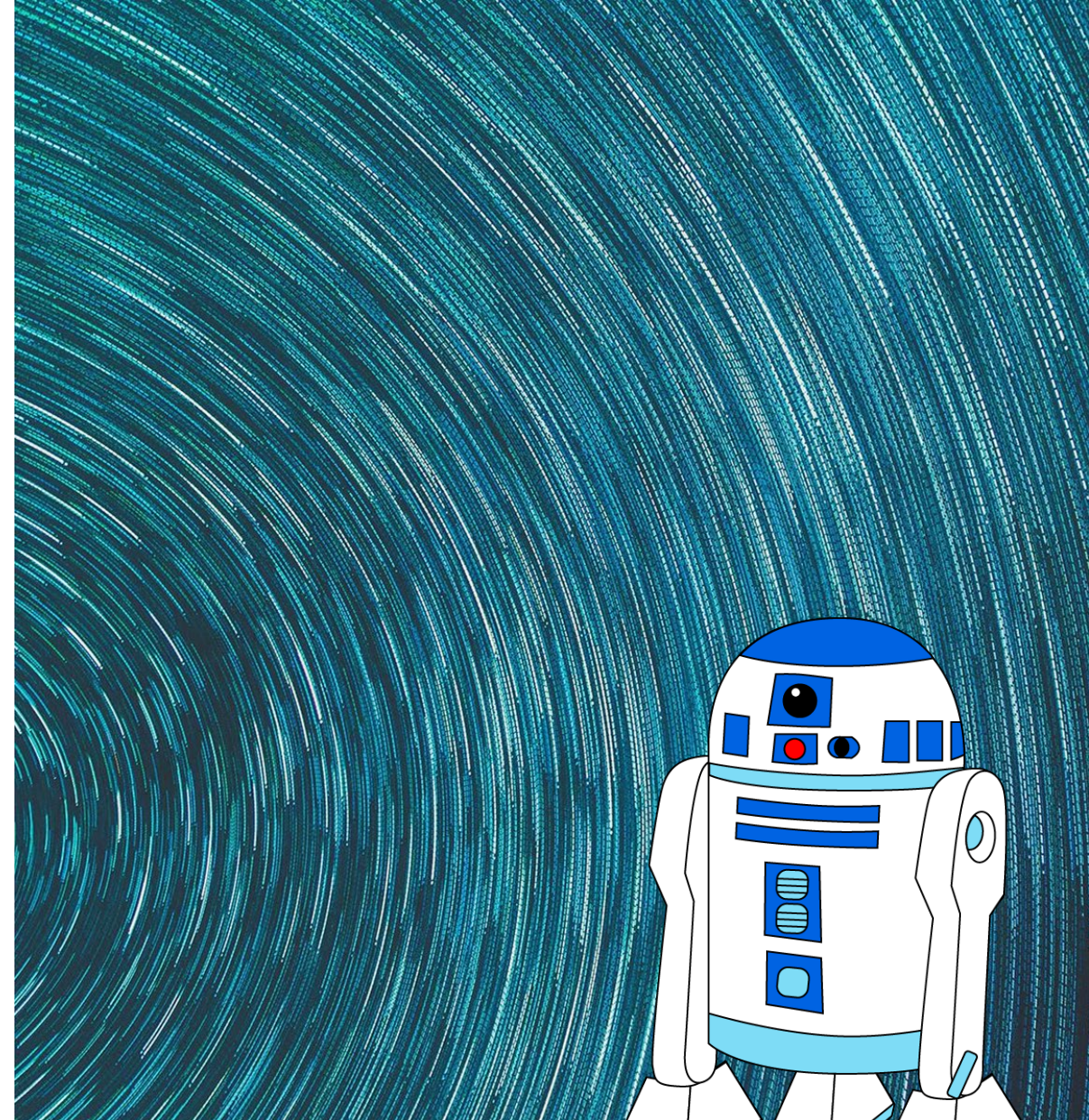CIS 521:
ARTIFICIAL INTELLIGENCE

# Expectimax and Utilities
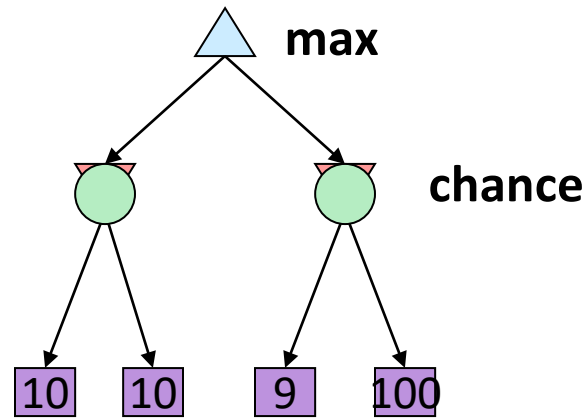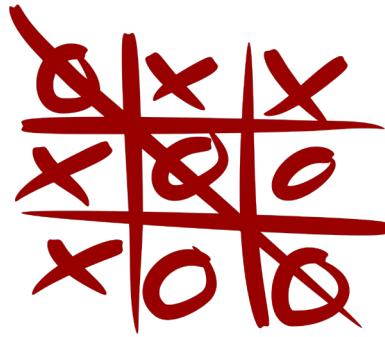
Harry Smith

Many of today's slides are courtesy of Dan Klein and
Pieter Abbeel of University of California, Berkeley

Penn
Engineering
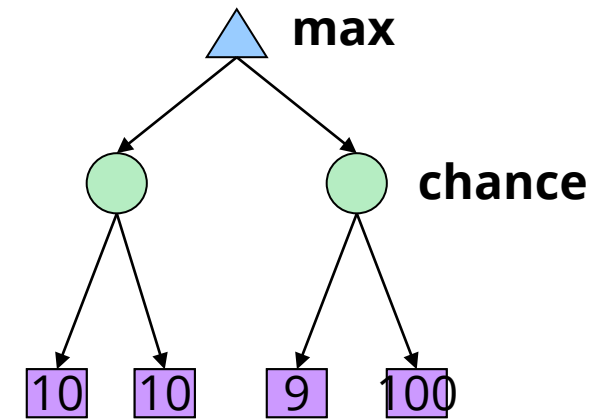UNIVERSITY of PENNSYLVANIA

# Uncertain Outcomes

Idea: Uncertain outcomes controlled by chance, not an adversary!

# Expectimax Search

o Why wouldn't we know what the result of an action will be?
- Explicit randomness: rolling dice
- Unpredictable opponents: the opponent isn't optimal
- Actions can fail: when moving a robot, wheels might slip

o Values should now reflect average-case (expectimax) outcomes, not worst-case (minimax) outcomes

o Expectimax search: compute the average score under optimal play
- Max nodes as in minimax search
- Chance nodes are like min nodes but the outcome is uncertain
- Calculate their expected utilities
- I.e. take weighted average (expectation) of children

o Later, we'll learn how to formalize the underlying uncertain-result problems as Markov Decision Processes

# Expectimax Pseudocode

def value(state):
    if the state is a terminal state: return the state's utility
    if the next agent is MAX: return max-value(state)
    if the next agent is EXP: return exp-value(state)

def max-value(state):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor))
    return v

def exp-value(state):
    initialize v = 0
    for each successor of state:
        p = probability(successor)
        v += p * value(successor)
    return v

# Expectimax Pseudocode

def exp-value(state):
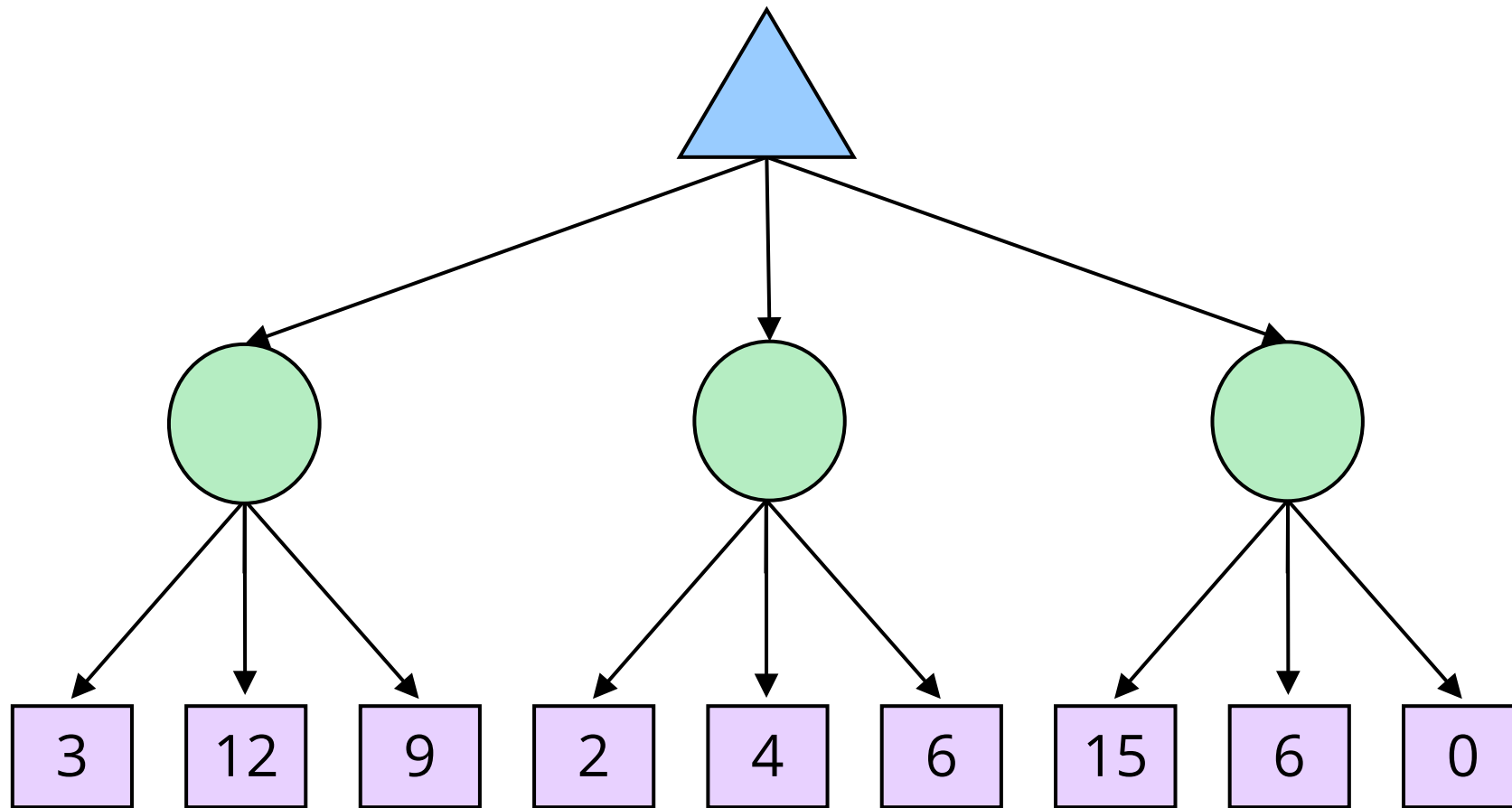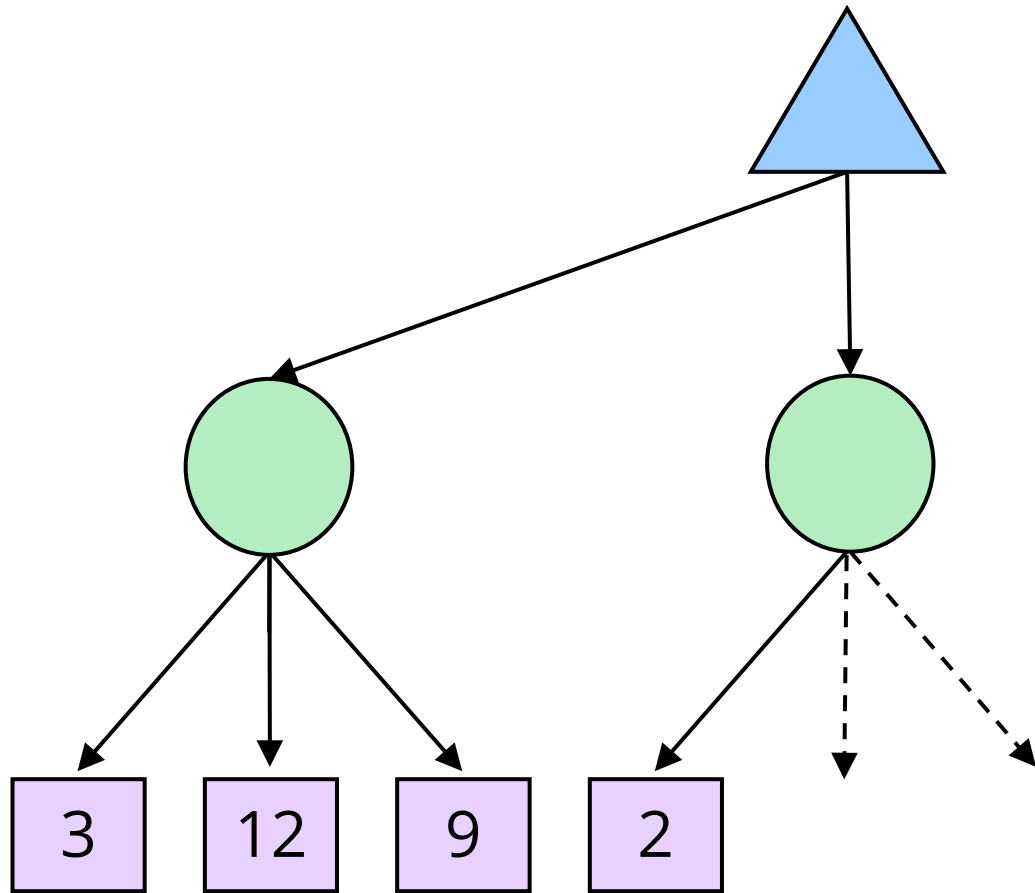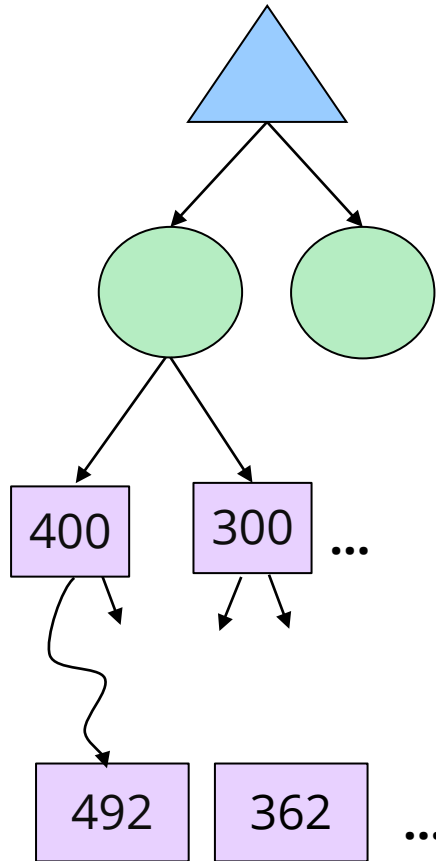    initialize v = 0
    for each successor of state:
        p = probability(successor)
        v += p * value(successor)
    return v



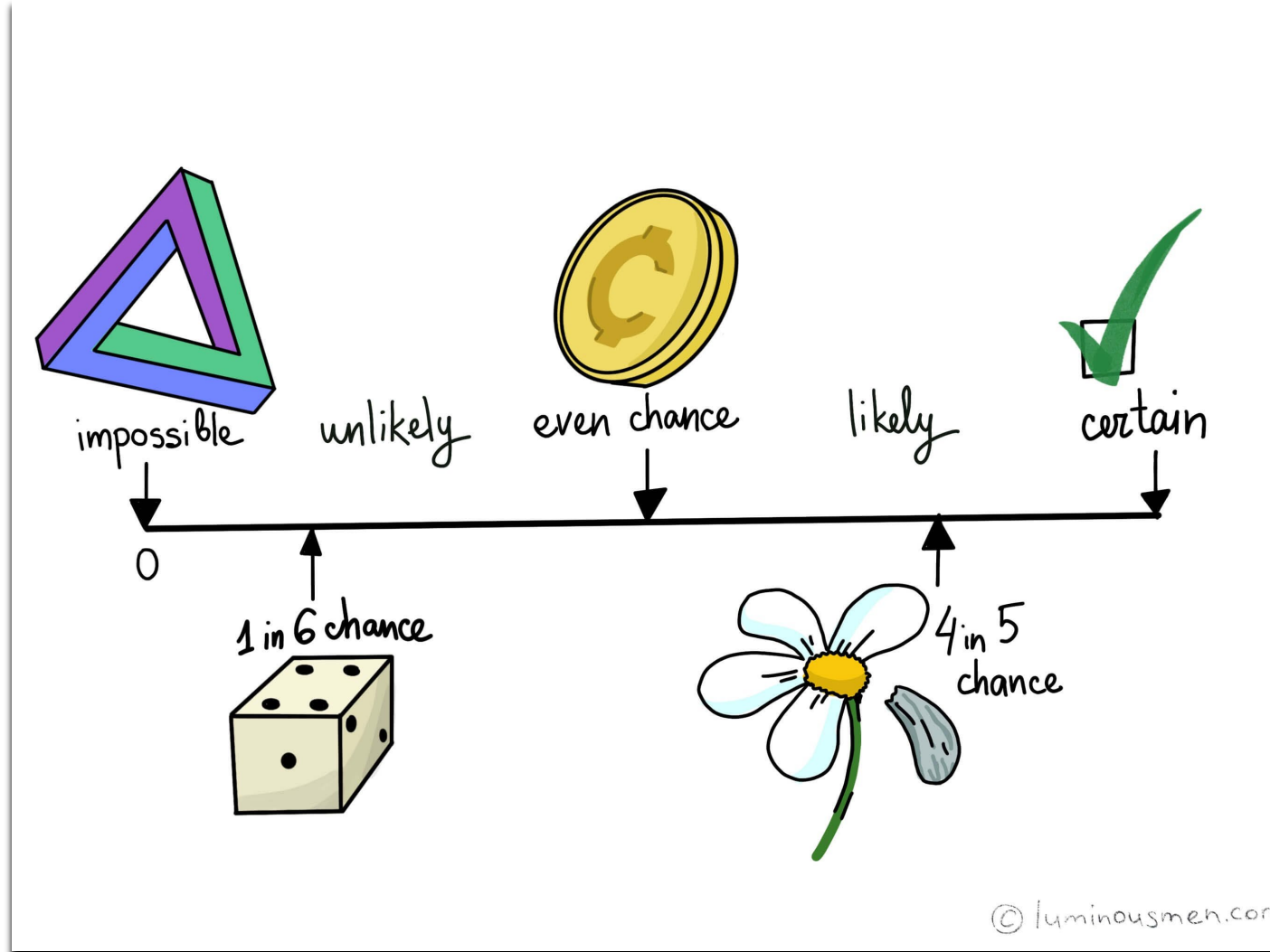$$v = \frac{1}{2} \cdot (8) + \frac{1}{3} \cdot (24) + \frac{1}{6} \cdot (-12)$$

Estimate of true expectimax value (which would require a lot of work to compute)
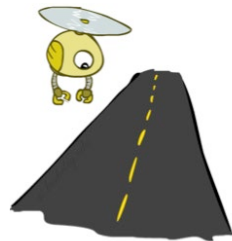
# Probabilities

# Probabilities

o   A random variable represents an event whose outcome is unknown

o   A probability distribution is an assignment of weights to outcomes

o   Example: Traffic on freeway
   - Random variable: T = whether there's traffic
   - Outcomes: T in {none, light, heavy}
   - Distribution: $P(T = \text{none}) = 0.25$, $P(T = \text{light}) = 0.50$, $P(T = \text{heavy}) = 0.25$

o   Some laws of probability (more later):
   - Probabilities are always non-negative
   - Probabilities over all possible outcomes sum to one

o   As we get more evidence, probabilities may change:
   - $P(T = \text{heavy}) = 0.25$, $P(T = \text{heavy} \mid \text{Hour} = 8\text{am}) = 0.60$
   - We'll talk about methods for reasoning and updating probabilities later
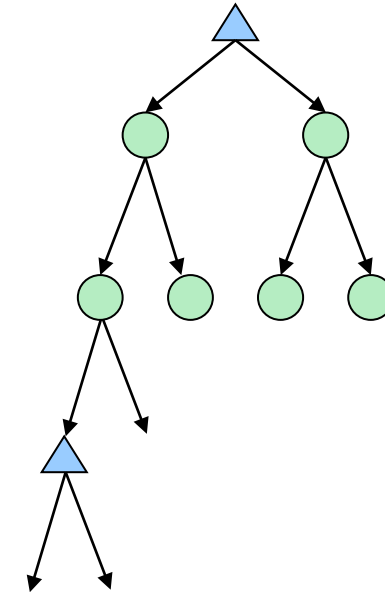
0.25

0.50

0.25

# Probabilities

o  The expected value of a function of a random variable is the average, weighted by the probability distribution over outcomes

o  Example: How long to get to the airport?

Time:  20 min  +  30 min  +  60 min  ⟹  35 min
       x             x             x

Probability:  0.25        0.50        0.25

# What Probabilities to Use?

o **In expectimax search, we have a probabilistic model of how the opponent (or environment) will behave in any state**

- Model could be a simple uniform distribution (roll a die)

- Model could be sophisticated and require a great deal of computation

- We have a chance node for any outcome out of our control: opponent or environment

- The model might say that adversarial actions are likely!

o **For now, assume each chance node magically comes along with probabilities that specify the distribution over its outcomes**

*Having a probabilistic belief about another agent's action does not mean that the agent is flipping any coins!*

o **Objectivist / frequentist answer:**
   - Averages over repeated *experiments*
   - E.g. empirically estimating P(rain) from historical observation
   - Assertion about how future experiments will go (in the limit)
   - New evidence changes the *reference class*
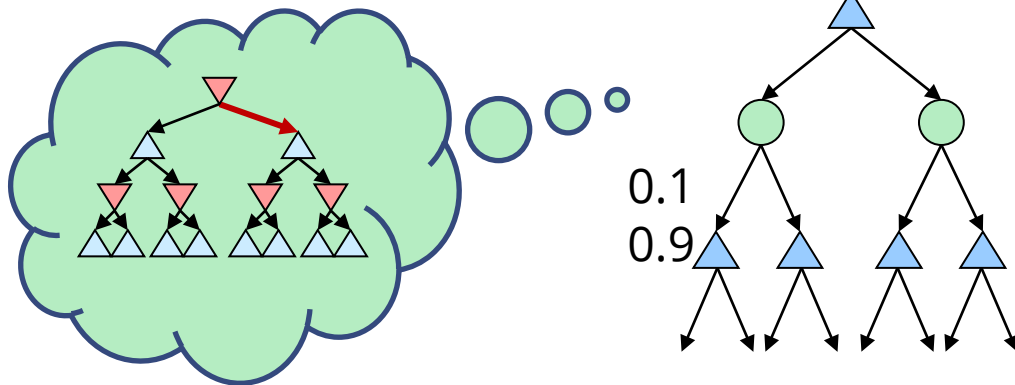   - Makes one think of *inherently random* events, like rolling dice

o **Subjectivist / Bayesian answer:**
   - Degrees of belief about unobserved variables
   - E.g. an agent's belief that it's raining, given the temperature
   - E.g. agent's belief how an opponent will behave, given the state
   - Often *learn* probabilities from past experiences (more later)
   - New evidence *updates beliefs* (more later)

# Quiz: Informed Probabilities

o   Let's say you know that your opponent is actually running a depth 2 minimax, using the result 80% of the time, and moving randomly otherwise

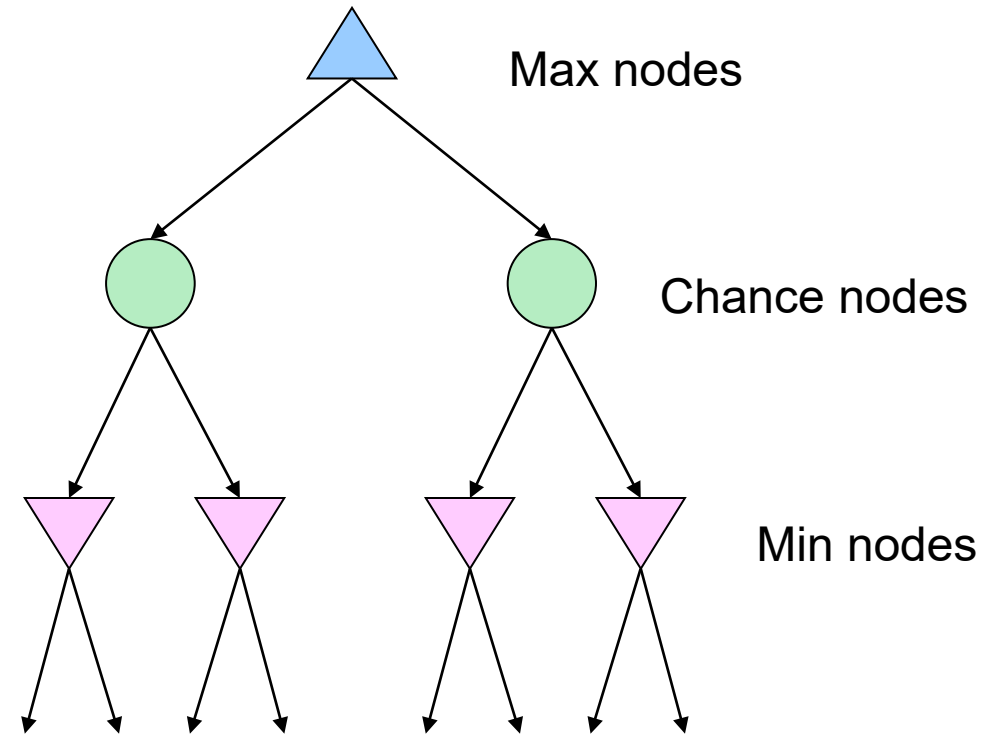o   Question: What tree search should you use?

0.1
0.9

- Answer: Expectimax!
  - To figure out EACH chance node's probabilities, you have to run a simulation of your opponent
  - This kind of thing gets very slow very quickly
  - Even worse if you have to simulate your opponent simulating you…
  - … except for minimax, which has the nice property that it all collapses into one game tree

o **Dice rolls increase _b_: 21 possible rolls with 2 dice**

   ▪ Backgammon $\approx$ 20 legal moves
   ▪ Depth 2 $\rightarrow$ $20 \times (21 \times 20)^3 = 1.2 \times 10^9$

o **As depth increases, probability of reaching a given search node shrinks**

   ▪ So usefulness of search is diminished
   ▪ So limiting depth is less damaging
   ▪ But pruning is trickier…

o **Historic AI: TDGammon uses depth-2 search + very good evaluation fuanction + reinforcement learning $\rightarrow$ world-champion level play**

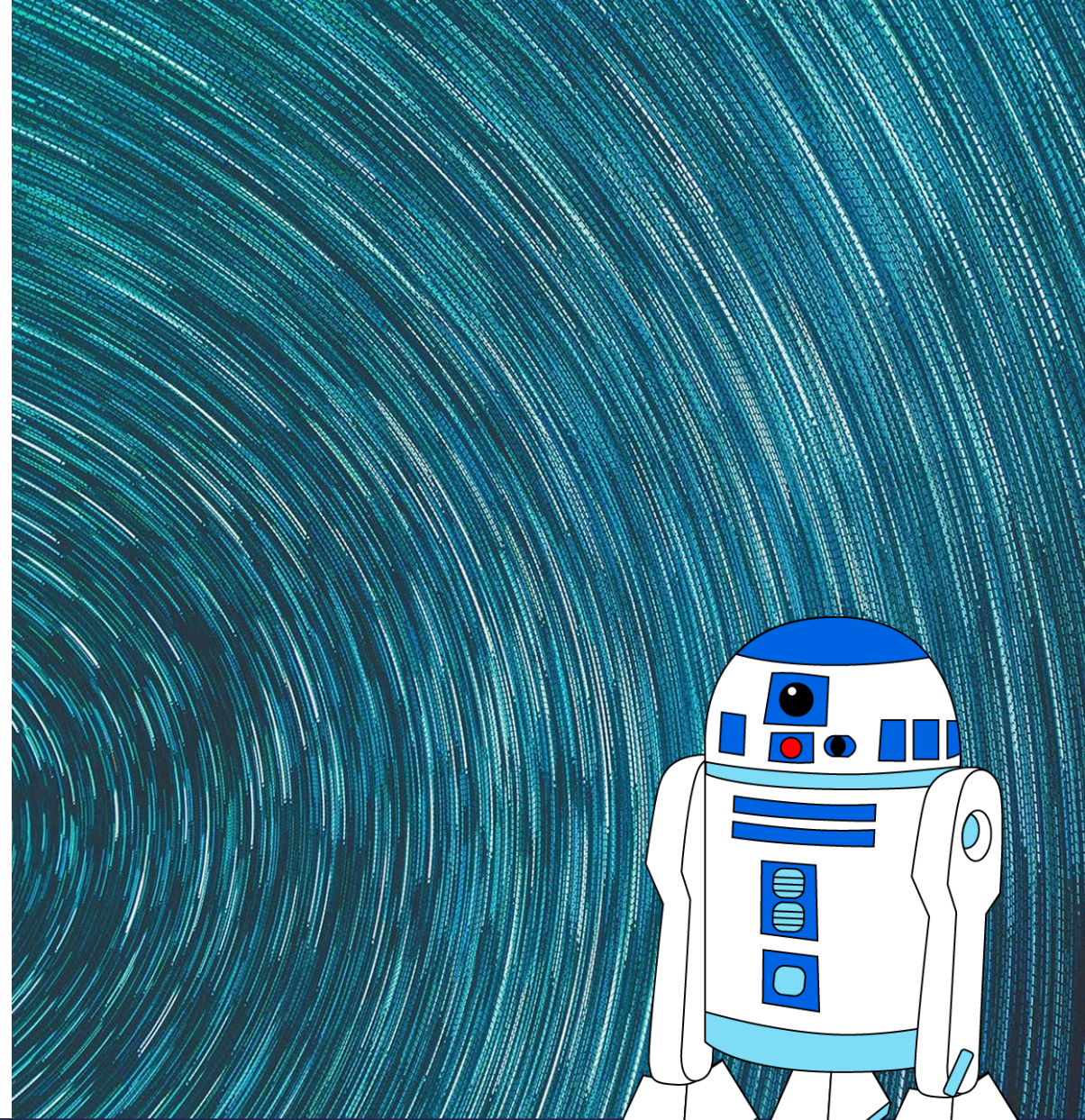o **1$^{st}$ AI world champion in any game!**

- o **E.g. Backgammon**
- o **Expectiminimax**
  - ▪ Environment is an extra "random agent" player that moves after each min/max agent
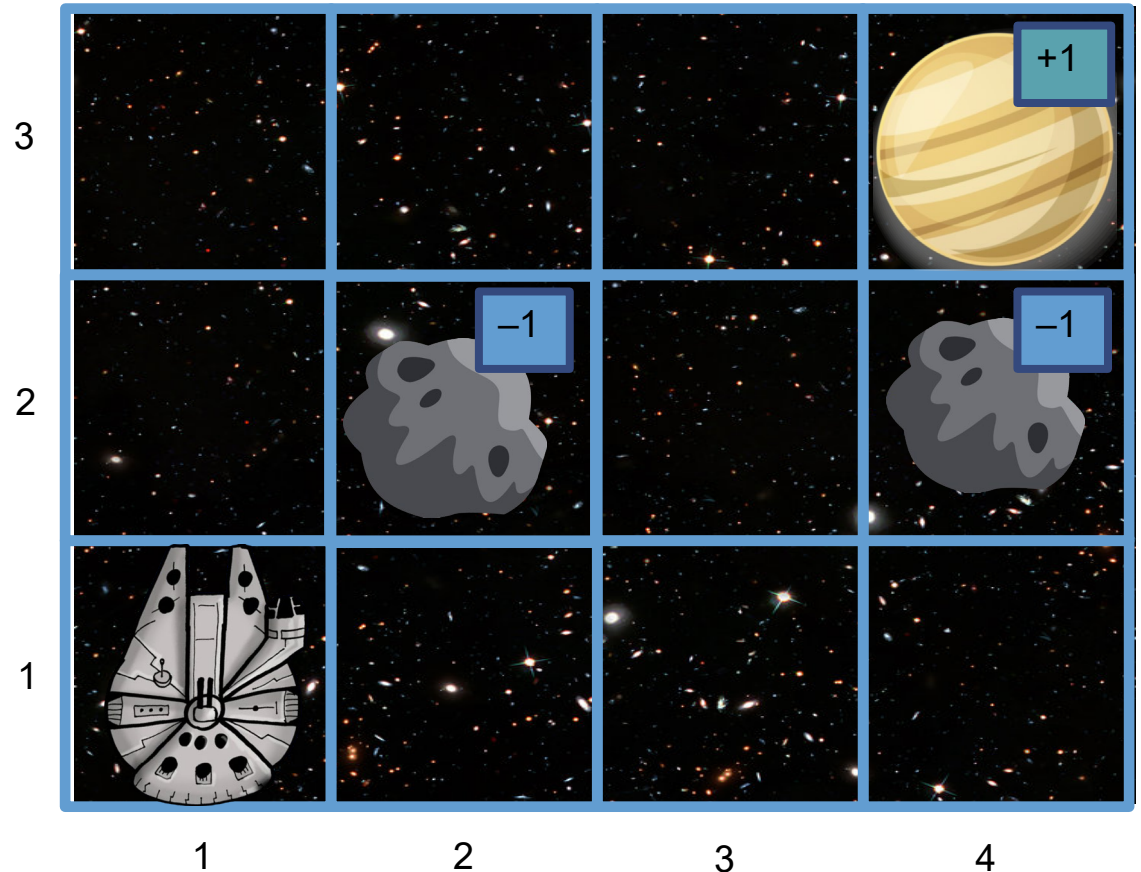  - ▪ Each node computes the appropriate combination of its children



Max nodes

Chance nodes

Min nodes

CIS 421/521:
ARTIFICIAL INTELLIGENCE

# Markov Decision Processes

Penn Engineering
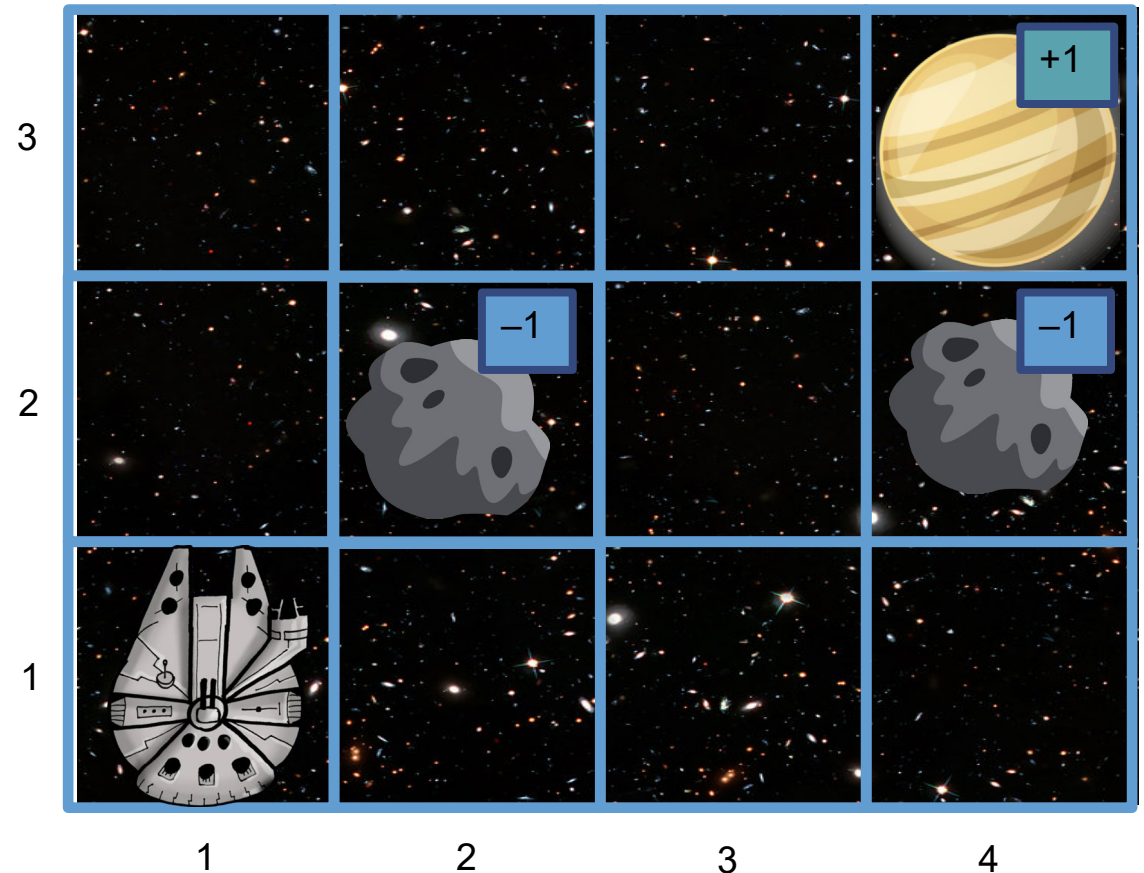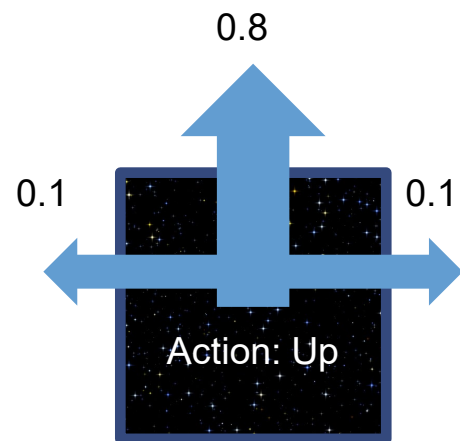UNIVERSITY of PENNSYLVANIA

# Navigating an Asteroid Field

○ Suppose we have a **fully-observable** 4x3 environment with goal states.

○ The millennium falcon begins in the start state and **picks an action at each time step**.

○ Actions: *Up, Down, Left, Right*

○ The game **terminates when it reaches a goal state** (+1 or -1).

○ If the environment were **deterministic**, the solution would be easy:

  ○ [*Up, Up, Right, Right, Right*]

# Navigating an Asteroid Field

- Instead of making the environment deterministic, we will make it **stochastic**.

- If the Falcon selects the action *Up* then it only moves up 80% of the time.

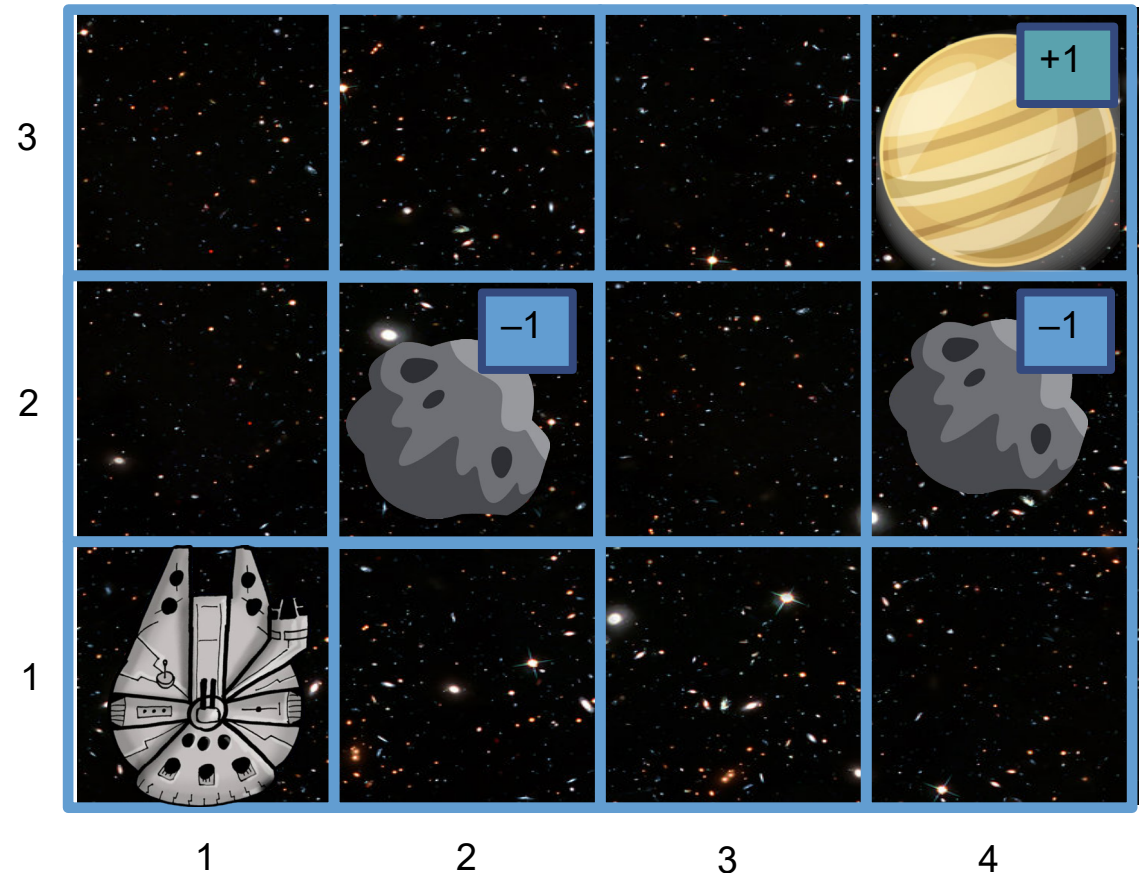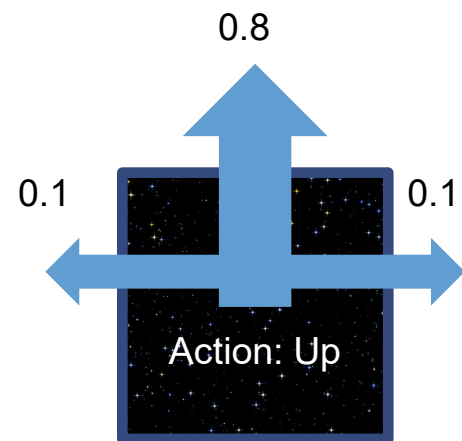- 10% of the time the weird gravity fields cause it to veer off to the left or right.

Transition Model:



0.8

0.1        0.1

Action: Up

# Navigating an Asteroid Field

o For action sequence

   o [*Up, Up, Right, Right, Right*],

o what's the probability that the millennium falcon reaches the intended goal?

Transition Model:
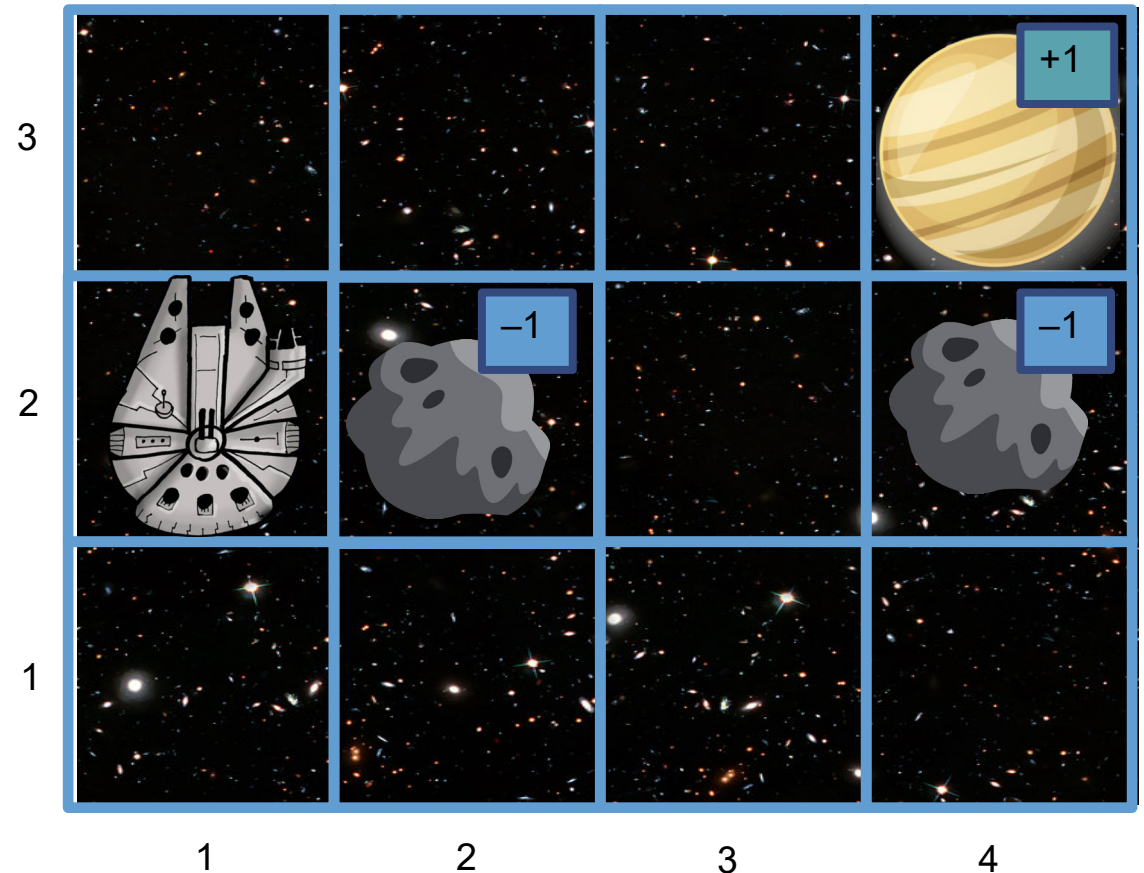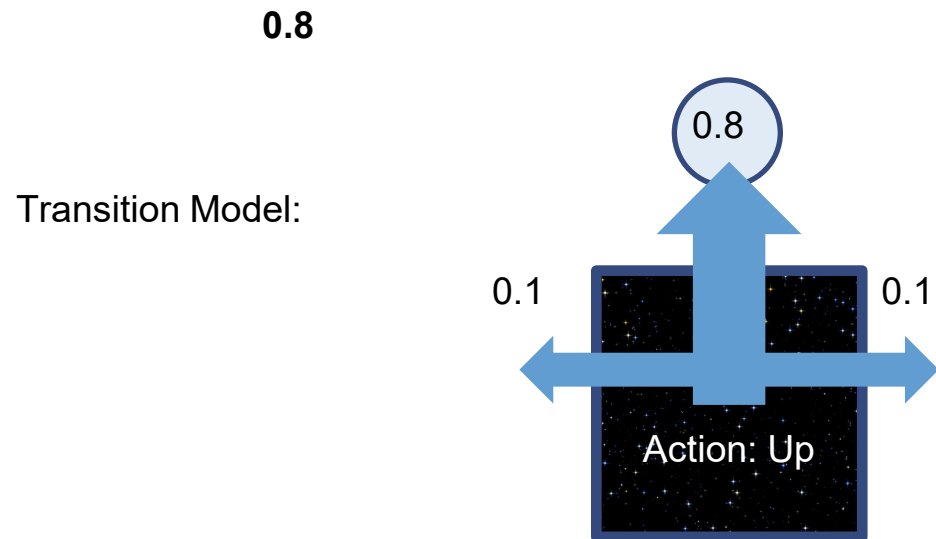
0.8

0.1    0.1

Action: Up

# Navigating an Asteroid Field

- For action sequence
  - [*Up, Up, Right, Right, Right*],
- what's the probability that the millennium falcon reaches the intended goal?

**0.8**

Transition Model:

# Navigating an Asteroid Field

○ For action sequence

    ○ [*Up, Up, Right, Right, Right*],

○ what's the probability that the millennium falcon reaches the intended goal?
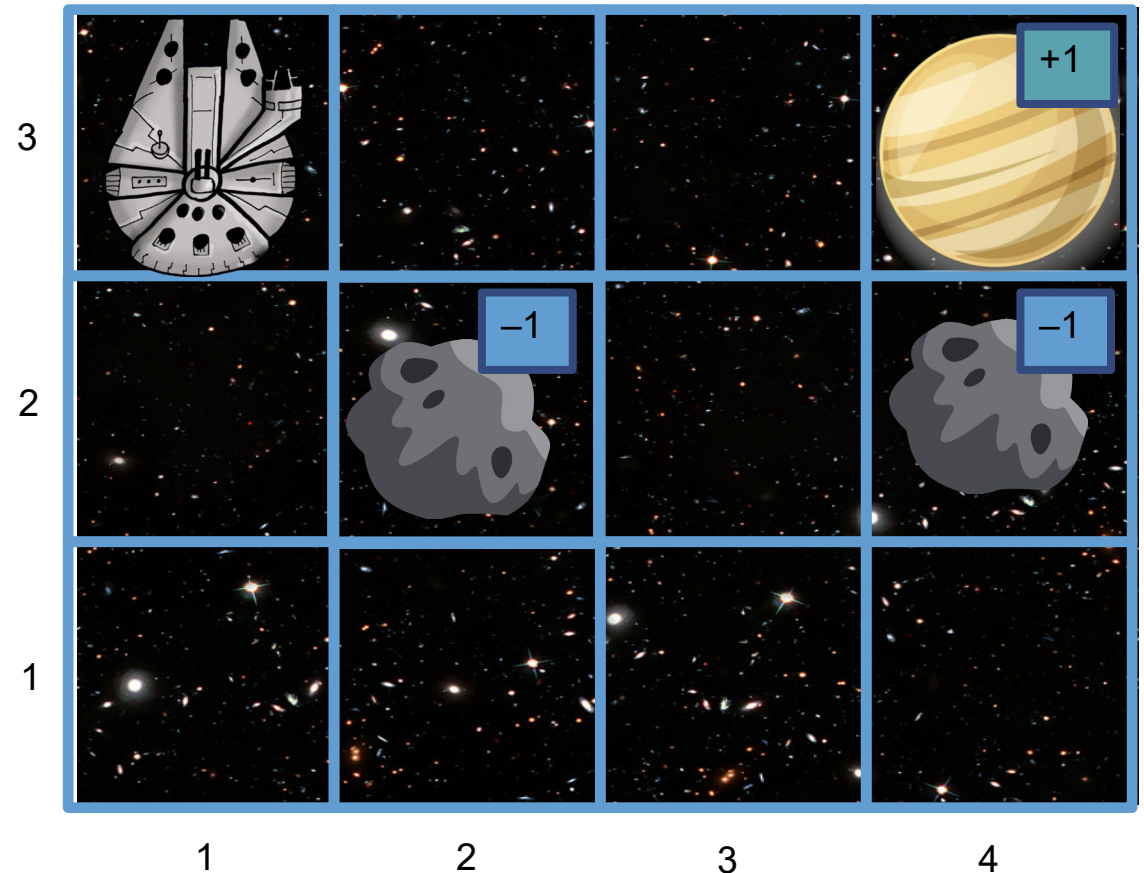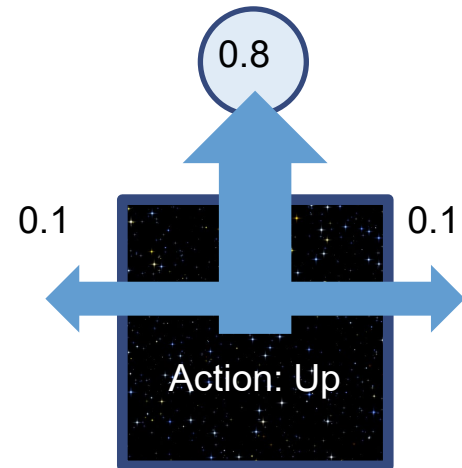
**0.8  * 0.8**

Transition Model:

# Navigating an Asteroid Field

○ For action sequence

　○ [*Up, Up, Right, Right, Right*],

○ what's the probability that the millennium falcon reaches the intended goal?

**0.8 * 0.8 * 0.8**

Transition Model:

0.1

Action:
Right

0.8

0.1

# Navigating an Asteroid Field

○ For action sequence

    ○ [*Up, Up, Right, Right, Right*],

○ what's the probability that the millennium falcon reaches the intended goal?

**0.8** * **0.8** * **0.8** * **0.8**
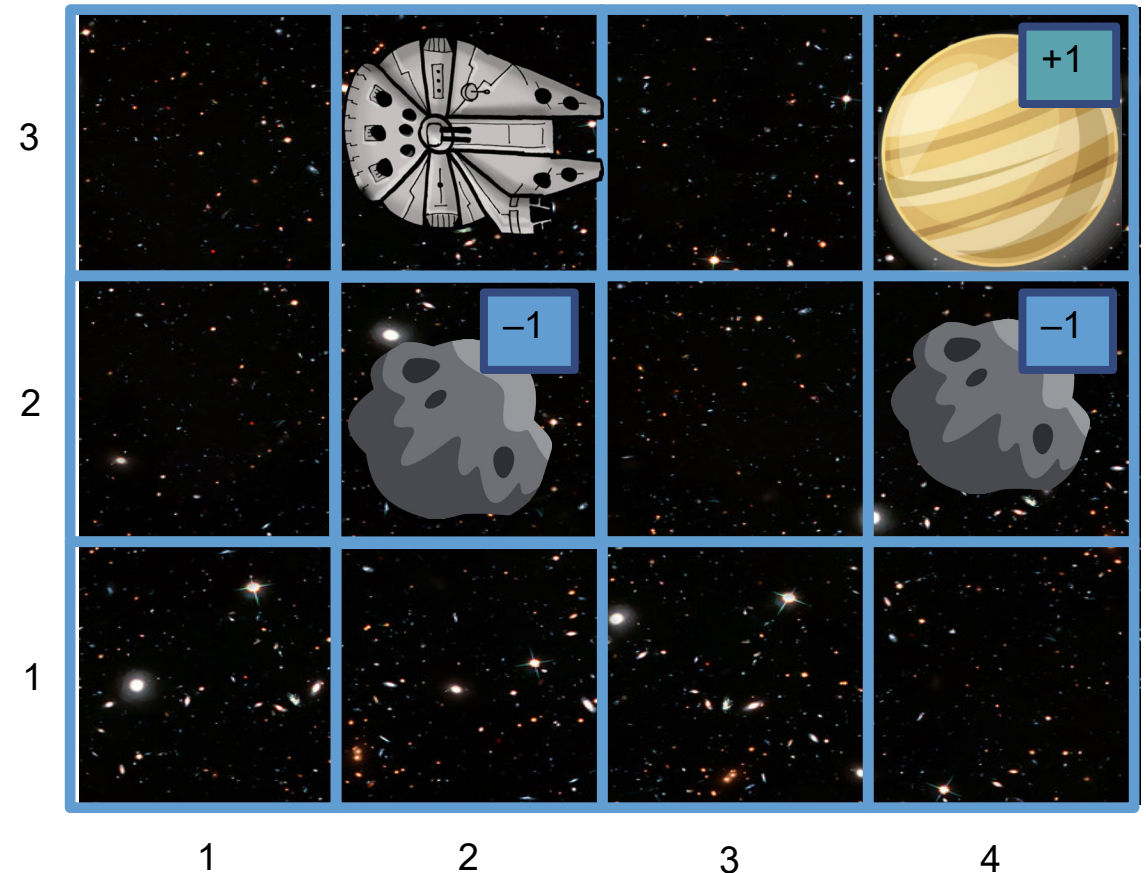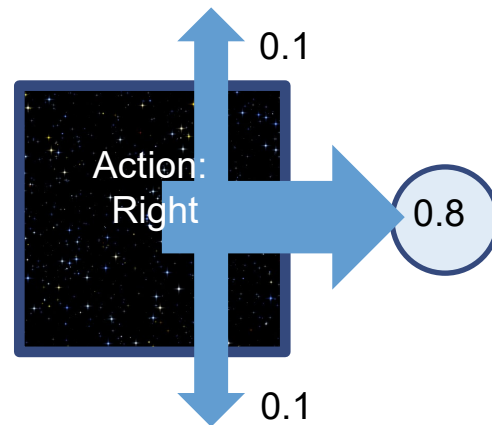
Transition Model:

# Navigating an Asteroid Field

○ For action sequence

  ○ [*Up, Up, Right, Right, Right*],

○ what's the probability that the millennium falcon reaches the intended goal?

**0.8 * 0.8 * 0.8 * 0.8 * 0.8**
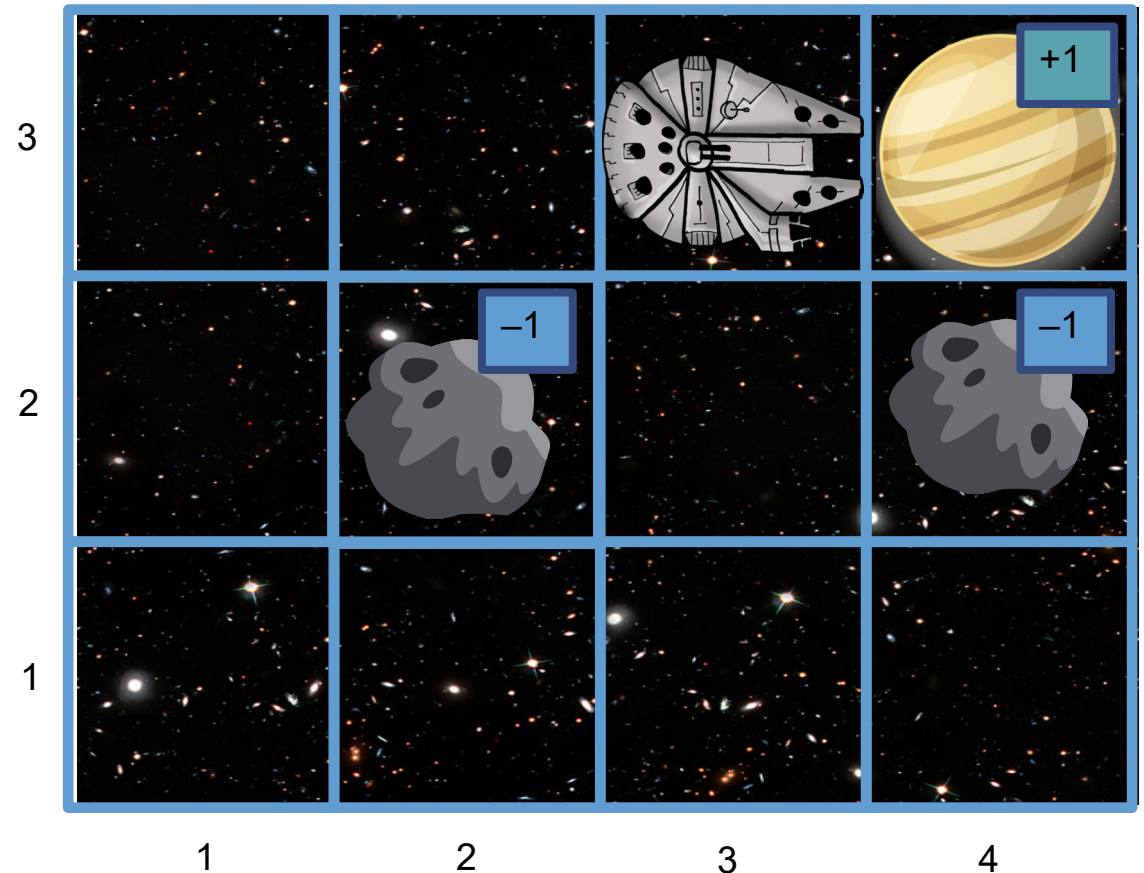
**= 0.32768**

Transition Model:

# Navigating an Asteroid Field

○ For action sequence

    ○ [*Up, Up, Right, Right, Right*],

○ what's the probability that the millennium falcon reaches the intended goal?

Transition Model:

0.8

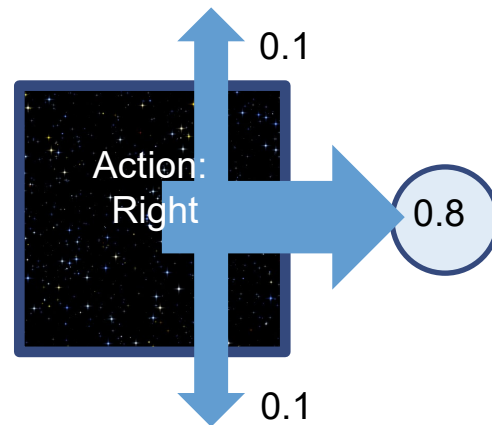0.1              0.1

Action: Up

# Navigating an Asteroid Field

- For action sequence
  - [*Up, Up, Right, Right, Right*],
- what's the probability that the millennium falcon reaches the intended goal?

**0.1**

Transition Model:

0.8
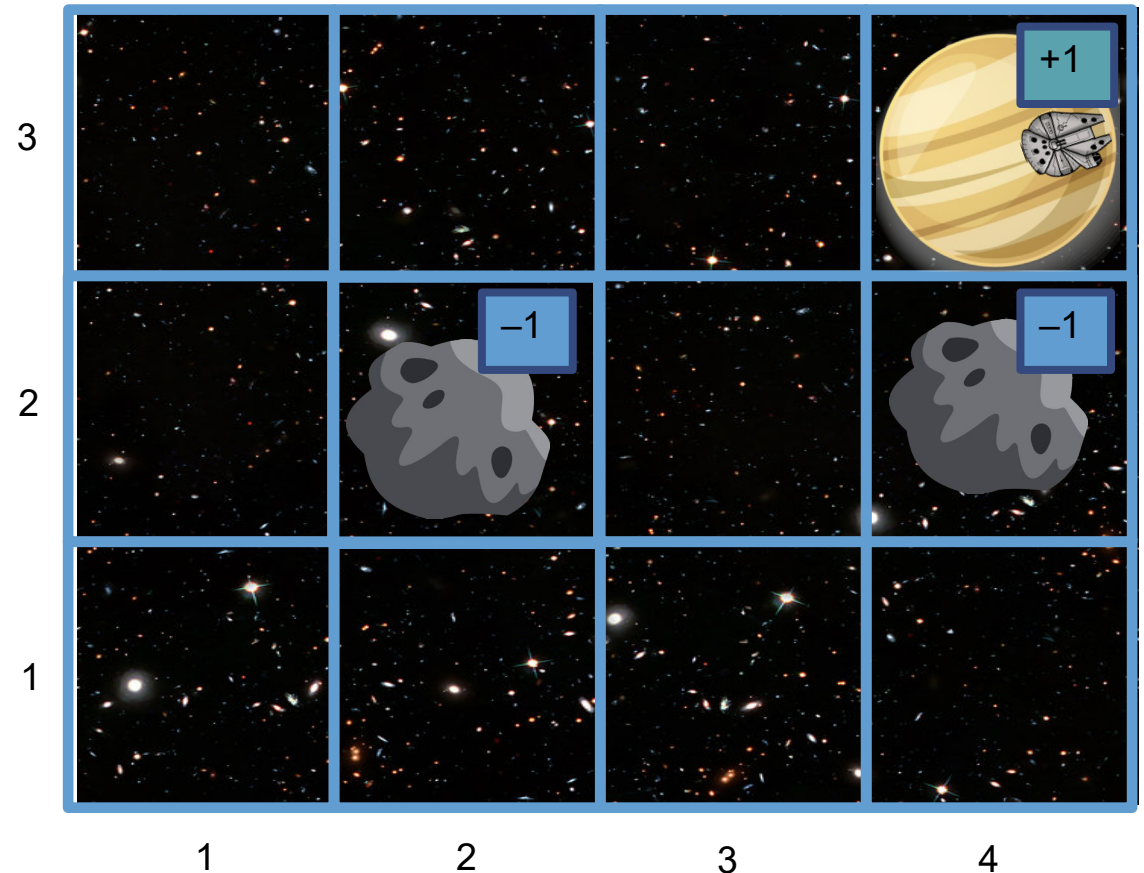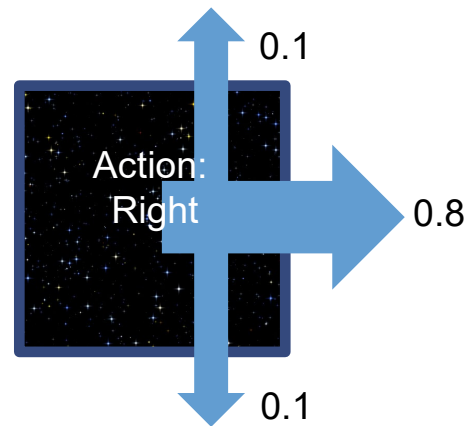
0.1          0.1

Action: Up

# Navigating an Asteroid Field

○ For action sequence

    ○ [*Up, Up, Right, Right, Right*],

○ what's the probability that the millennium falcon reaches the intended goal?
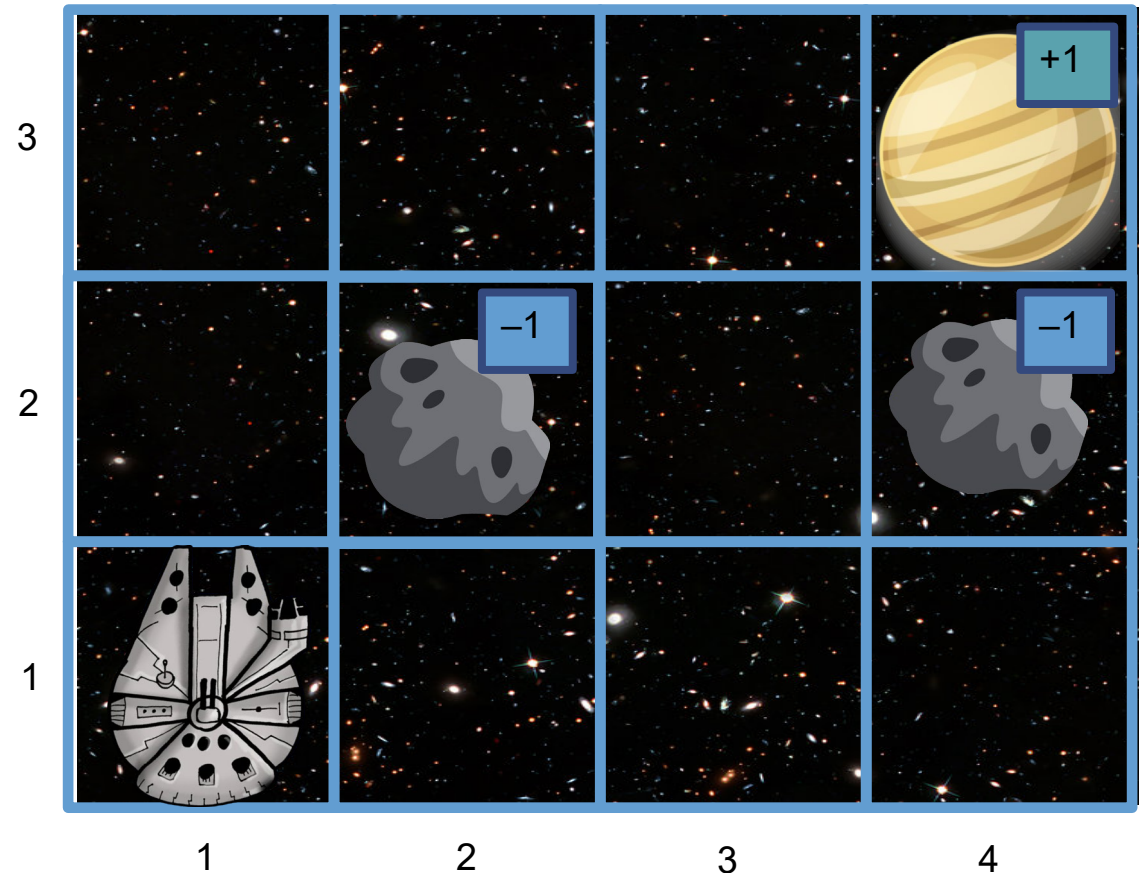
**0.1 * 0.1**

Transition Model:

# Navigating an Asteroid Field

- For action sequence
  - [*Up, Up, Right, Right, Right*],
- what's the probability that the millennium falcon reaches the intended goal?

**0.1** * **0.1** * **0.1**

Transition Model:



Action: Right — 0.1 (up), 0.8 (right), 0.1 (down)

# Navigating an Asteroid Field

○ For action sequence

    ○ [*Up, Up, Right, Right, Right*],

○ what's the probability that the millennium falcon reaches the intended goal?
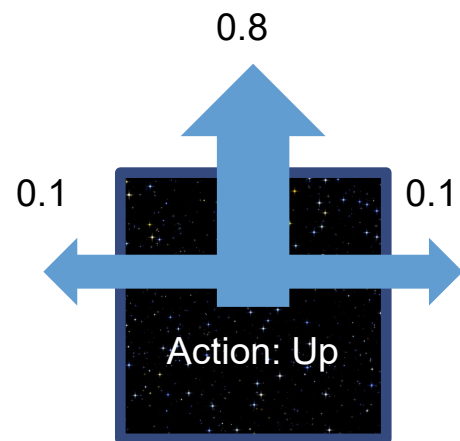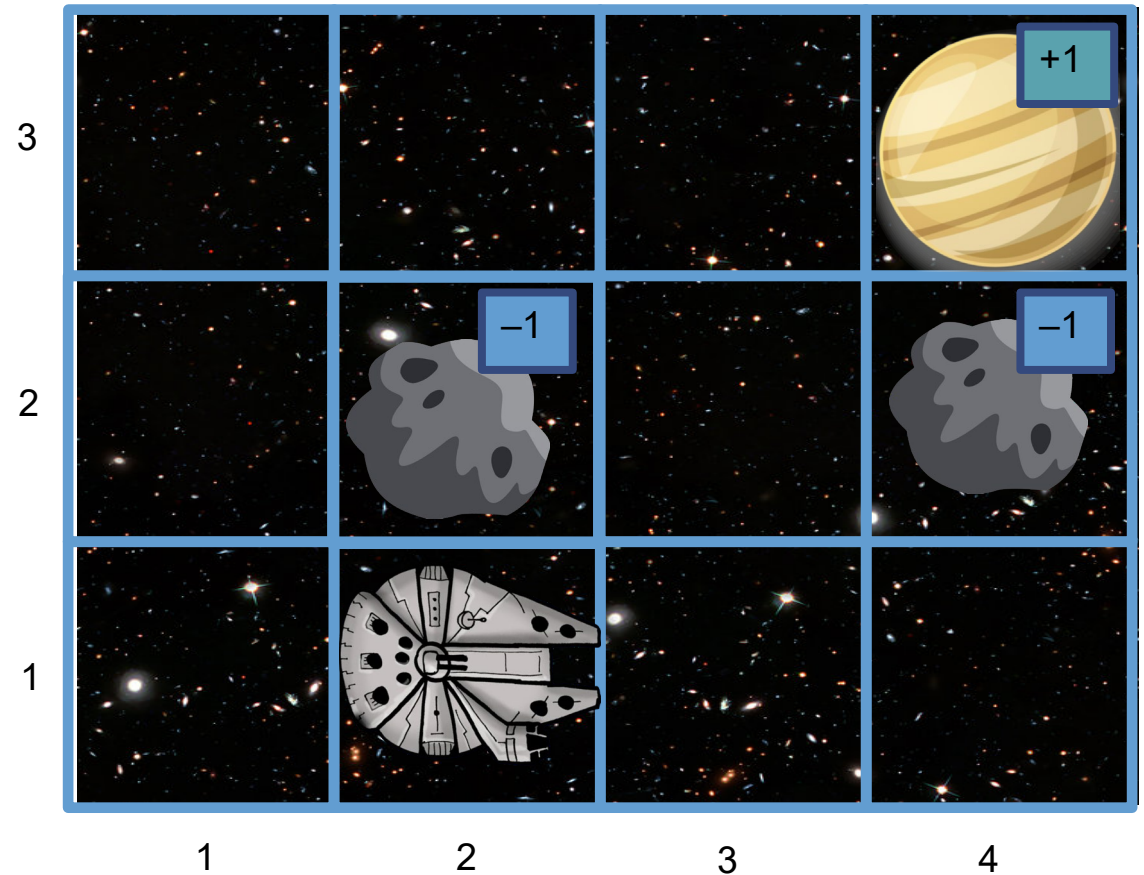
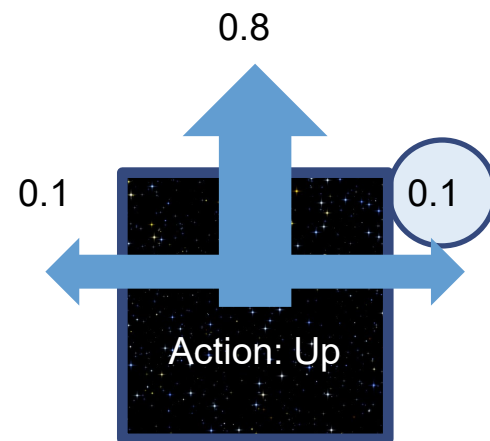**0.1 * 0.1 * 0.1 * 0.1**

Transition Model:

# Navigating an Asteroid Field

○ For action sequence

  ○ [*Up, Up, Right, Right, Right*],

○ what's the probability that the millennium falcon reaches the intended goal?

**0.1 * 0.1 * 0.1 * 0.1 * 0.8**
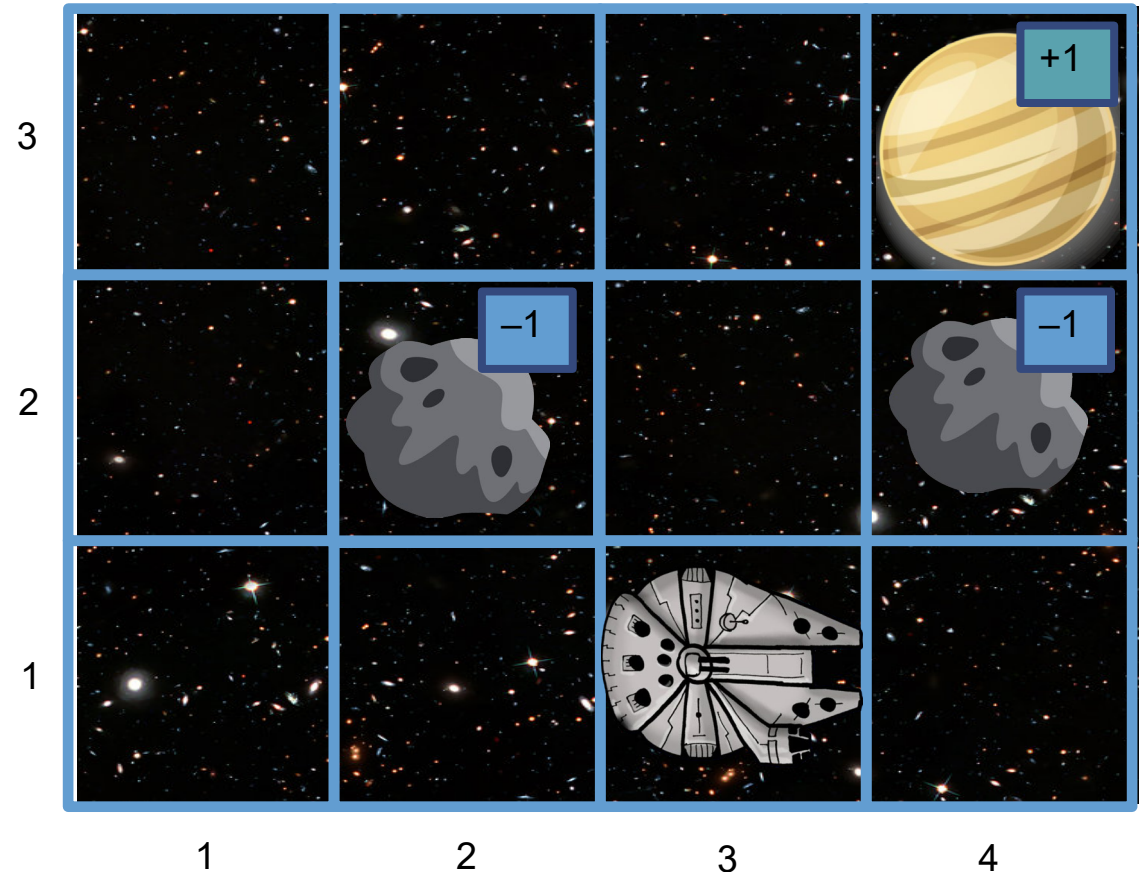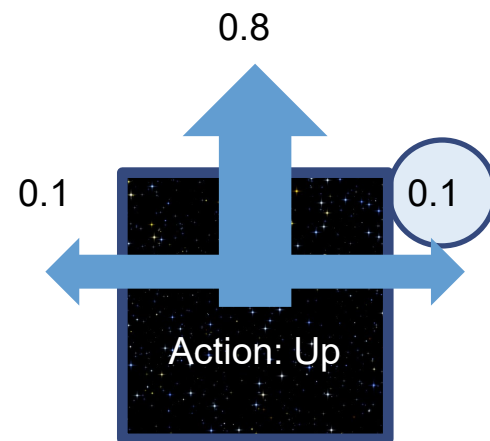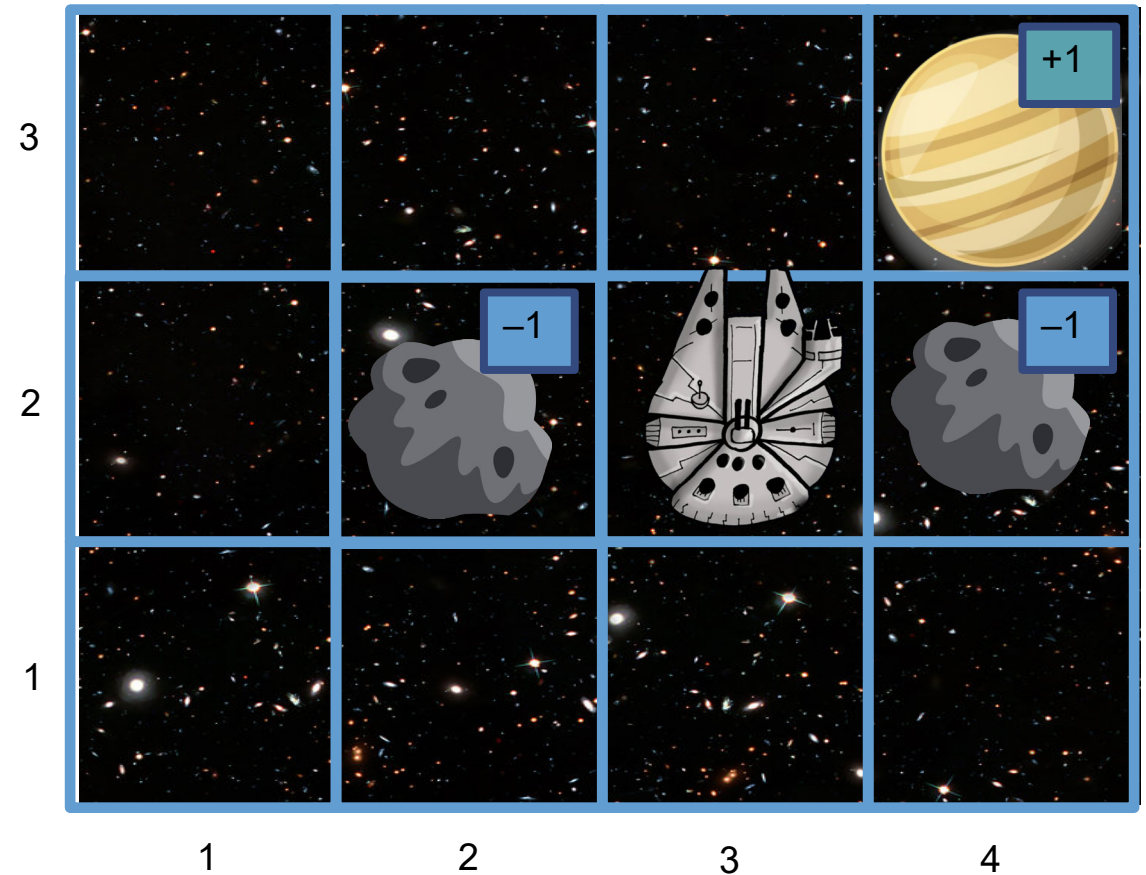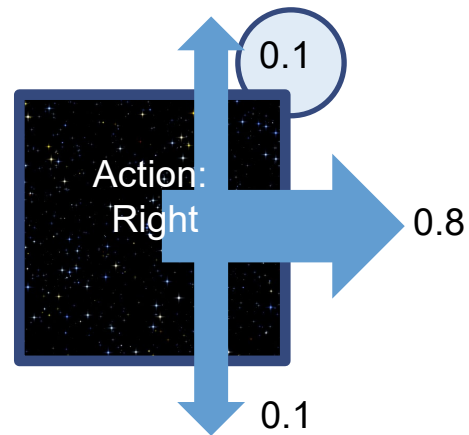
**= 0.00008**

Transition Model:

# Navigating an Asteroid Field

○ For action sequence

    ○ [*Up, Up, Right, Right, Right*],

○ what's the probability that the millennium falcon reaches the intended goal?

**0.32768 + 0.00008**

**= 0.32776**

Transition Model:

The odds of successfully navigating an asteroid field

Action: Right

0.1

0.8

0.1

+1

−1

−1

3

1

1    2    3    4

# Stochastic Transition Model

○ In our search algorithms so far, the transition model was deterministic and described the outcome of each action in each state.

○ The transition function is sometimes written as T(s, a, s'), or explicitly as a probability:

The probability of        arriving in state **s'**

p( s' | s, a )

given that        we are in state **s** and we selected action **a**

0.8

0.1        0.1

Action: Up

# Stochastic Transition Model

○ In our search algorithms so far, the transition model was deterministic and described the outcome of each action in each state.

○ The transition function is sometimes written as T(s, a, s'), or explicitly as a probability:

Andrey Markov (1856-1922)

p( s' | s, a )

Transitions are **Markovian**: the probability of arriving in s' only depends on s and not the history of earlier states.

0.8

0.1

0.1

Action: Up

# Reward function

o   We will specify a **utility or reward function** for the agent.

o   The "rewards" can be **positive** or **negative** but are bounded by some maximum value.

o   Because the decision process is **sequential,** we must specify the utility function on a sequence of states and actions.

o   Instead of only giving a reward at the goal states, the agent can **receive a reward at each time step**, based on its transition from **s** to **s'** via action **a**.

o   This is defined by a reward function

o      R( s, a, s' )

o   *For example, we could give the Millennium Falcon a small negative reward of -0.04 for every transition except for entering the terminal states (+1 for entering the planet's orbit or -1 for smashing into an asteroid).*

o   The **rewards are additive**, so if the Millennium Falcon takes 4 steps before entering the planet's orbit, it gets *-0.04 + -0.04 + -0.04 + -0.04 + 1 = 0.84* for that solution.

# Markov Decision Pr...

Expectimax node: outcome is uncertain. In expectimax search we calculate their expected utilities.

○ A **Markov decision process** or **MDP** is

- a **sequential** decision problem
- for a **fully observable** environment
- with a **stochastic** transition model
- that has **additive rewards**

○ **MDPs are non-deterministic search problems.** One way of solving them is via **expectimax** search.

3  12  9  2  4  6  15  6  0

# Markov Decision Process

o   To find a solution to an MDP, you need to define the following things:

- **A set of states** $s \in S$

- **A set of actions** $a \in A$

- A transition function **T(s, a, s')**
  - Probability that executing action **a** in **s** will lead to **s' P(s' | s, a)**
  - The probability is called **the model**

- A reward function **R(s, a, s')**
  - Sometimes just R(s) or R(s')

- An **initial state** $s_0$

- Optionally, one or more **terminal states**

# Solution == Policy

○ In search problems a solution was a sequence of action that corresponded to the shortest path.

○ Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

○ Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.

    ○ $\pi(s) = a$

Policy $\pi$ tells the agent what action to take at state s.

This is an example policy for a grid world

# Solution == Policy

○ In search problems a solution was **a plan**: a sequence of action that corresponded to the shortest path from the start to a goal.

○ Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

○ Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.

    ○ **π**(s) = a

# Solution == Policy

○ In search problems a solution was **a plan**: a sequence of action that corresponded to the shortest path from the start to a goal.

○ Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

○ Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.

    ○ $\pi(s) = a$

# Solution == Policy

○ In search problems a solution was **a plan**: a sequence of action that corresponded to the shortest path from the start to a goal.

○ Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

○ Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.
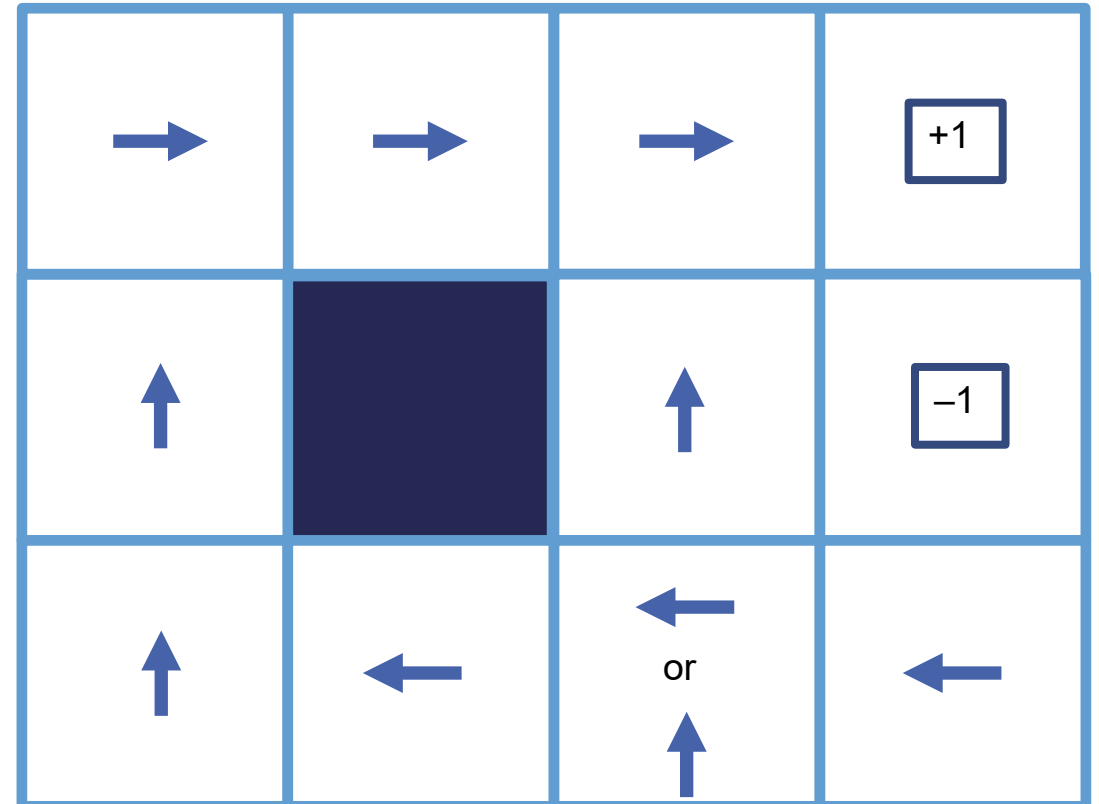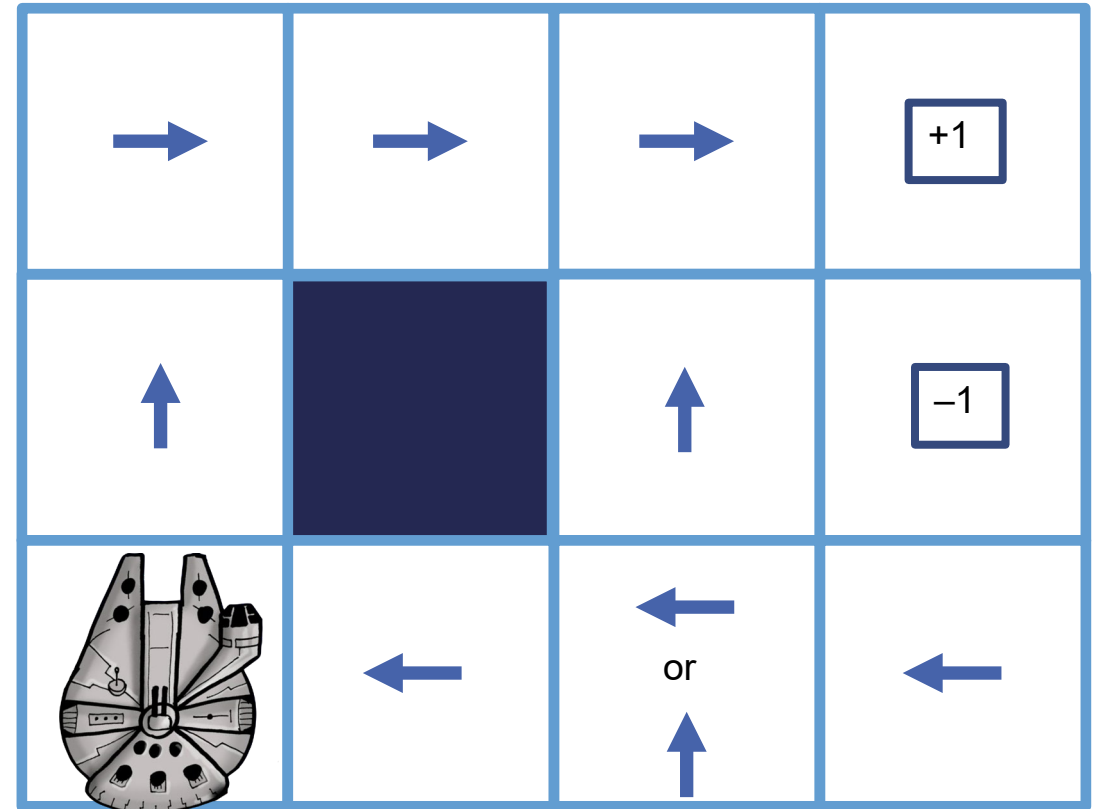
   ○ $\pi(s) = a$

# Solution == Policy

○ In search problems a solution was **a plan**: a sequence of action that corresponded to the shortest path from the start to a goal.

○ Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

○ Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.

    ○ **π**(s) = a

# Solution == Policy

- In search problems a solution was **a plan**: a sequence of action that corresponded to the shortest path from the start to a goal.

- Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

- Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.
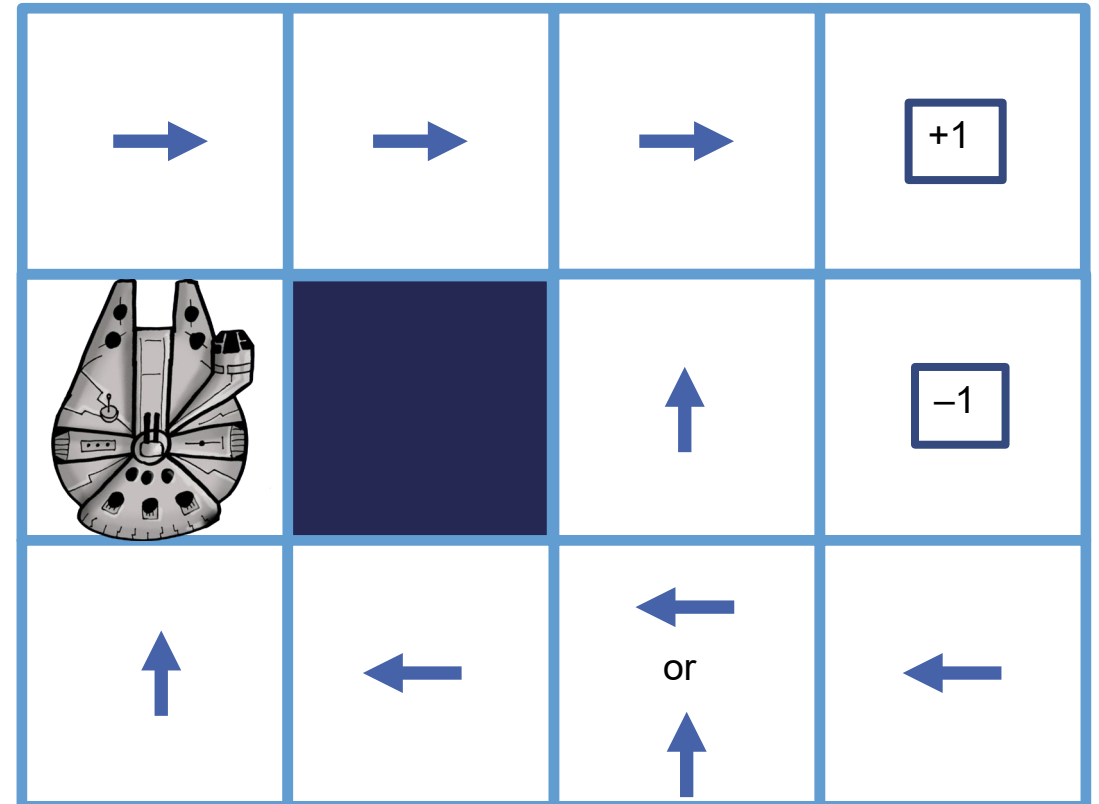
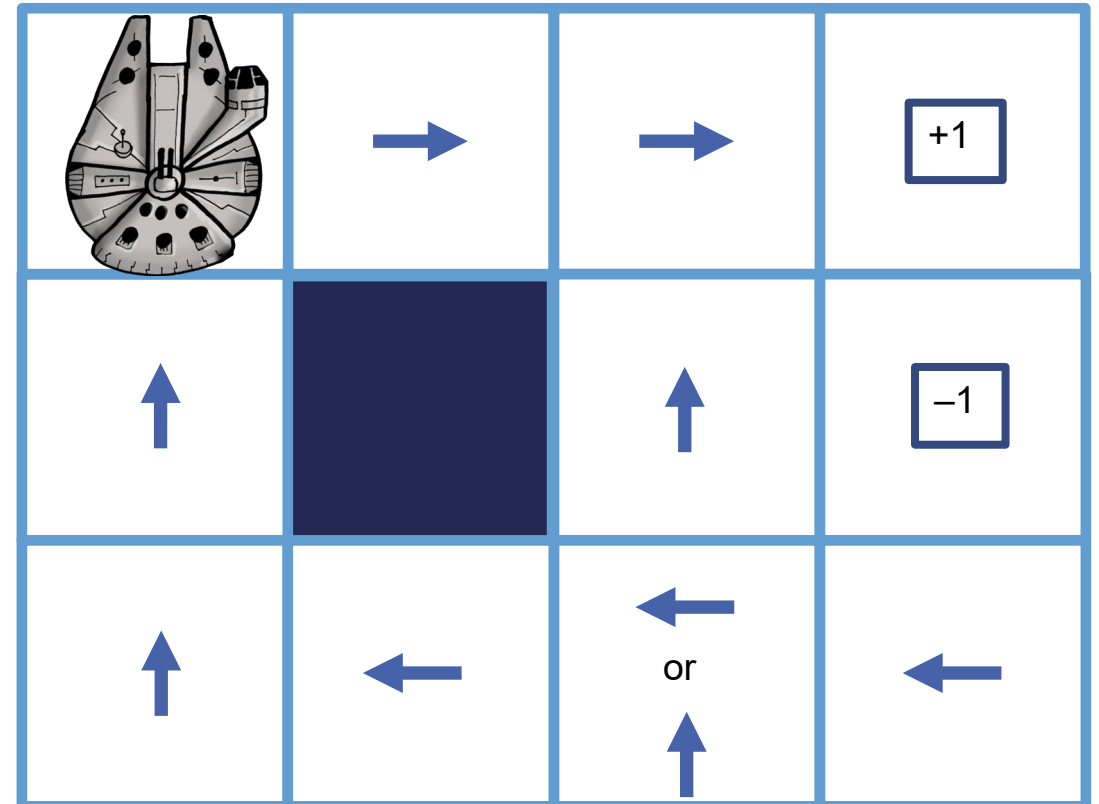  - $\pi(s) = a$

# Solution == Policy

- In search problems a solution was **a plan**: a sequence of action that corresponded to the shortest path from the start to a goal.

- Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

- Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.
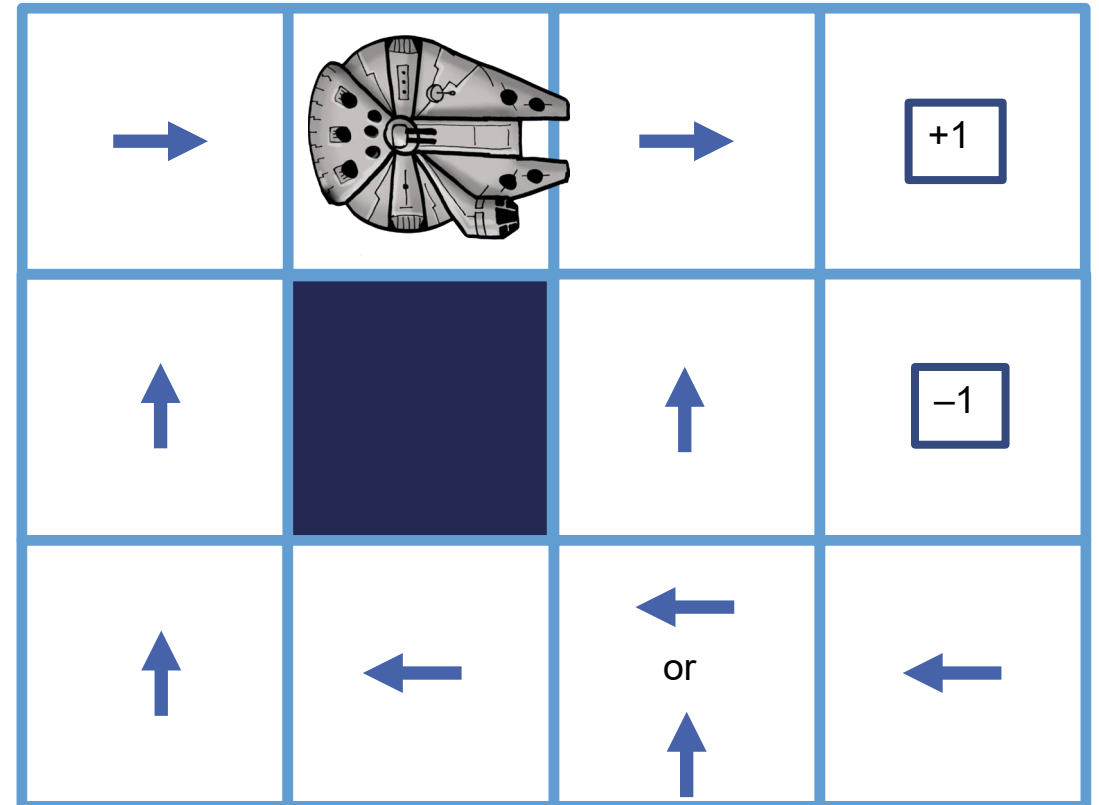
  - $\pi(s) = a$

# Solution == Policy

Even though the policy told me to go right here, there's no guarantee that me picking the action Right will result in me moving right. It's stochastic!

- In search problems a solution was **a plan**: a sequence of action that corresponded to the shortest path from the start to a goal.

- Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

- Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.
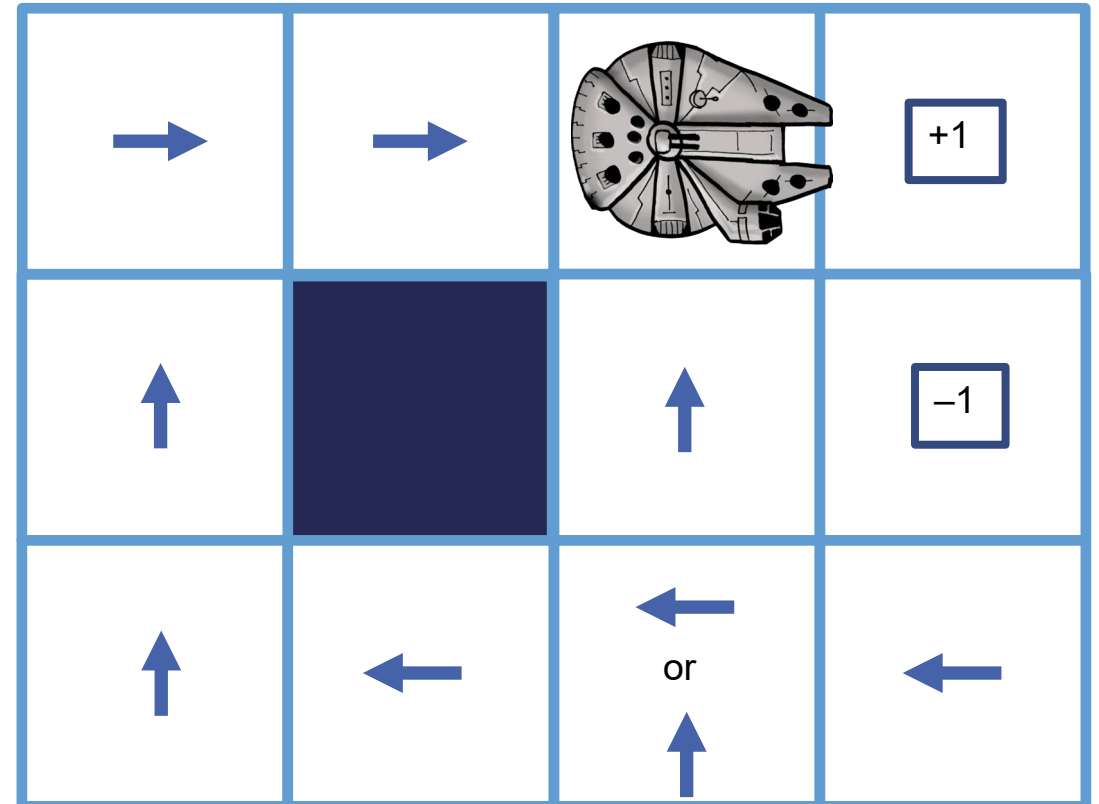
  - $\pi(s) = a$

# Solution == Policy

○ In search problems a solution was **a plan**: a sequence of action that corresponded to the shortest path from the start to a goal.

○ Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

○ Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.
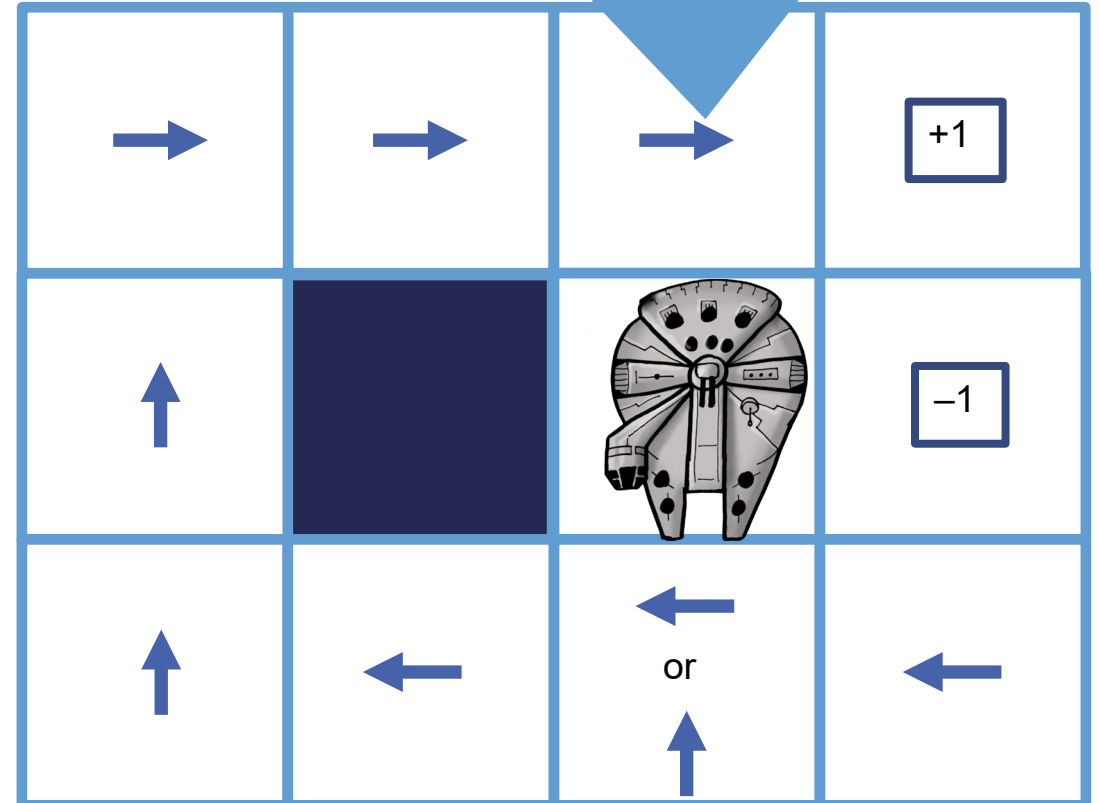
    ○ $\pi(s) = a$

# Solution == Policy

o   In search problems a solution was **a plan**: a sequence of action that corresponded to the shortest path from the start to a goal.

o   Because of the non-determinism in MDPs we cannot simply give a sequence of actions.

o   Instead, the solution to an MDP is a **policy.** A policy maps from a state onto the action to take if the agent is in that state.
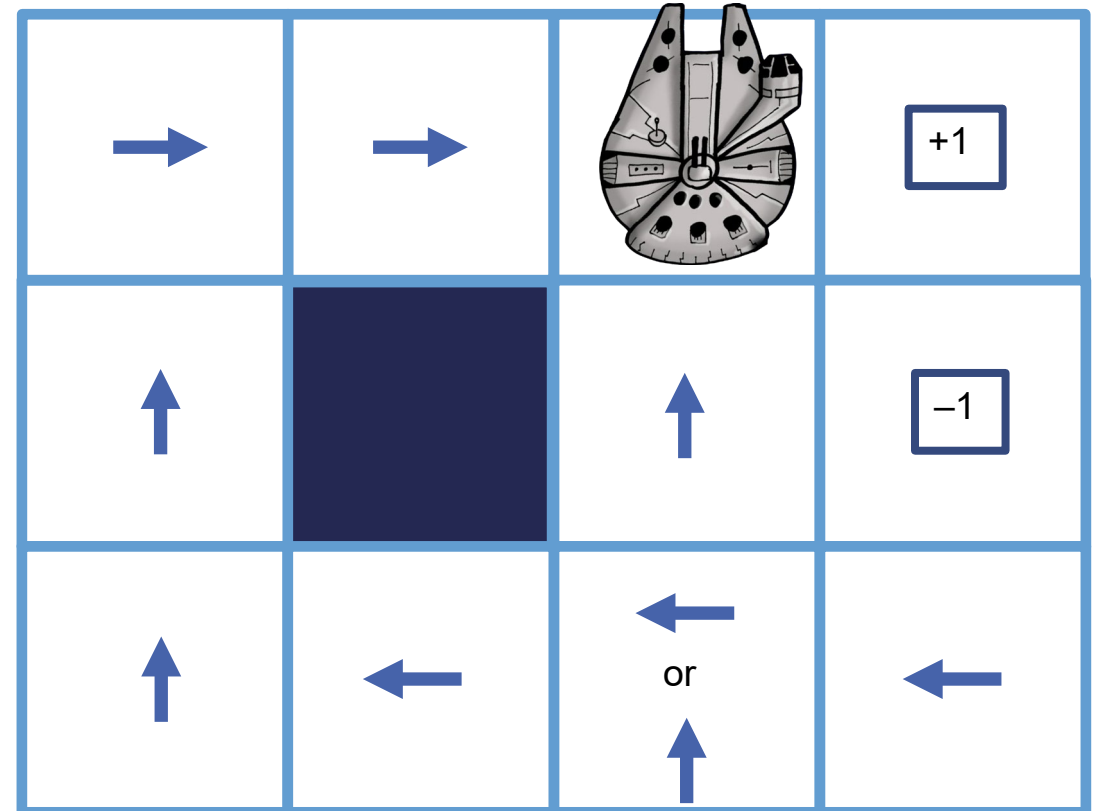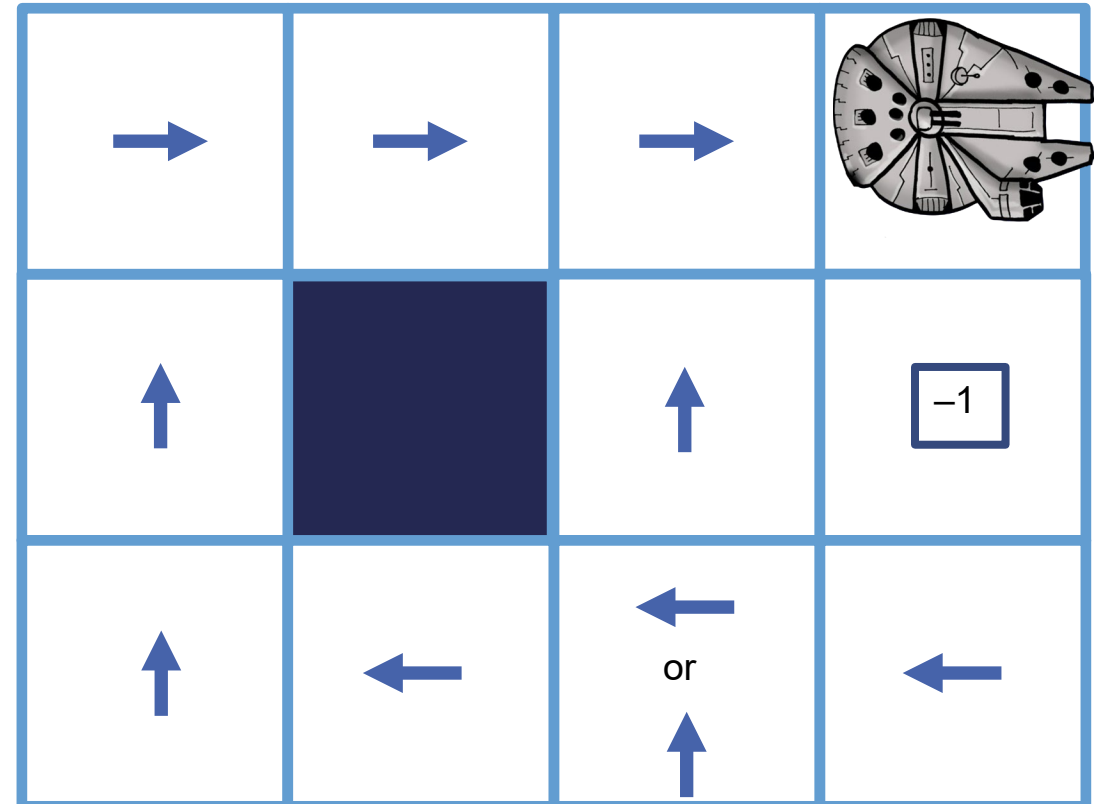
o   **π**(s) = a

We will use **π\*** to denote the **optimal policy**.

# Policies and Rewards



r = -0.04

o Even if the **same policy** is executed multiple times by the agent, this may lead to different sequence of states and actions (**environment history**), and thus a **different score** under the reward function.

o Therefore we need to compute the **expected utility** of all the possible paths generated by a policy.

r is the reward per action

r < -1.65



-0.03 < r < 0



r > 0



Penn Engineering

# Sequences of Rewards

○ The performance of an agent in an MDP is the sum of the rewards for the transitions it takes.

○ $U_h([s_0, a_0, s_1, a_1 ..., s_n])$

$= R(s_0, a_0, s_1) + R(s_1, a_1, s_2) + ... + R(s_{n-1}, a_{n-1}, s_n)$

Utility function on an environment history.

Sequence of states and actions

Bounce around forever, and avoid the exits ... **infinite rewards!!**

**r > 0**

# Utilities of Sequences

Penn Engineering

# 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]

# Discounting

o It's reasonable to maximize the sum of rewards

o It's also reasonable to prefer rewards now to rewards later

o One solution: values of rewards decay exponentially

$1$

$\gamma$

$\gamma^2$

Worth Now

Worth Next Step

Worth In Two Steps

# Discounting

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

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

o Example: discount of 0.5
  ▪ $U([1,2,3]) = 0.5^0 * 1 + 0.5^1 * 2 + 0.5^2 * 3$
    $= 1*1 + 0.5*2 + 0.25*3$
  ▪ $U([1,2,3]) < U([3,2,1])$



$1$

$\gamma$

$\gamma^2$

# Stationary Preferences

o Theorem: if we assume stationary preferences:

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

$$\updownarrow$$

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

o 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:
$$U([r_0, r_1, r_2, \ldots]) = r_0 + \gamma r_1 + \gamma^2 r_2 \cdots$$

# 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 $0 < \gamma < 1$

$$U([r_0, \ldots r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{\max}/(1 - \gamma)$$

  - 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

o 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 $\gamma$)

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

s

a

s, a

s,a,s'

s'

# Example Hyperdrive MDP

The Millennium Falcon needs to travel far far away, quickly
Three states: Cruising, Hyperspace, Crashed
Two actions: *Maintain speed*, *Punch it*
Going faster gets double reward



0.5    +1

*Punch It*

1.0

+1    *Maintain*    0.5    -10

**Hyperspace**

*Maintain*    0.5

*Punch It*    0.5    +2

1.0

+1    0.5    +2

**Cruising**

**Crashed**

# MDP Search Trees

○ Each MDP state projects an expectimax-like search tree



s is a state

(s, a) is a q-state

(s,a,s') is called a transition
T(s,a,s') = P(s'|s,a)
R(s,a,s')

s

a

s, a

s,a,s'

s'

# Hyperdrive Search Tree

# Optimal Quantities

- **The value (utility) of a state s:**

  $V^*(s)$ = expected utility starting in s and acting optimally

- **The value (utility) of a q-state (s,a):**

  $Q^*(s,a)$ = expected utility starting out having taken action a from state s and (thereafter) acting optimally

- **The optimal policy:**

  $\pi^*(s)$ = optimal action from state s

s is a *state*

(s, a) is a *q-state*

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

s

a

s, a

s,a,s'

s'

# Values of States

- Fundamental operation: compute the (expectimax) value of a state
  - Expected utility under optimal action
  - Average sum of (discounted) rewards
  - This is just what expectimax computed!

- Recursive definition of value:

$$V^*(s) = \max_a Q^*(s,a)$$

$$Q^*(s,a) = \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V^*(s') \right]$$

$$V^*(s) = \max_a \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V^*(s') \right]$$



s

a

s, a

s,a,s'

s'

# Racing Search Tree

# Racing Search Tree

# Racing Search Tree

o  We're doing way too much work with expectimax!

o  Problem: States are repeated
   ▪  Idea: Only compute needed quantities once

o  Problem: Tree goes on forever
   ▪  Idea: Do a depth-limited computation, but with increasing depths until change is small
   ▪  Note: deep parts of the tree eventually don't matter if $\gamma <$ 1

# Time-Limited Values

○ Key idea: time-limited values

○ Define $V_k(s)$ to be the optimal value of s if the game ends in k more time steps
  ▪ Equivalently, it's what a depth-k expectimax would give from s



$V_2(\text{🚗})$

# Reminders

- 21 days until the American election.  I voted.  Did you?

- Deadline to register to vote in PA is  **Monday, Oct 19.**

- HW4 due tonight at 11:59pm Eastern.

- Quiz 5 on Adversarial Search is due tomorrow.

- HW5 has been released.  It will be due on Tuesday Oct 20.

- No lecture on Thursday.

- Midterm details:

- * No HW from Oct 20-27.
  * Tues Oct 20: Practice midterm released (for credit)
  * Saturday Oct 24: Practice midterm is due.
  * Midterm available Monday Oct 26 and Tuesday Oct 27.
  * 3 hour block.  Open book, open notes, no collaboration.

# Computing Time-Limited Values

# Value Iteration

# Value Iteration

- Start with $V_0(s) = 0$: no time steps left means an expected reward sum of zero

- Given vector of $V_k(s)$ values, do one ply of expectimax from each state:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V_k(s') \right]$$

- Repeat until convergence

- Complexity of each iteration: $O(S^2 A)$

- Theorem: will converge to unique optimal values
  - Basic idea: approximations get refined towards optimal values
  - Policy may converge long before values do

$V_{k+1}(s)$

a

s, a

s,a,s'

$V_k(s')$

# Example: Value Iteration



$V_2$

$V_1$

$V_0$

Slow 0.5 +1

Slow +1 0.5

Fast 1.0 -10

Warm

Slow

Cool

Fast 0.5 +2

0.5

+2

1.0 +1

Overheated

*Assume no discount!*

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$

# Example: Value Iteration



$V_2$     3.5     2.5     0

$V_1$     2     1     0

$V_0$     0     0     0

0.5     +1

Slow

0.5

Slow     Warm     Fast     1.0

-10

Fast     0.5     +2

Slow     0.5

Cool

1.0     +1     +2

Overheated

*Assume no discount!*

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$

Penn Engineering

# Convergence*

o How do we know the $V_k$ vectors are going to converge?

o Case 1: If the tree has maximum depth M, then $V_M$ holds the actual untruncated values

o Case 2: If the discount is less than 1
  ▪ Sketch: For any state $V_k$ and $V_{k+1}$ can be viewed as depth k+1 expectimax results in nearly identical search trees
  ▪ The difference is that on the bottom layer, $V_{k+1}$ has actual rewards while $V_k$ has zeros
  ▪ That last layer is at best all $R_{MAX}$
  ▪ It is at worst $R_{MIN}$
  ▪ But everything is discounted by $\gamma^k$ that far out
  ▪ So $V_k$ and $V_{k+1}$ are at most $\gamma^k \max|R|$ different
  ▪ So as k increases, the values converge

$$V_k(s) \qquad V_{k+1}(s)$$

# The Bellman Equations

How to be optimal:

Step 1: Take correct first action

Step 2: Keep being optimal

# The Bellman Equations

o Definition of "optimal utility" via expectimax recurrence gives a simple one-step lookahead relationship amongst optimal utility values

$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^*(s') \right]$$

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^*(s') \right]$$

s

a

s, a

s,a,s'

s'

o These are the Bellman equations, and they characterize optimal values in a way we'll use over and over

# Policy Methods

# Policy Evaluation

# Fixed Policies

Do the optimal action

Do what $\pi$ says to do



o Expectimax trees max over all actions to compute the optimal values

o If we fixed some policy $\pi(s)$, then the tree would be simpler – only one action per state
  ▪ … though the tree's value would depend on which policy we fixed

Penn Engineering

# Utilities for a Fixed Policy

o  Another basic operation: compute the utility of a state s under a fixed (generally non-optimal) policy

o  Define the utility of a state s, under a fixed policy $\pi$:

   $V^\pi(s)$ = expected total discounted rewards starting in s and following $\pi$

o  Recursive relation (one-step look-ahead / Bellman equation):

$$V^\pi(s) = \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V^\pi(s')]$$

s

$\pi(s)$

s, $\pi(s)$

s, $\pi(s)$,s'

s'

# Example: Policy Evaluation

Always Go Right                    Always Go Forward

# Example: Policy Evaluation



Always Go Right

Always Go Forward

# Policy Evaluation

S

○ How do we calculate the V's for a fixed policy π?

π(s)

○ Idea 1: Turn recursive Bellman equations into updates
(like value iteration)

s, π(s)

$$V_0^\pi(s) = 0$$

s, π(s),s'

$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

s'

○ Efficiency: $O(S^2)$ per iteration

○ Idea 2: Without the maxes, the Bellman equations are just a linear system
▪ Solve with Matlab (or your favorite linear system solver)

# Policy Extraction

# Computing Actions from Values



- Let's imagine we have the optimal values V

- How should we act?
  - It's not obvious!

- We need to do a mini-expectimax (one step

$$\pi^*(s) = \arg\max_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma V^*(s')]$$

- This is called policy extraction, since it gets the policy implied by the values

# Computing Actions from Q-Values

o Let's imagine we have the optimal q-values:

o How should we act?
  ▪ Completely trivial to decide!

$$\pi^*(s) = \arg\max_a Q^*(s, a)$$

o Important lesson: actions are easier to select from q-values than values!

# Policy Iteration

# Problems with Value Iteration

○ Value iteration repeats the Bellman updates:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$

○ Problem 1: It's slow – O(S²A) per iteration

○ Problem 2: The "max" at each state rarely changes

○ Problem 3: The policy often converges long before the values

# k=0



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=1



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=2



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=3



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=4



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=5



Noise = 0.2
Discount = 0.9
Living reward = 0

k=6



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=7



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=8



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=9



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=10



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=11



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=12



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=100



Noise = 0.2
Discount = 0.9
Living reward = 0
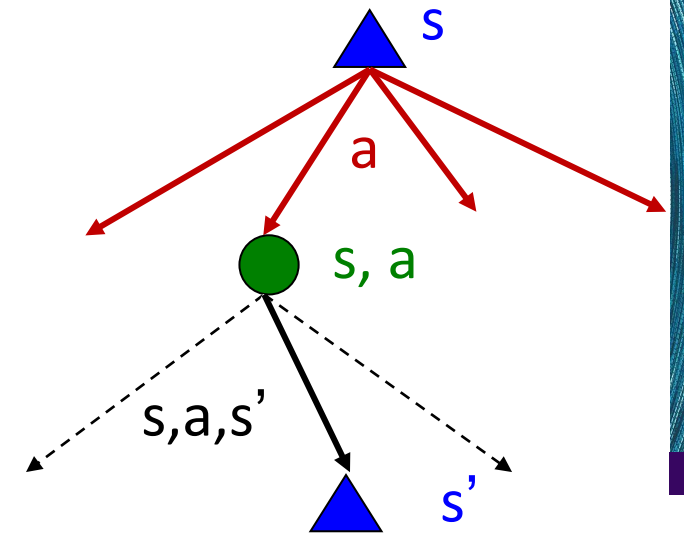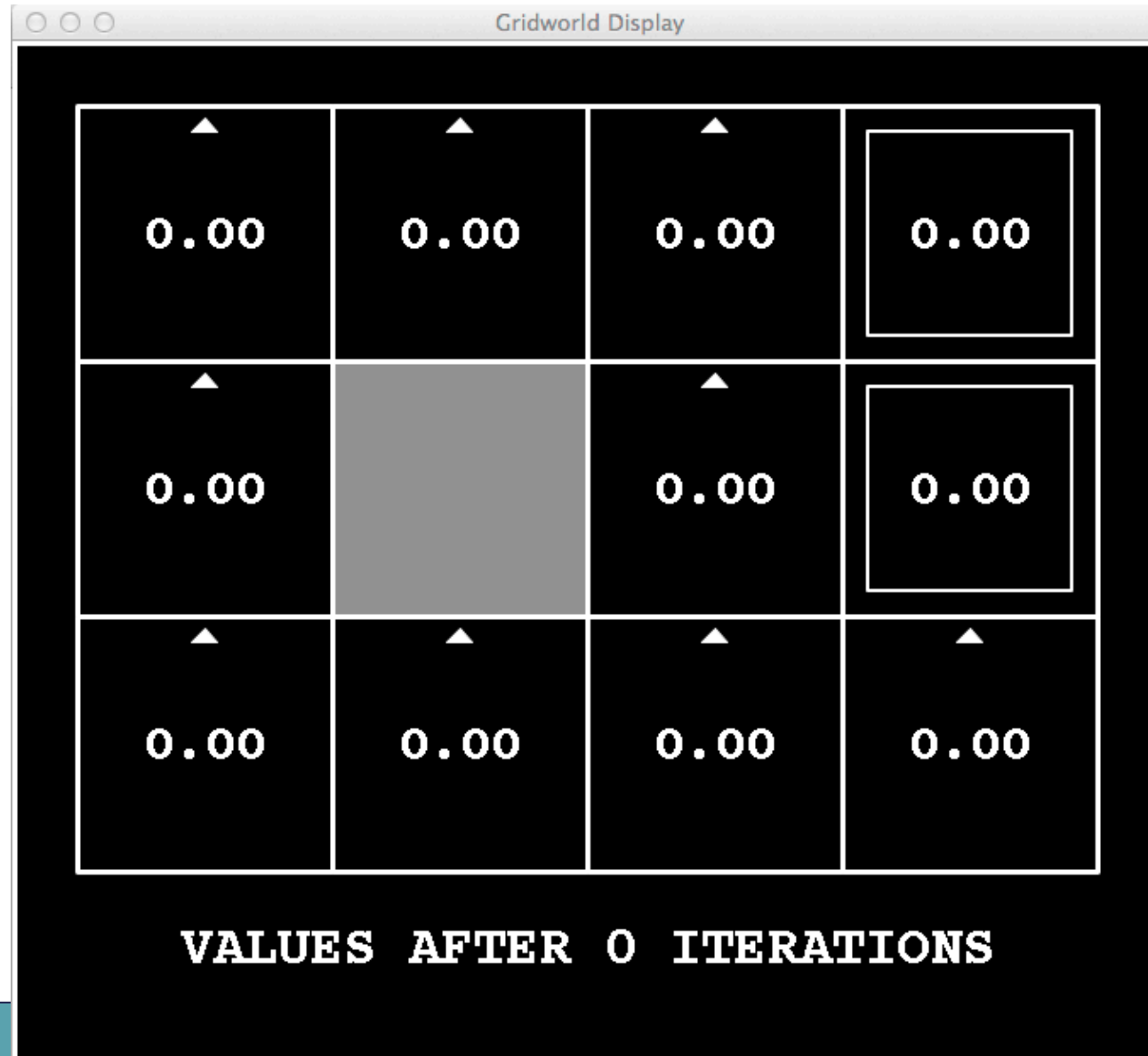
Penn Engineering

# Problems with Value Iteration

- Value iteration repeats the Bellman updates:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$

- Problem 1: It's slow – $O(S^2 A)$ per iteration

- Problem 2: The "max" at each state rarely changes

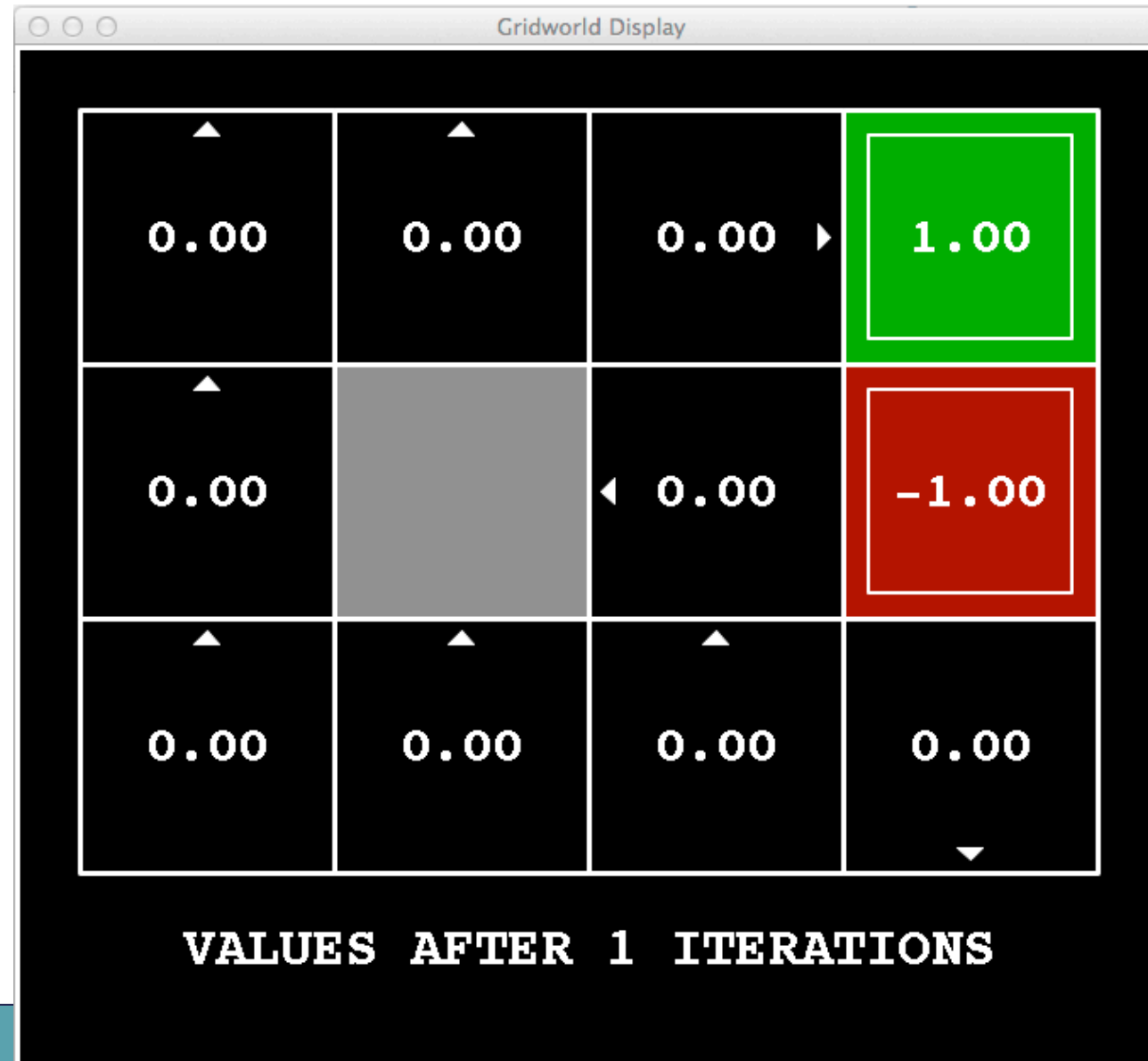- Problem 3: The policy often converges long before the values



s

a

s, a

s,a,s'

s'

# Policy Iteration

○ Alternative approach for optimal values:
- Step 1: Policy evaluation: calculate utilities for some fixed policy (not optimal utilities!) until convergence
- Step 2: Policy improvement: update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values
- Repeat steps until policy converges

○ This is policy iteration
- It's still optimal!
- Can converge (much) faster under some conditions

# Policy Iteration

- Evaluation: For fixed current policy $\pi$, find values with policy evaluation:

  - It $$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') \left[ R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s') \right]$$

- Improvement: For fixed values, get a better policy using policy

  extra $$\pi_{i+1}(s) = \arg\max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^{\pi_i}(s') \right]$$

  - One-step look-ahead.

# Comparison

o Both value iteration and policy iteration compute the same thing (all optimal values)

o In value iteration:
- Every iteration updates both the values and (implicitly) the policy
- We don't track the policy, but taking the max over actions implicitly recomputes it

o In policy iteration:
- We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
- After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
- The new policy will be better (or we're done)

o Both are dynamic programs for solving MDPs

# Summary: MDP Algorithms

- So you want to….
  - Compute optimal values: use value iteration or policy iteration
  - Compute values for a particular policy: use policy evaluation
  - Turn your values into a policy: use policy extraction (one-step lookahead)

- These all look the same!
  - They basically are – they are all variations of Bellman updates
  - They all use one-step lookahead expectimax fragments
  - They differ only in whether we plug in a fixed policy or max over actions

# Maximum Expected Utility

- Why should we average utilities?  Why not minimax?

- Principle of maximum expected utility:
  - A rational agent should choose the action that <span style="color:red">maximizes its expected utility, given its knowledge</span>
- Questions:
  - Where do utilities come from?
  - How do we know such utilities even exist?
  - How do we know that averaging even makes sense?
  - What if our behavior (preferences) can't be described by utilities?
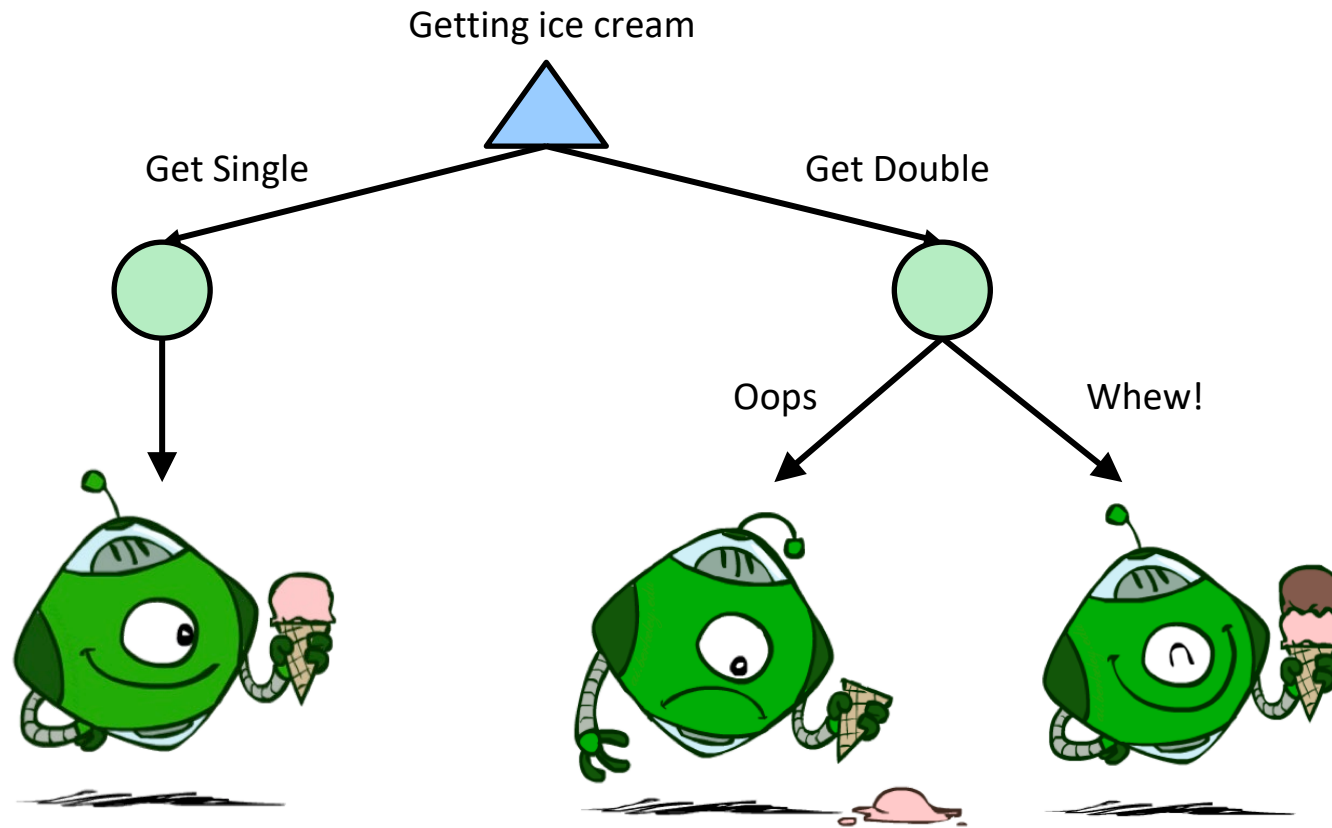
# What Utilities to Use?



- For worst-case minimax reasoning, terminal function scale doesn't matter
  - We just want better states to have higher evaluations (get the ordering right)
  - We call this insensitivity to monotonic transformations

- For average-case expectimax reasoning, we need *magnitudes* to be meaningful

# Utilities

o **Utilities are functions from outcomes (states of the world) to real numbers that describe an agent's preferences**

o Where do utilities come from?
  ▪ In a game, may be simple (+1/-1)
  ▪ Utilities summarize the agent's goals
  ▪ Theorem: any "rational" preferences can be summarized as a utility function

o We hard-wire utilities and let behaviors emerge
  ▪ Why don't we let agents pick utilities?
  ▪ Why don't we prescribe behaviors?

# Utilities: Uncertain Outcomes
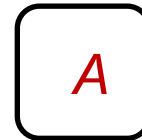
Getting ice cream

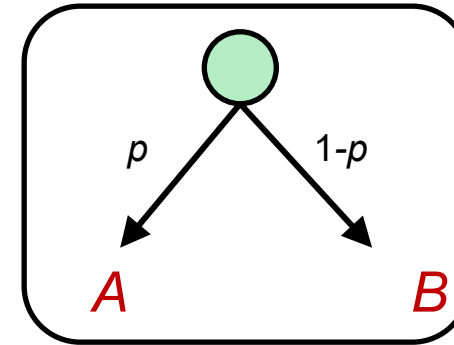Get Single          Get Double

Oops          Whew!

# Preferences

o An agent must have preferences among:

- Prizes: $A, B$, etc.
- Lotteries: situations with uncertain prizes
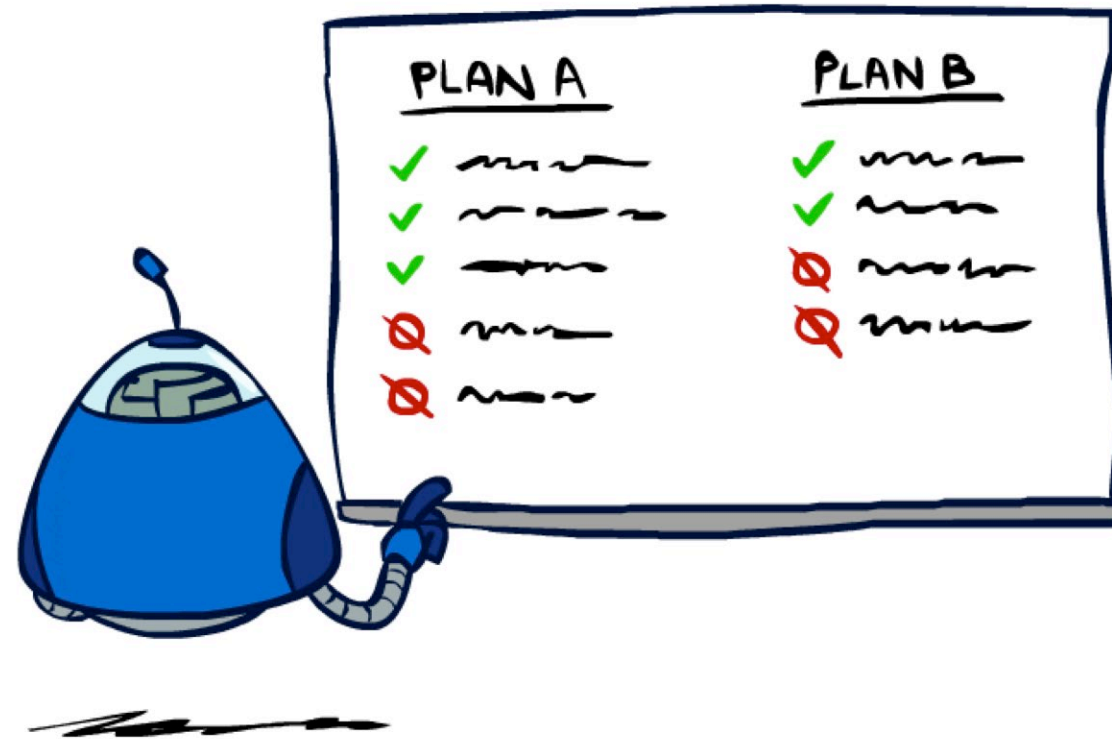$$L = [p, A;\ (1-p), B]$$

A Prize



A Lottery



o Notation: $A \succ B$

- Preference $A \sim B$
- Indifference:

# Rationality

# Rational Preferences

o We want some constraints on preferences before we call them rational, such as:

Axiom of Transitivity: $(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C)$

o For example: an agent with intransitive preferences can be induced to give away all of its money
- If B > C, then an agent with C would pay (say) 1 cent to get B
- If A > B, then an agent with B would pay (say) 1 cent to get A
- If C > A, then an agent with A would pay (say) 1 cent to get C

# Rational Preferences

The Axioms of Rationality

Orderability
$$(A \succ B) \vee (B \succ A) \vee (A \sim B)$$
Transitivity
$$(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C)$$
Continuity
$$A \succ B \succ C \Rightarrow \exists p \; [p, A; \; 1-p, C] \sim B$$
Substitutability
$$A \sim B \Rightarrow [p, A; \; 1-p, C] \sim [p, B; 1-p, C]$$
Monotonicity
$$A \succ B \Rightarrow$$
$$(p \geq q \Leftrightarrow [p, A; \; 1-p, B] \succeq [q, A; \; 1-q, B])$$

RATIONALITY CONFIRMED

Theorem: Rational preferences imply behavior describable as maximization of expected utility
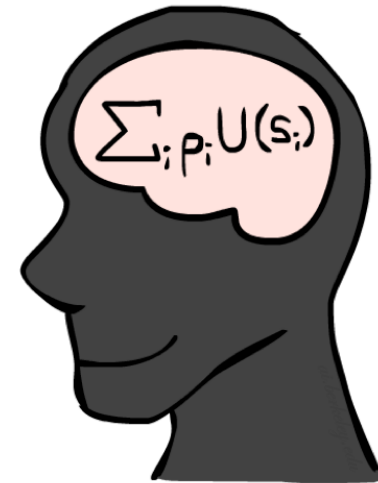
# MEU Principle

o Theorem [Ramsey, 1931; von Neumann & Morgenstern, 1944]
- Given any preferences satisfying these constraints, there exists a real-valued function U such that:
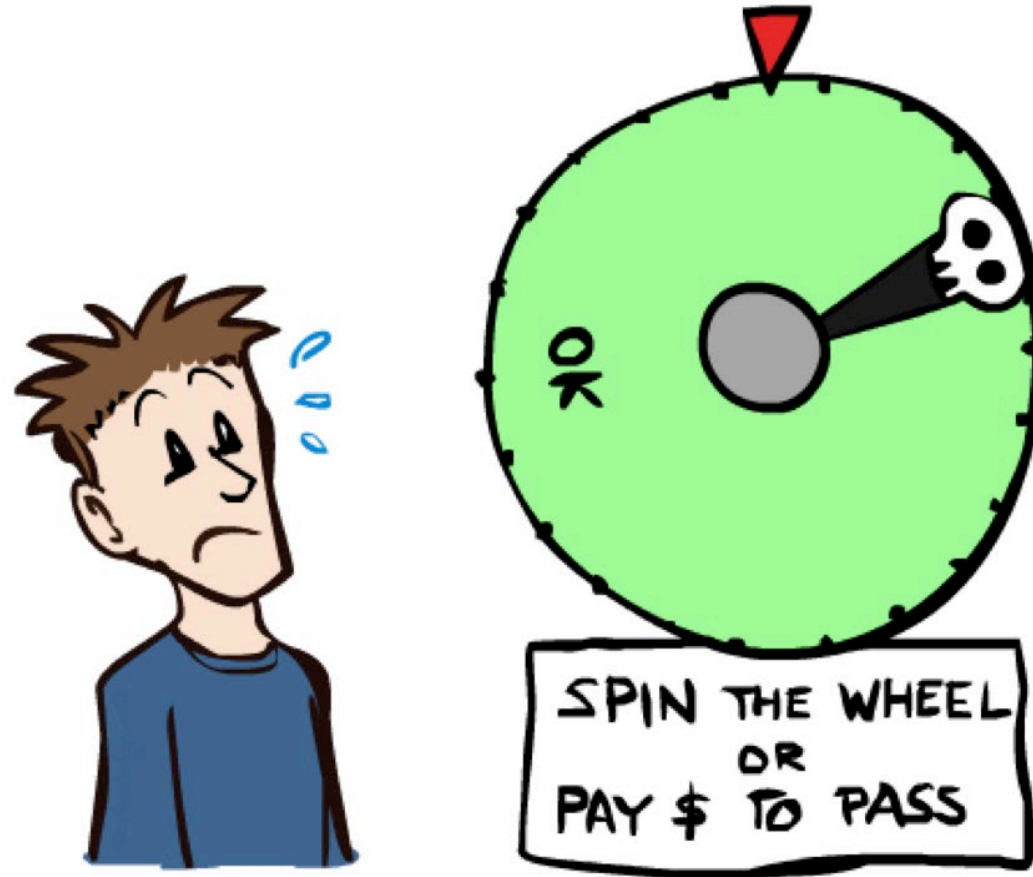
$$U(A) \geq U(B) \iff A \succeq B$$

$$U([p_1, S_1; \ \ldots \ ; \ p_n, S_n]) = \sum_i p_i U(S_i)$$

- I.e. values assigned by U preserve preferences of both prizes and lotteries!

o Maximum expected utility (MEU) principle:
- Choose the action that maximizes expected utility
- Note: an agent can be entirely rational (consistent with MEU) without ever representing or manipulating utilities and probabilities
- E.g., a lookup table for perfect tic-tac-toe, a reflex vacuum cleaner

# Human Utilities

# Utility Scales

o **Normalized utilities**: $u_+ = 1.0$, $u_- = 0.0$

o **Micromorts**: one-millionth chance of death, useful for paying to reduce product risks, etc.

o **QALYs**: quality-adjusted life years, useful for medical decisions involving substantial risk

o Note: behavior is invariant under positive linear transformation

$$U'(x) = k_1 U(x) + k_2 \quad \text{where } k_1 > 0$$

o With deterministic prizes only (no lottery choices), only **ordinal utility** can be determined, i.e., total order on prizes
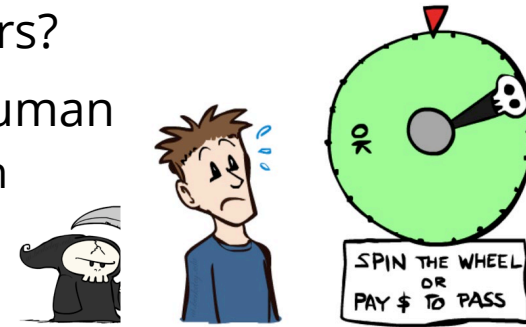
# Micromort examples

| Death from | Micromorts per exposure |
|---|---|
| Scuba diving | 5 per dive |
| Skydiving | 7 per jump |
| Base-jumping | 430 per jump |
| Climbing Mt. Everest | 38,000 per ascent |

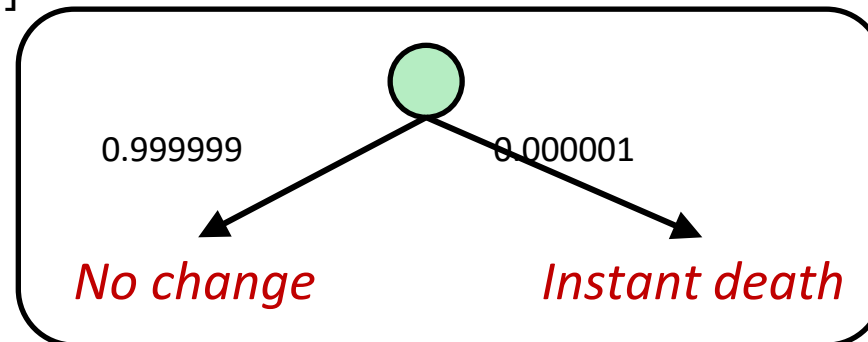| 1 Micromort | |
|---|---|
| Train travel | 6000 miles |
| Jet | 1000 miles |
| Car | 230 miles |
| Walking | 17 miles |
| Bicycle | 10 miles |
| Motorbike | 6 miles |

# Human Utilities

o Utilities map states to real numbers. Which numbers?

o Standard approach to assessment (elicitation) of human
  - Compare a prize A to a standard lottery $L_p$ between
    - "best possible prize" $u_+$ with probability p
    - "worst possible catastrophe" $u_-$ with probability 1-p
  - Adjust lottery probability p until indifference: A ~ $L_p$
  - Resulting p is a utility in [0,1]



Pay $30    ~    0.999999        0.000001
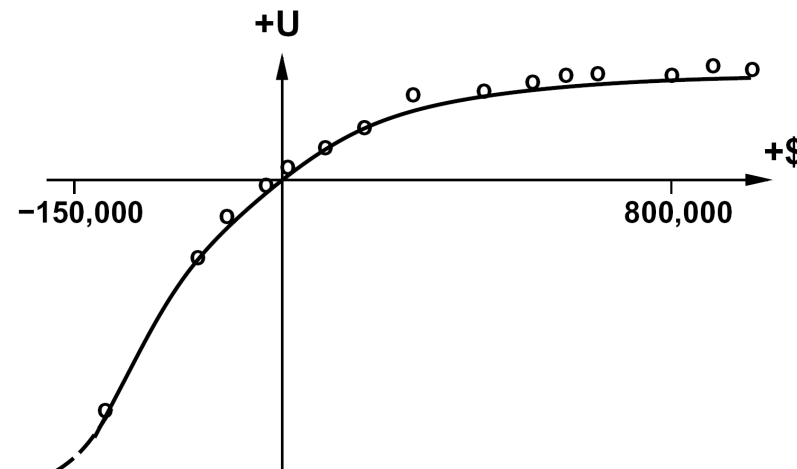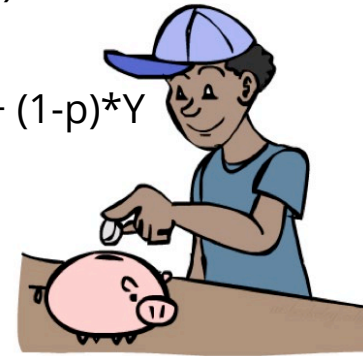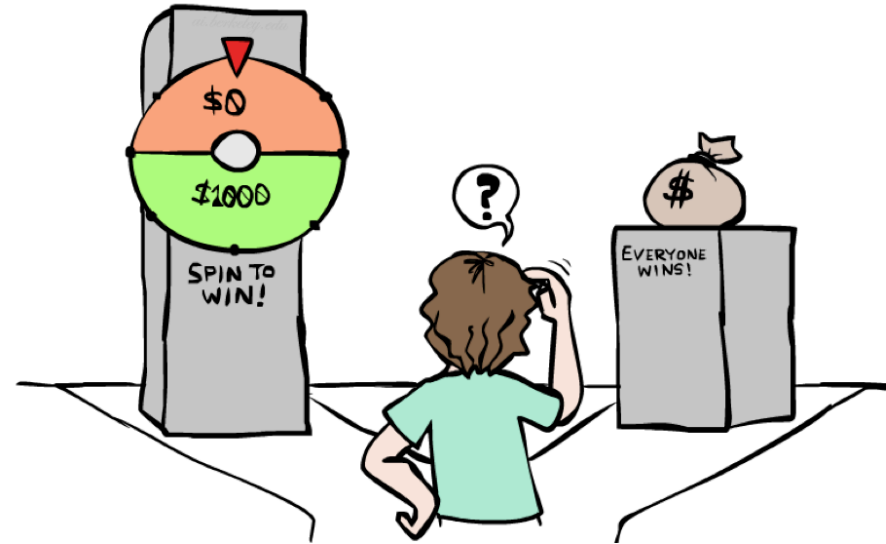                No change       Instant death

# Money

o Money <u>does</u> <u>not</u> behave as a utility function, but we can talk about the utility of having money (or being in debt)

o Given a lottery L = [p, $X; (1-p), $Y]
- The expected monetary value EMV(L) is p*X + (1-p)*Y
- U(L) = p*U($X) + (1-p)*U($Y)
- Typically, U(L) < U( EMV(L) )
- In this sense, people are risk-averse
- When deep in debt, people are risk-prone

# Example: Insurance

- Consider the lottery [0.5, $1000;  0.5, $0]
  - What is its expected monetary value? ($500)
  - What is its certainty equivalent?
    - Monetary value acceptable in lieu of lottery
    - $400 for most people
  - Difference of $100 is the insurance premium
    - There's an insurance industry because people will pay to reduce their risk
    - If everyone were risk-neutral, no insurance needed!
  - It's win-win: you'd rather have the $400 and the insurance company would rather have the lottery (their utility curve is linear and they have many lotteries)

# Example: Human Rationality?

o Famous example of Allais (1953) ⬅

- A: [0.8, $4k;   0.2, $0]
- B: [1.0, $3k;   0.0, $0]

- C: [0.2, $4k;   0.8, $0]
- D: [0.25, $3k;   0.75, $0]

o Most people prefer B > A, C > D

o But if U($0) = 0, then
- B > A $\Rightarrow$ U($3k) > 0.8 U($4k)
- C > D $\Rightarrow$ 0.8 U($4k) > U($3k)