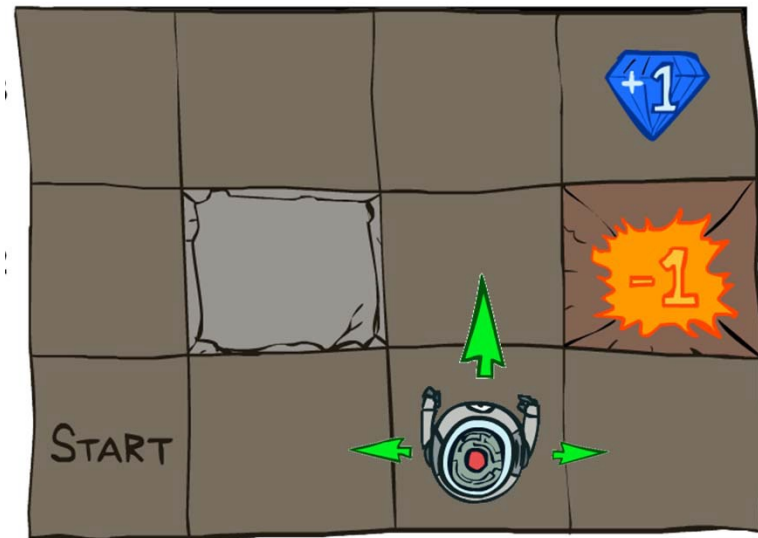


# CS440/ECE448 Lecture 22: Markov Decision Processes

Mark Hasegawa-Johnson

These slides are in the public domain.



Grid World

Invented and drawn by Peter Abbeel and Dan Klein, UC Berkeley CS 188

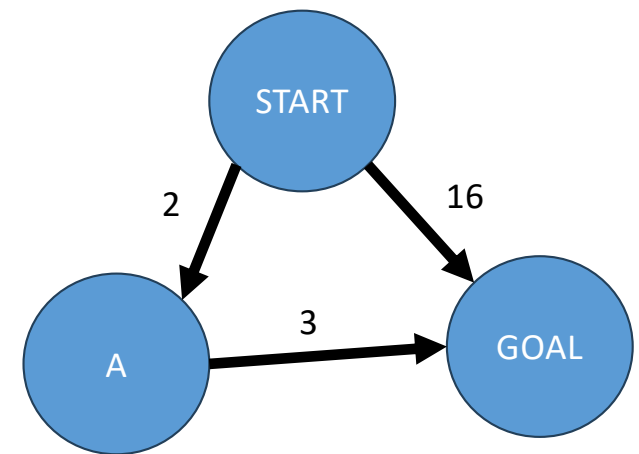
# Outline

- Problem statement
- Bellman equation
- Discount factor
- Policy Iteration
- Value Iteration

# How does an intelligent agent plan its actions?

- If there is no randomness: Use Dijkstra or A\* to plan the best path

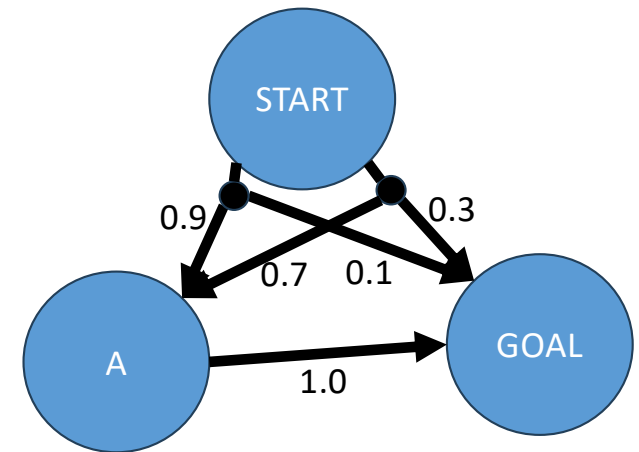
EXPAND	FRONTIER (state, cost, predecessor)
START	(A,2,START), (GOAL,16,START)
A	(GOAL,5,A)
GOAL	



# How does an intelligent agent plan its actions?

- If there is no randomness: Use A\* search to plan the best path
- What if our movements are affected by randomness?

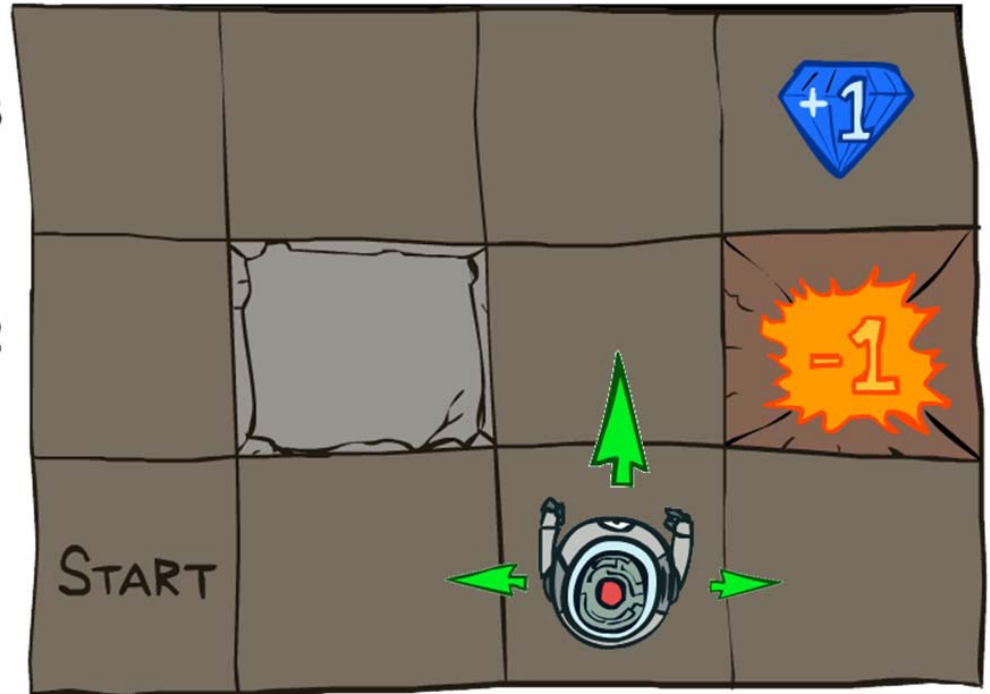
STATE	ACTION	P(go to A)	P(go to GOAL)	REWARD
START	A	0.9	0.1	0
START	GOAL	0.3	0.7	0
A	GOAL	0.0	1.0	0
GOAL	-	-	-	100



# Example: Grid World

Invented by Peter Abbeel and Dan Klein

- Maze-solving problem: state is  $s = (i, j)$ , where  $0 \leq i \leq 2$  is the row and  $0 \leq j \leq 3$  is the column.
- The robot is trying to find its way to the diamond.
- If it reaches the diamond, it gets a reward of  $r((0,3)) = +1$  and the game ends.
- If it falls in the fire it gets a reward of  $r((1,3)) = -1$  and the game ends.

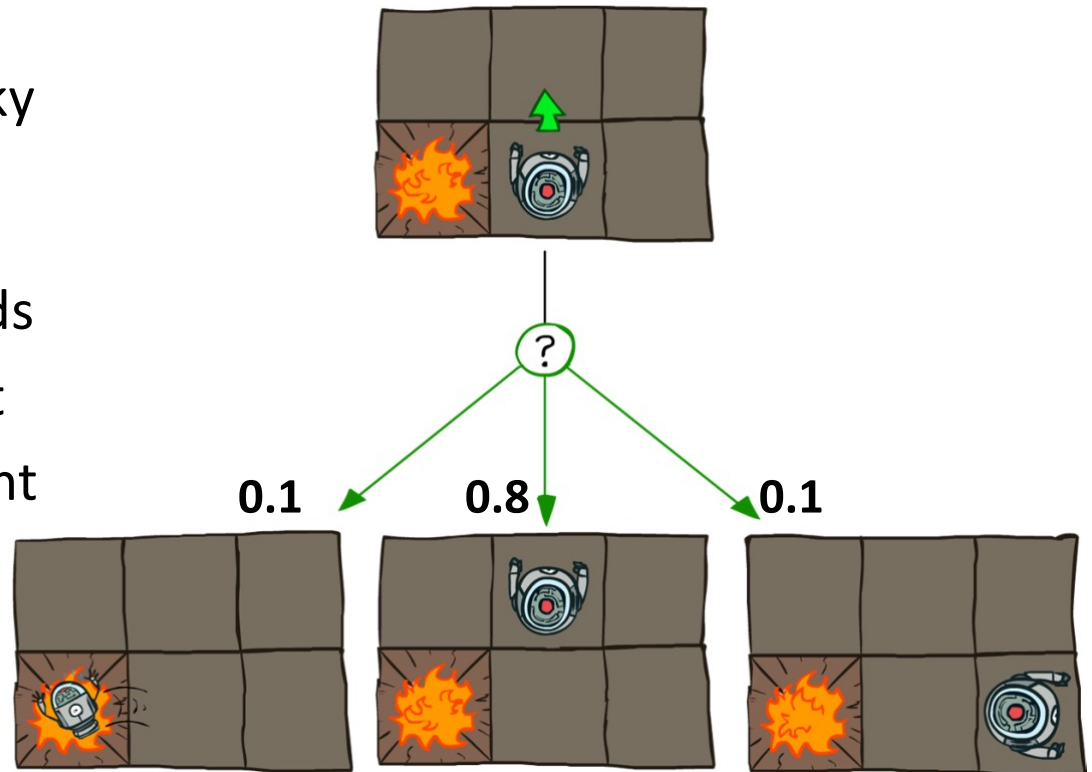


# Example: Grid World

Invented by Peter Abbeel and Dan Klein

Randomness: the robot has shaky actuators. If it tries to move forward,

- With probability 0.8, it succeeds
- With probability 0.1, it falls left
- With probability 0.1, it falls right



# Markov Decision Process

A Markov Decision Process (MDP) is defined by:

- A set of states,  $s \in \mathcal{S}$
- A set of actions,  $a \in \mathcal{A}$
- A transition model,  $P(S_{t+1} = s_{t+1} | S_t = s_t, a_t)$ 
  - $S_t$  is the state at time  $t$
  - $a_t$  is the action taken at time  $t$  (not random)
- A reward function,  $r(s)$

# Solving an MDP: The Policy

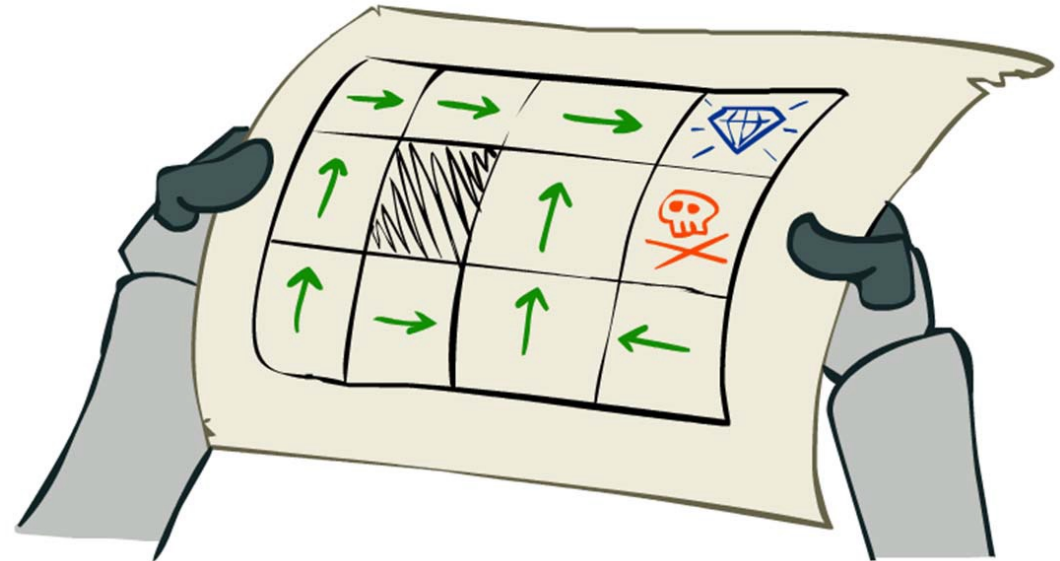
- The solution to a maze is a path: the shortest path from start to goal
- In MDP, finding 1 path is not enough: randomness might cause us to accidentally deviate from the optimal path.

# Solving an MDP: The Policy

Since  $P(S_{t+1} = s_{t+1} | S_t = s_t, a_t)$  and  $r(s)$  depend only on the state (the model is Markov), a complete solution can be expressed as follows:

*If you find yourself in state  $s$ , what should you do?*

A policy,  $a = \pi(s)$ , is a function telling you, for any state  $s$ , what is the best action to take in that state.



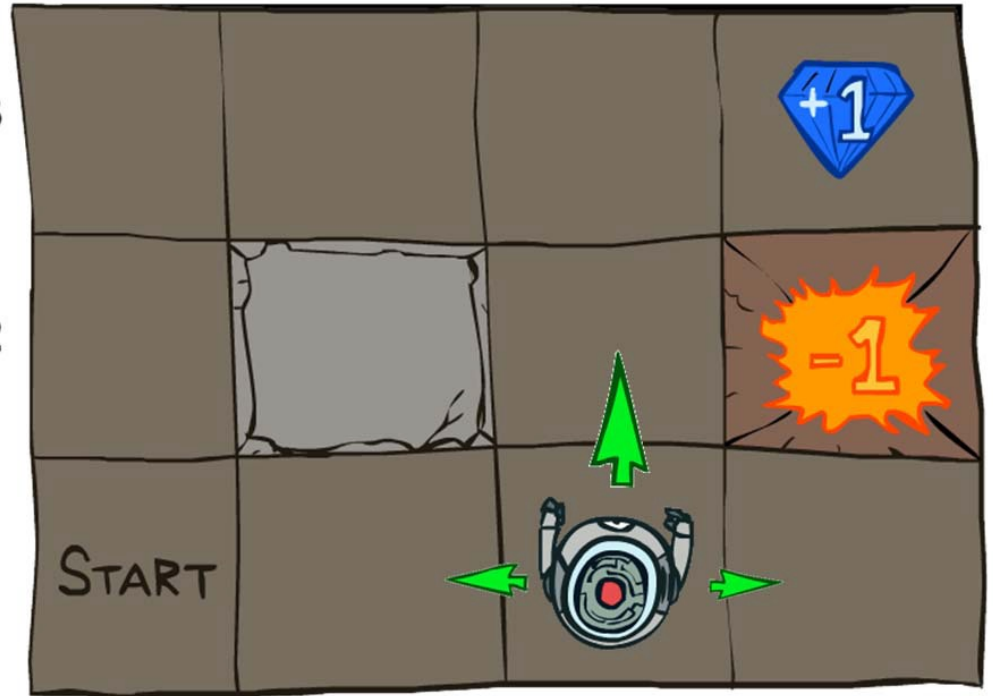
# Outline

- Problem statement
- Bellman equation
- The discount factor
- Policy Iteration
- Value Iteration

# Utility

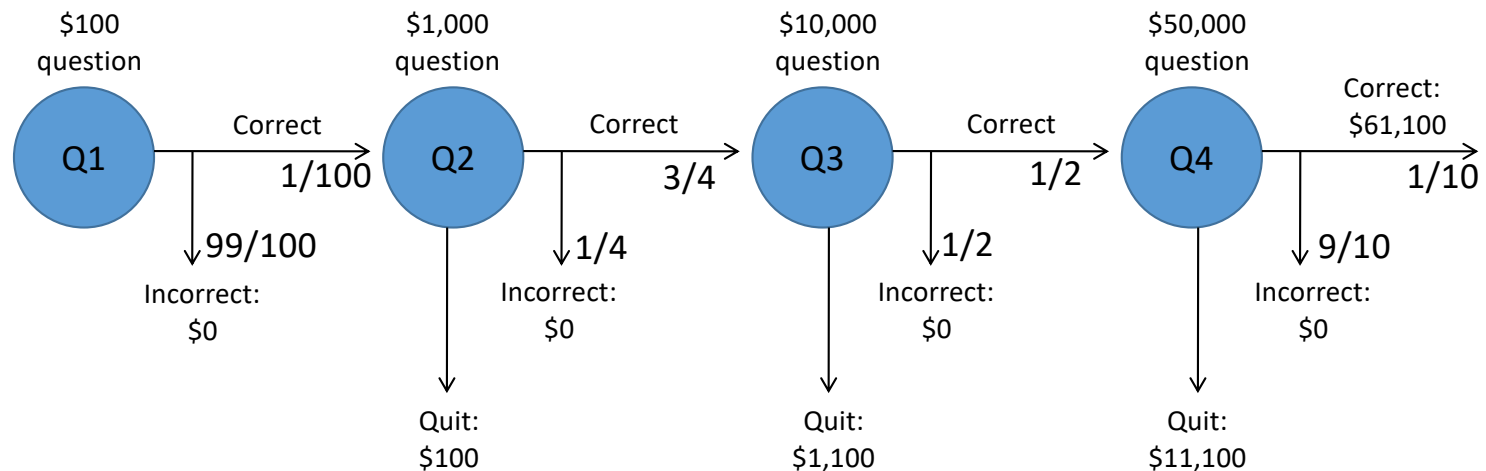
The utility of a state,  $u(s)$ , is defined to be:

- the sum of all current and future rewards that can be achieved if we start in state  $s$ ,
- ...if we choose the best possible sequence of actions,
- ...and if we average over all possible results of those actions.



# Example: Game show

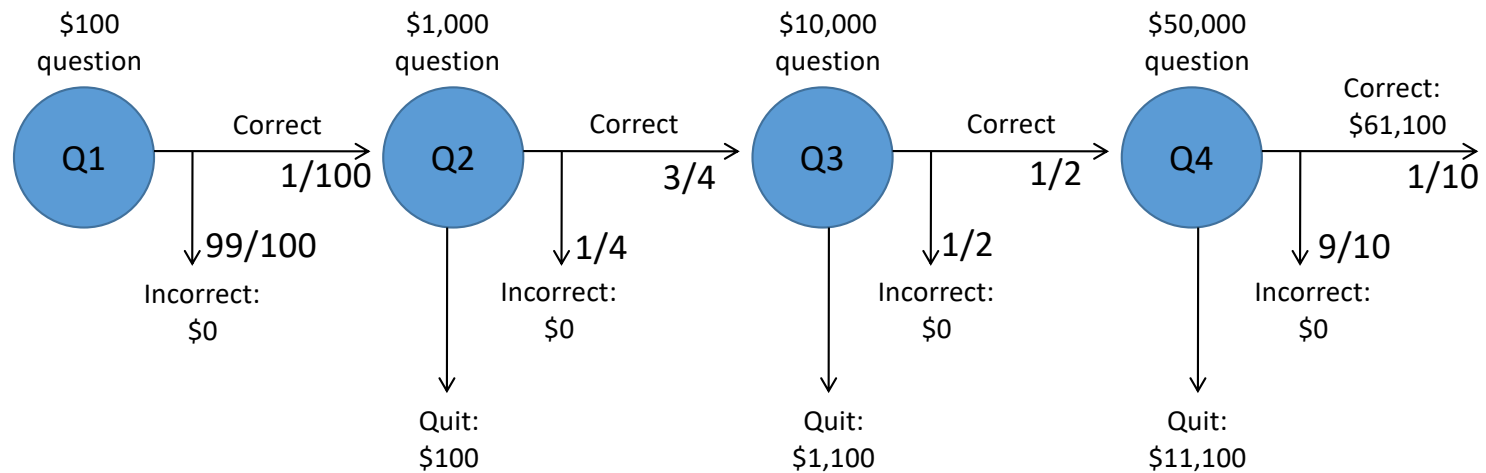
- You've been offered a spot as a contestant in a game show.
- Reward: you receive successively larger prizes for each question you answer correctly, but if you answer any question incorrectly, you lose it all.
- Transition: the questions become harder to answer.
- Actions: after each question you can decide whether to take another question or stop.



# Example: Game show

Policy:

- If you've correctly answered  $N-1$  questions, should you attempt question  $QN$ , or stop?



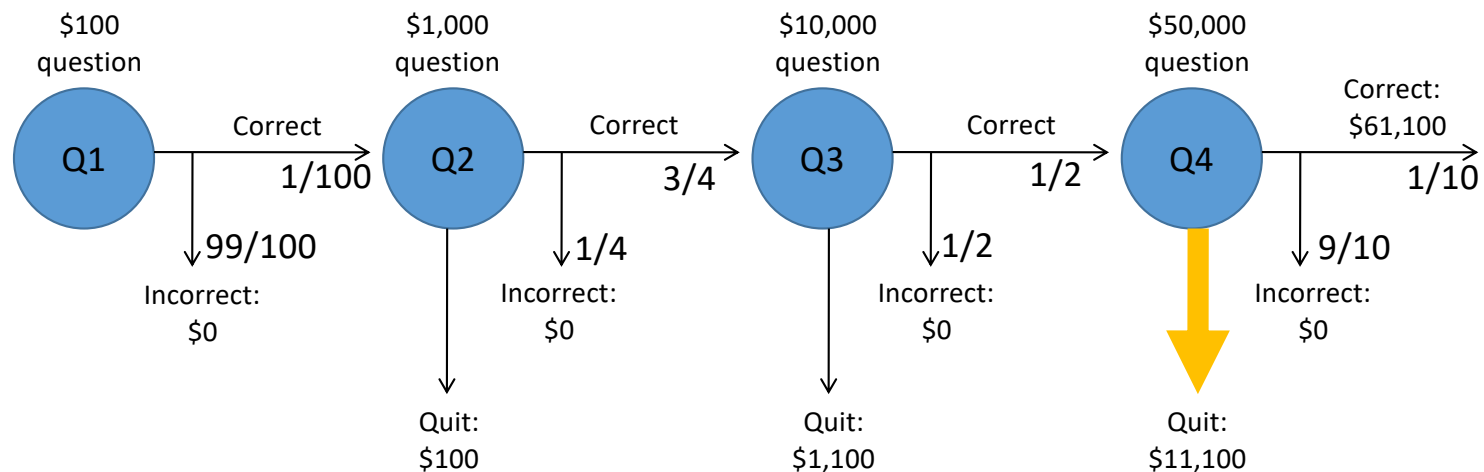
# Example: Game show

Policy  $\pi(Q4)$ : If you've correctly answered 3 questions, should you attempt question Q4, or stop?

- If you stop: total reward is \$11,100
- If you attempt Q4: expected total reward is  $\frac{1}{10} \times 61100 + \frac{9}{10} \times 0 = \$6110$

Policy:  $\pi(Q4) = \text{stop}$ .

Utility:  $u(Q4) = \$11,100$



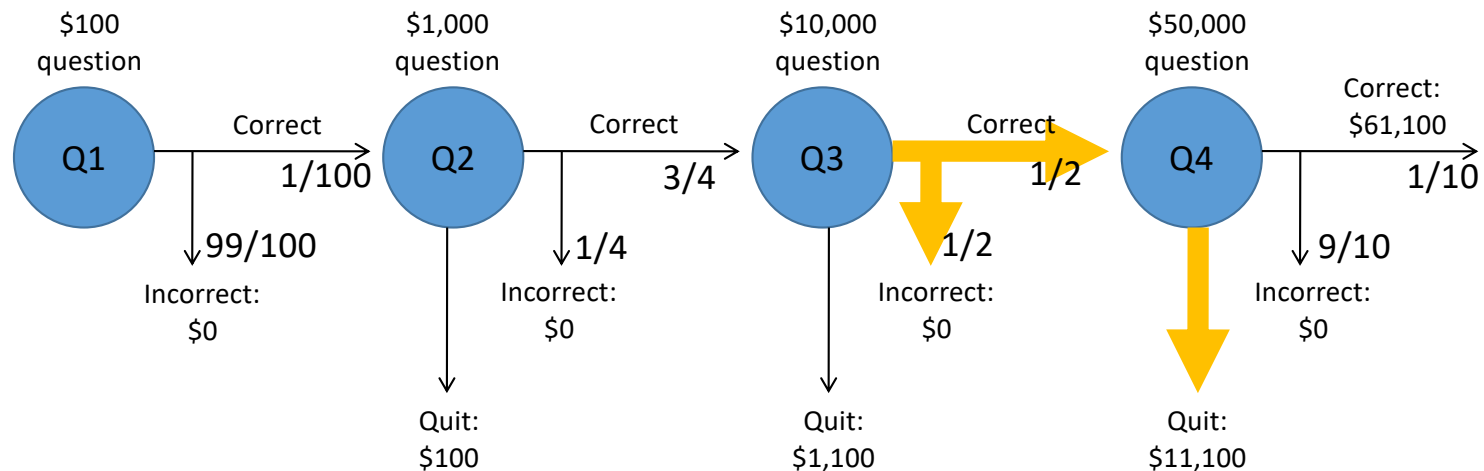
# Example: Game show

Policy  $\pi(Q3)$ : If you've correctly answered 2 questions, should you attempt question Q3, or stop?

- If you stop: total reward is \$1,100
- If you attempt Q3: expected total reward is  $\frac{1}{2} \times \$11,100 + \frac{1}{2} \times 0 = \$5550$

Policy:  $\pi(Q3) = \text{continue}$ .

Utility:  $u(Q3) = \$5550$



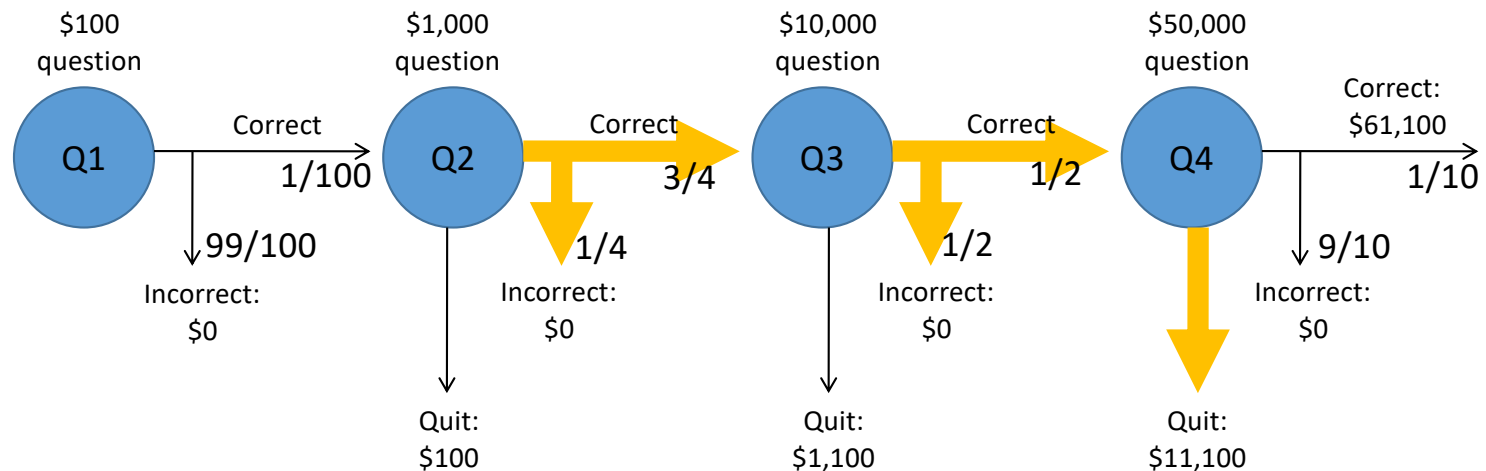
# Example: Game show

Policy  $\pi(Q2)$ : If you've correctly answered 1 question, should you attempt question Q2, or stop?

- If you stop: total reward is \$100
- If you attempt Q2: expected total reward is  $\frac{3}{4} \times \$5550 + \frac{1}{4} \times 0 = \$4162.50$

Policy:  $\pi(Q2) = \text{continue}$ .

Utility:  $u(Q2) = \$4162.50$

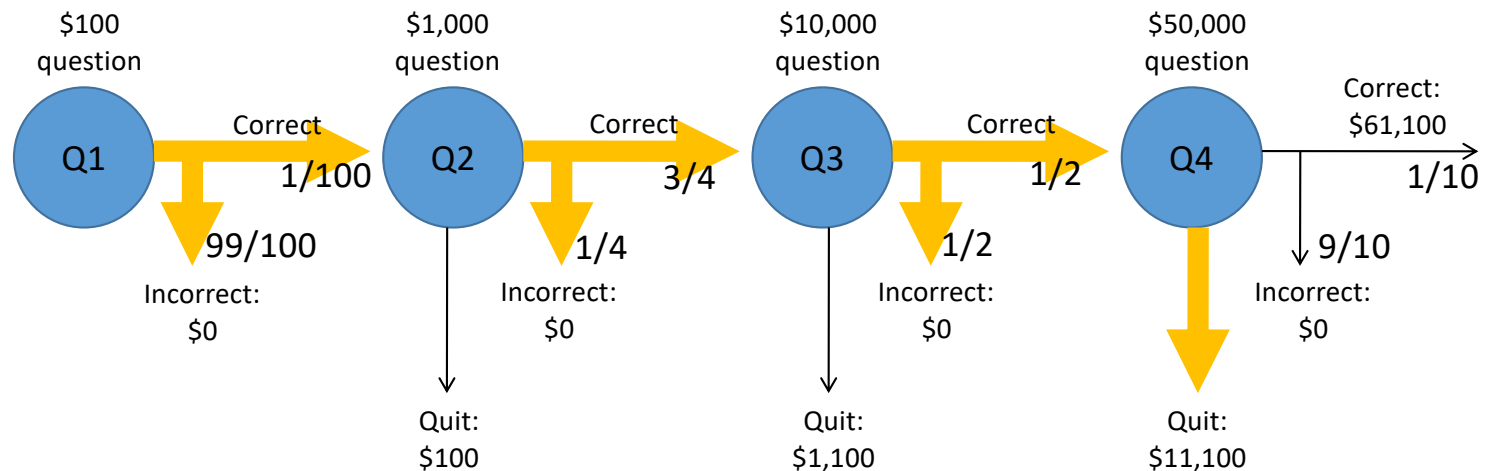


# Example: Game show

Policy  $\pi(Q1)$ : If you've correctly answered no questions, then you have nothing to lose, so even though the chance of success is very small, you might as well try it!

Policy:  $\pi(Q1) = \text{continue}$ .

Utility:  $u(Q1) = \$41.63$



# Utility

The utility of a state,  $u(s)$ , is

- ...the maximum, over all possible sequences of actions, of
- ...the expected value, over all possible results of those actions, of
- ...the total of all future rewards.

$$u(s) = r(s) + \max_a \sum_{s'} P(s'|s, a) \left( r(s') + \max_{a'} \sum_{s''} P(s''|s', a') (r(s'') + \dots \dots \dots) \right)$$

# Utility

The utility of a state,  $u(s)$ , is

- ...the maximum, over all possible actions, of
- ...the expected value, over all possible results of those actions, of
- ...the utility of the resulting state.

$$u(s) = r(s) + \max_a \sum_{s'} P(S_{t+1} = s' | S_t = s, a) u(s')$$

# Bellman Equation

$$u(s) = r(s) + \gamma \max_a \sum_{s'} P(S_{t+1} = s' | S_t = s, a) u(s')$$

- The Bellman equation specifies the utility of every state
- It also tells you the policy, i.e., the solution to the problem!
- The best action to take, in any given state, is the action that maximizes utility, i.e.,

$$\pi(s) = \operatorname{argmax}_a \sum_{s'} P(S_{t+1} = s' | S_t = s, a) u(s')$$

# Outline

- Problem statement
- Bellman equation
- The discount factor
- Policy Iteration
- Value Iteration

# Discount factor

You have just won a contest sponsored by the Galaxia Foundation.  
They offer you the choice of two options:

- \$60,000 right now, or...
- \$1000 per year, paid to you and your heirs annually forever.

Which option is better?

# Discount factor

- Inflation has averaged 3.8% annually from 1960 to 2025.
- Equivalently, \$1000 received one year from now is worth approximately \$962 today.
- A reward of \$1000 annually forever (starting today,  $t=0$ ) is equivalent to an immediate reward of

$$r = \sum_{t=0}^{\infty} 1000(0.962)^t = \frac{1000}{1 - 0.962} = \$26,316$$

We call the factor  $\gamma = 0.962$  the discount factor.

# Discount factor

Why is a dollar tomorrow worth less than a dollar today?

- A dollar will buy less tomorrow
- The person paying you might go out of business
- You might have to go into hiding and become unable to collect

The discount factor,  $\gamma$ , is our model of the unknowable uncertainty of promised future rewards.



Public domain image of J. Wellington Wimpy, the character who popularized the saying "I will gladly pay you Tuesday for a hamburger today."

[https://commons.wikimedia.org/wiki/File:Wimpyh\\_otdog.png](https://commons.wikimedia.org/wiki/File:Wimpyh_otdog.png)

# The Bellman Equation

$$u(s) = r(s) + \gamma \max_a \sum_{s'} P(S_{t+1} = s' | S_t = s, a) u(s') \quad \forall s \in \{1, \dots, n\}$$

- The Bellman equation specifies the utilities of all  $n$  states,  $1 \leq s \leq n$ .
- In solving the Bellman equation, we also find the optimum action, which is the policy.
- However, the “max” makes this  $n$  nonlinear equations ( $1 \leq s \leq n$ ) in  $n$  unknowns (the unknowns are  $u(s)$ ). It has no closed form solution!

# Outline

- Problem statement
- Utility
- The discount factor
- Policy Iteration
- Value Iteration

# Policy iteration step 1: Policy evaluation

The Bellman Equation:

$$u(s) = r(s) + \gamma \max_a \sum_{s'} P(s'|s, a) u(s')$$

Policy iteration starts with a fixed policy. With a fixed policy, the Bellman equation becomes n **linear** equations in n unknowns:

$$u_i(s) = r(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) u_i(s')$$

...which can be solved by inverting a matrix.

# Step 1: Policy Evaluation

**Bellman equation:** n nonlinear equations in n unknowns:

$$\begin{bmatrix} u(1) \\ \vdots \\ u(n) \end{bmatrix} = \begin{bmatrix} r(1) \\ \vdots \\ r(n) \end{bmatrix} + \gamma \max_a \begin{bmatrix} P(1|1, a) & \cdots & P(n|1, a) \\ \vdots & \ddots & \vdots \\ P(1|n, a) & \cdots & P(n|n, a) \end{bmatrix} \begin{bmatrix} u(1) \\ \vdots \\ u(n) \end{bmatrix}$$

**Policy Evaluation:** n linear equations in n unknowns:

$$\begin{bmatrix} u_i(1) \\ \vdots \\ u_i(n) \end{bmatrix} = \begin{bmatrix} r(1) \\ \vdots \\ r(n) \end{bmatrix} + \gamma \begin{bmatrix} P(1|1, \pi_i(1)) & \cdots & P(n|1, \pi_i(1)) \\ \vdots & \ddots & \vdots \\ P(1|n, \pi_i(n)) & \cdots & P(n|n, \pi_i(n)) \end{bmatrix} \begin{bmatrix} u_i(1) \\ \vdots \\ u_i(n) \end{bmatrix}$$

The difference is that policy evaluation is linear, so it can be solved by inverting a matrix:



$$\begin{aligned} \mathbf{u}_i &= \mathbf{r} + \gamma \mathbf{P}_i \mathbf{u}_i \\ \mathbf{u}_i &= (\mathbf{I} - \gamma \mathbf{P}_i)^{-1} \mathbf{r} \end{aligned}$$

# Example: Grid World



**Policy Evaluation:**  $u_i(s) = r(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) u_i(s')$

- Assume the initial policy is  $\pi_1(s) = \text{“Go Right”}$  for all states
- Solve the linear equations to find  $u_1(s)$

$u_1(s)$

+0.50	+0.69	+0.74	
-0.65		-0.90	
-1.40	-1.44	-1.39	-1.40

$\pi_1(s)$

→	→	→	
→		→	
→	→	→	→

## Policy iteration step 2: Policy improvement

The Bellman Equation:

$$u(s) = r(s) + \gamma \max_a \sum_{s'} P(s'|s, a) u(s')$$

Policy evaluation assumes a fixed policy  $\pi_i(s)$  and computes the resulting  $u_i(s)$ . But now that we know the utility of each state, we can improve the policy:

$$\pi_{i+1}(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) u_i(s')$$



...which can be solved by inverting a matrix.

# Policy Improvement



**Policy Evaluation:**  $u_i(s) = r(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) u_i(s')$

**Policy Improvement:**  $\pi_{i+1}(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) u_i(s')$



$\pi_2(s)$

→	→	→	
↑	■	↑	
↑	→	↑	↑

$u_1(s)$

+0.50	+0.69	+0.74	
-0.65	■	-0.90	
-1.40	-1.44	-1.39	-1.40

$\pi_1(s)$

→	→	→	
→	■	→	
→	→	→	→

# Policy Iteration

- **Policy Evaluation:**  $u_i(s) = r(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) u_i(s')$ 
  - Given a **fixed** policy  $\pi_i(s)$ ,
  - Calculate the resulting utility  $u_i(s)$ .
- **Policy Improvement:**  $\pi_{i+1}(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) u_i(s')$ 
  - Given a **fixed** utility  $u_i(s)$ ,
  - Find an improved  $\pi_{i+1}(s)$ .
- Guaranteed to converge in a finite but possibly large number of steps (less than or equal to the number of distinct policies).
- Takes a long time to converge if your initial guess,  $\pi_1(s)$ , is bad.

# Outline

- Problem statement
- Bellman equation
- The discount factor
- Policy Iteration
- **Value Iteration**

# Value iteration

The Bellman Equation:

$$u(s) = r(s) + \gamma \max_a \sum_{s'} P(s'|s, a) u(s')$$

Value iteration solves the Bellman equation iteratively. In iteration number  $i$ , for  $i = 0, 1, \dots$ ,

- For all states  $s$ ,  $u_i(s)$  is an estimate of  $u(s)$
- Start out with  $u_0(s) = 0$  for all states
- In the  $i^{\text{th}}$  iteration,

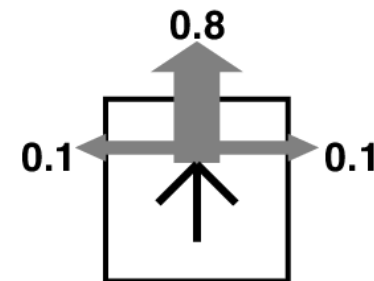
$$u_i(s) = r(s) + \gamma \max_a \sum_{s'} P(s'|s, a) u_{i-1}(s')$$

# Example: Grid world

$r(s)$

3	-0.04	-0.04	-0.04	+1
2	-0.04		-0.04	-1
1	START	-0.04	-0.04	-0.04
	1	2	3	4

Transition model  $P(s'|s, a)$ :





Assume a “loitering penalty” of  $r(s) = -0.04$  for all non-terminal states.



# Value Iteration: Iteration 1

$$u_1(s) = r(s) + \gamma \max_a \sum_{s'} P(s'|s, a) u_0(s')$$



$u_1(s)$

-0.04	-0.04	-0.04	
-0.04		-0.04	
-0.04	-0.04	-0.04	-0.04

$r(s)$

-0.04	-0.04	-0.04	
-0.04		-0.04	
-0.04	-0.04	-0.04	-0.04



$u_0(s)$

0	0	0	
0		0	
0	0	0	0

# Value Iteration: Iteration 2



$$u_2(s) = r(s) + \gamma \max_a \sum_{s'} P(s'|s, a) u_1(s')$$

$u_2(s)$

-0.08	-0.08	+0.75	
-0.08		-0.08	
-0.08	-0.08	-0.08	-0.08



=

$r(s)$



-0.04	-0.04	-0.04	
-0.04		-0.04	
-0.04	-0.04	-0.04	-0.04

+  $\gamma \max_a$



$$\sum_{s'} P(s'|s, \text{down}) u_1(s')$$

-0.04	-0.04	+0.06	
-0.04		-0.14	
-0.04	-0.04	-0.04	-0.04



$$\sum_{s'} P(s'|s, \text{up}) u_1(s')$$

-0.04	-0.04	+0.06	
-0.04		-0.14	
-0.04	-0.04	-0.04	-0.81

$$\sum_{s'} P(s'|s, \text{left}) u_1(s')$$

-0.04	-0.04	-0.04	
-0.04		-0.04	
-0.04	-0.04	-0.04	-0.14

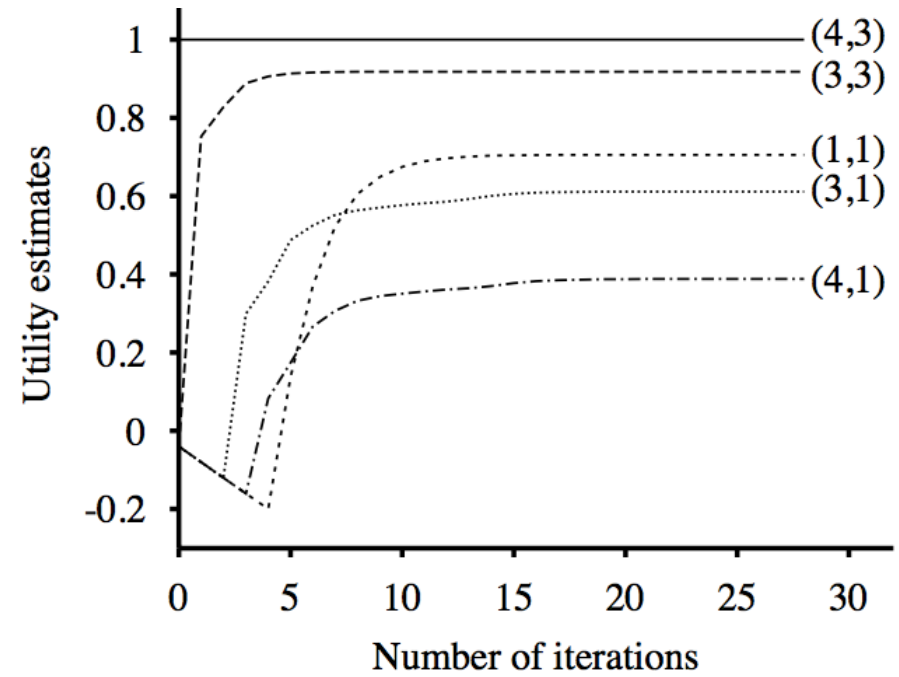
$$\sum_{s'} P(s'|s, \text{right}) u_1(s')$$

-0.04	-0.04	+0.79	
-0.04		-0.81	
-0.04	-0.04	-0.04	-0.14

# Value iteration

Optimal utilities with discount factor 1  
(Result of value iteration)

<b>3</b>	<b>0.812</b>	<b>0.868</b>	<b>0.918</b>	<b>+1</b>
<b>2</b>	<b>0.762</b>		<b>0.660</b>	<b>-1</b>
<b>1</b>	<b>0.705</b>	<b>0.655</b>	<b>0.611</b>	<b>0.388</b>
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>



Final policy

<b>3</b>	→	→	→	<b>+1</b>
<b>2</b>	↑		↑	<b>-1</b>
<b>1</b>	↑	←	←	←
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>

# Value iteration

$$u_i(s) = r(s) + \gamma \max_a \sum_{s'} P(s'|s, a) u_{i-1}(s')$$

- After  $i$  iterations,  $u_i(s)$  has information about the rewards earned in the first  $i$  steps after the agent starts the maze
- A policy designed based on  $u_i(s)$  maximizes reward in the first  $i$  steps of the maze
- In this sense, it's kind of like BFS: each iteration explores one step farther into the future.
- If the maze contains infinite paths, value iteration might never completely converge, however...
- If the biggest rewards/costs are  $M$  steps away, value iteration converges exponentially quickly for  $i > M$

# Quiz

Try the quiz!

# Summary

- Bellman equation: n nonlinear equations in n unknowns

$$u(s) = r(s) + \gamma \max_a \sum_{s'} P(S_{t+1} = s' | S_t = s, a) u(s')$$

- Policy iteration: converges quickly if your initial guess  $\pi_1(s)$  is good

$$u_i(s) = r(s) + \gamma \sum_{s'} P(S_{t+1} = s' | S_t = s, \pi_i(s)) u_i(s')$$

$$\pi_{i+1}(s) = \operatorname{argmax}_a \sum_{s'} P(S_{t+1} = s' | S_t = s, a) u_i(s')$$

- Value iteration: converges exponentially quickly after at least one path reaches each reward

$$u_i(s) = r(s) + \gamma \max_a \sum_{s'} P(S_{t+1} = s' | S_t = s, a) u_{i-1}(s')$$