, a terminal state will eventually be reached (like "overheated" for racing)
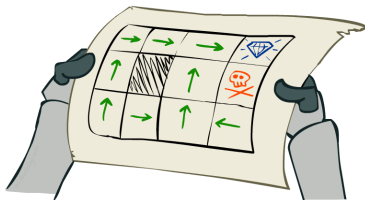
## Recap: Defining MDPs

- Markov decision processes:
  - Set of states S
  - Start state $s_0$
  - Set of actions A
  - Transitions P(s'|s,a) (or T(s,a,s'))
  - Rewards R(s,a,s') (and discount $\gamma$)

- MDP quantities so far:
  - Policy = Choice of action for each state
  - Utility = sum of (discounted) rewards

## Solving MDPs