

Uof  
SC



# Machine Learning: Dynamic Programming

Forest Agostinelli

University of South Carolina

# Topics Covered in This Class

- **Part 1: Search**

- Pathfinding
  - Uninformed search
  - Informed search
- Adversarial search
- Optimization
  - Local search
  - Constraint satisfaction

- **Part 2: Knowledge Representation and Reasoning**

- Propositional logic
- First-order logic
- Prolog

- **Part 3: Knowledge Representation and Reasoning Under Uncertainty**

- Probability
- Bayesian networks

- **Part 4: Machine Learning**

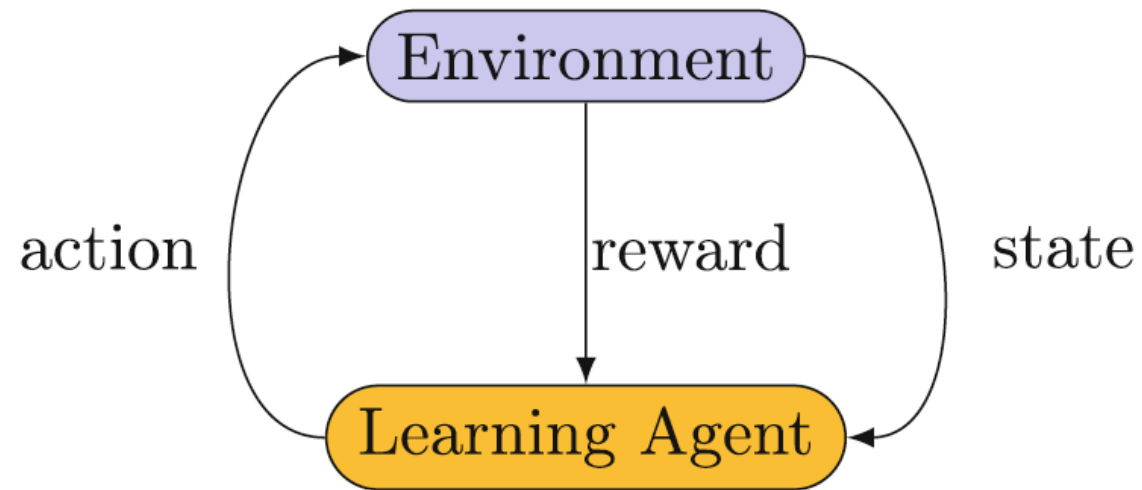
- Supervised learning
  - Inductive logic programming
  - Linear models
  - Deep neural networks
  - PyTorch
- Reinforcement learning
  - Markov decision processes
  - Dynamic programming
  - Model-free RL
- Unsupervised learning
  - Clustering
  - Autoencoders

# Outline

- Background
- Policy Evaluation
- Policy Improvement
- Policy Iteration
  - Modified Policy Iteration
  - Value Iteration
- Approximate value iteration

# Reinforcement Learning

- **Reinforcement learning:** learning to map **states** to **actions** so that we maximize the expected future **reward** we receive from the **environment**.
- This mapping of states to actions is called a **policy function**.
  - Deterministic:  $a = \pi(s)$
  - Stochastic:  $\pi(a|s) = P(A = a|S = s)$
- At each time step  $t$ 
  - In state  $S_t$ , agent takes action  $A_t$
  - Based on state  $s_t$  and action  $a_t$ , the environment transitions to state  $S_{t+1}$  and outputs reward  $R_{t+1}$



# Markov Decision Processes (MDPs)

- **States**
- **Actions**
- **Transition Probabilities:**  $P(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$ 
  - Defines the dynamics of the MDP
- The state-transition probabilities can be obtained from the transition probabilities
  - $p(s' | s, a) = \sum_{r \in \mathcal{R}} p(s', r | s, a)$
- The expected reward can be obtained from the transition probabilities
  - $r(s, a) = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a) = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$
- For now, we assume the state and actions are discrete and finite, however, this restriction can be relaxed to be continuous and infinite

# MDPs: Returns and Value

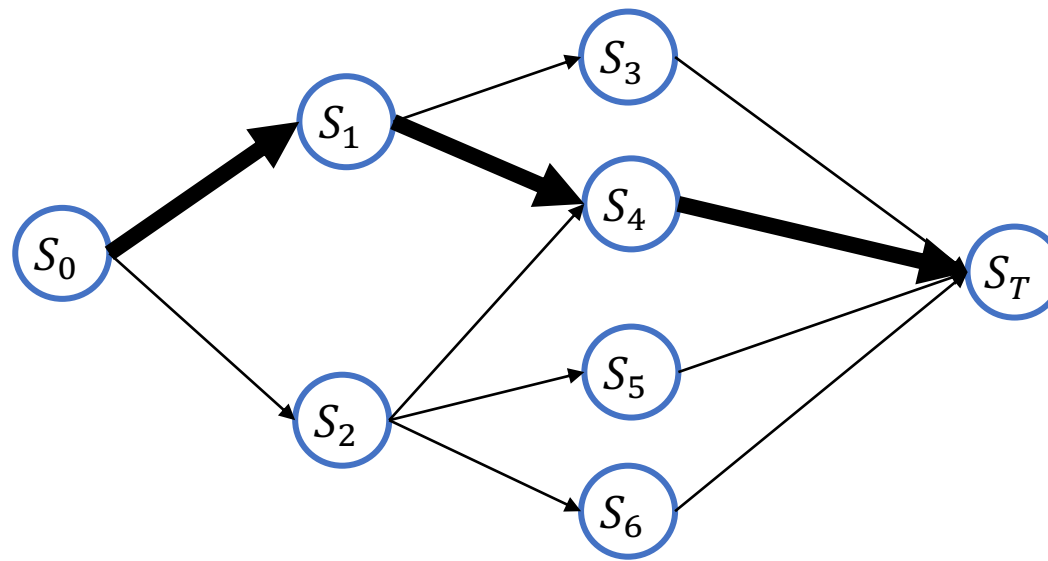
- **Return:** the sum of rewards after timestep  $t$ 
  - $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots$
  - We seek to maximize the expected return
- **State-value function**
  - $v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+1+k} | S_t = s]$
  - $v_*(s) = \max_{\pi} v_\pi(s)$
- **Action-value function**
  - $q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+1+k} | S_t = s, A_t = a]$
  - $q_*(s, a) = \max_{\pi} q_\pi(s, a)$
- Value functions are specific to a given policy  $\pi$

# Dynamic Programming

- Solves problems by recursively breaking them down into simpler subproblems
- Requires
  - **Optimal substructure**: Can construct an optimal solution from optimal solutions of subproblems

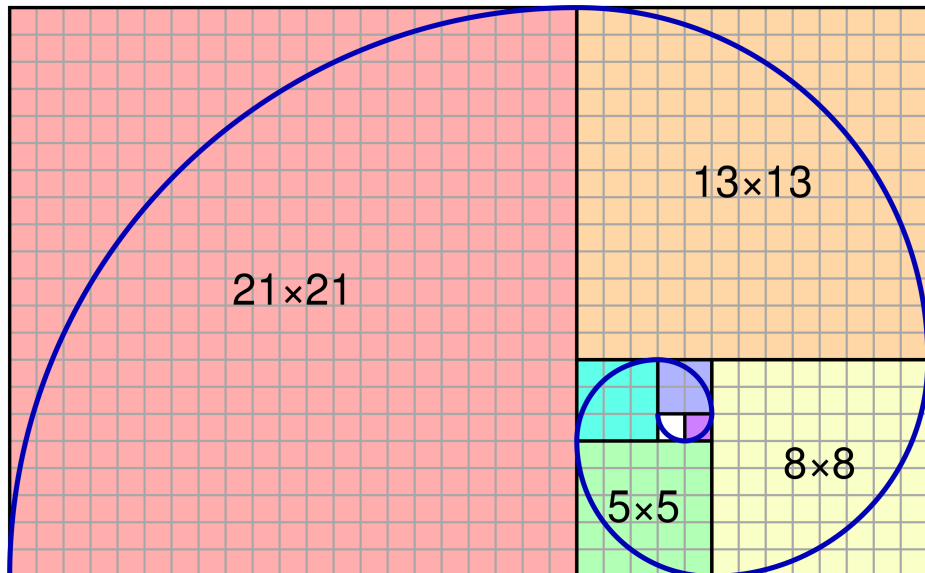
$$v_*(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_*(s'))$$

- Principle of optimality
- **Overlapping subproblems**: Solutions to subproblems are re-used
  - Value functions

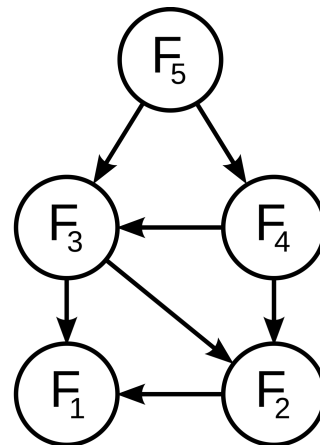


# Dynamic Programming

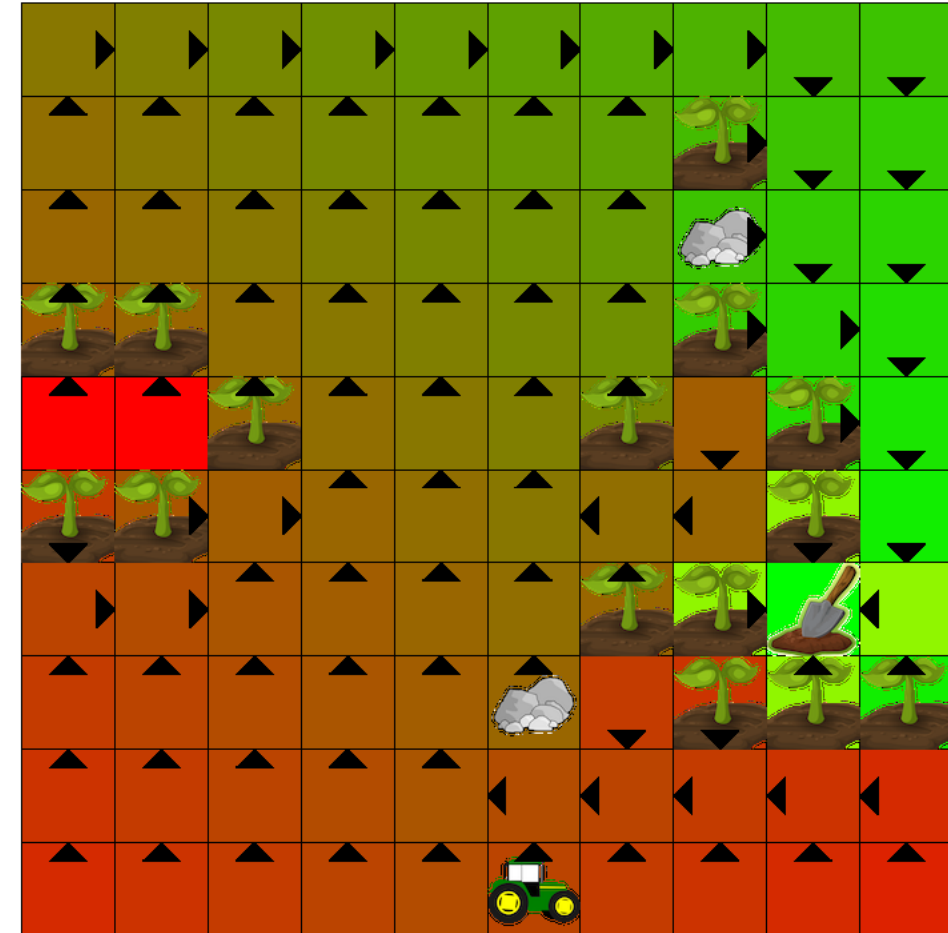
- Fibonacci sequence
- Scheduling
- Sequence alignment (DNA)
- AI Farm



By Jahobr - Own work, CC0, <https://commons.wikimedia.org/w/index.php?curid=58460223>



By en:User:Dcoatze, traced by User:Stannered - en:Image:Fibonacci dynamic programming.png, Public Domain, <https://commons.wikimedia.org/w/index.php?curid=3325402>





# Dynamic Programming: History

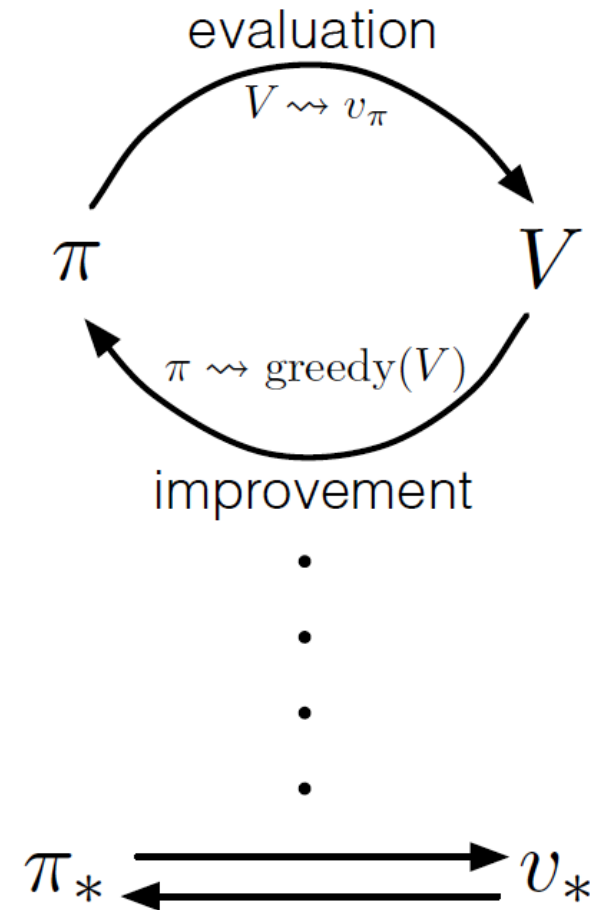
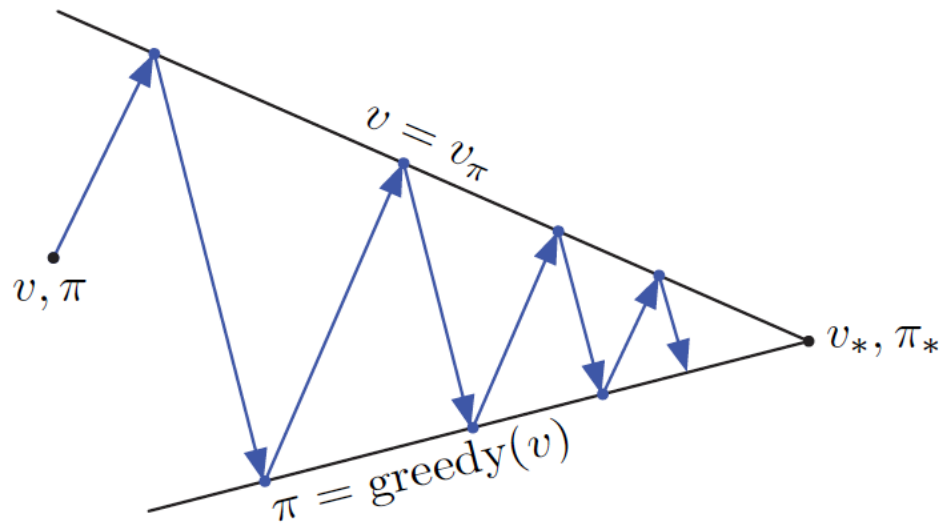
- The term *dynamic programming* was originally used in the 1940s by Richard Bellman to describe the process of solving problems where one needs to find the best decisions one after another.
- By 1953, he refined this to the modern meaning, referring specifically to nesting smaller decision problems inside larger decisions.
- The word *dynamic* was chosen by Bellman to capture the time-varying aspect of the problems, and because it sounded impressive.
- The word *programming* referred to the use of the method to find an optimal *program*.

# Dynamic Programming

- We will use it to evaluate a policy and compute an optimal policy given a perfect model an MDP
  - $p(s', r|s, a)$
- Foundational for reinforcement learning
- Using dynamic programming, we will do **policy iteration** by iterating between **policy evaluation** and **policy improvement**

# Generalized Policy Iteration

- **Policy Evaluation:** Estimate the expected future reward when following policy  $\pi$
- **Policy Improvement:** Improve policy  $\pi$  so that it obtains a greater expected future reward
- We can obtain an optimal policy by iterating between **policy evaluation** and **policy improvement**

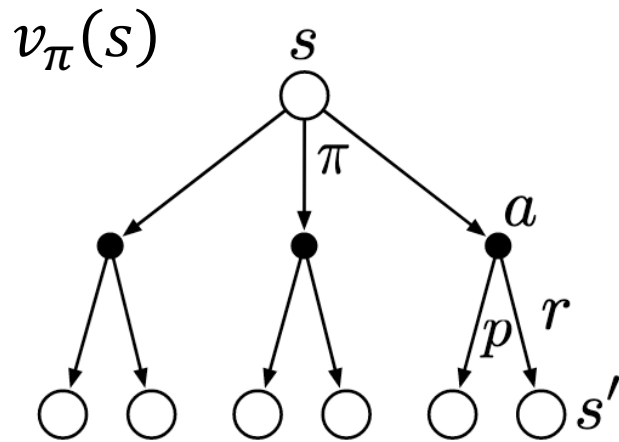


# Dynamic Programming: Applications

- We assume that we are given a model of the MDP that characterizes our problem
- While this assumption does not hold in many contexts, it does in many others
  - Organic chemistry
  - Puzzles
  - Quantum computing
  - Theorem proving
- Furthermore, it builds the foundation that we will use to explore model-free algorithms

# Bellman Equation

- $v_{\pi}(s) = \sum_a \pi(a|s)(r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\pi}(s'))$
- $q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) \sum_a \pi(a|s') q_{\pi}(s', a)$
- **Optimal substructure**: Can construct an optimal solution from optimal solutions of subproblems



# Outline

- Background
- Policy Evaluation
- Policy Improvement
- Policy Iteration
  - Modified Policy Iteration
  - Value Iteration
- Approximate value iteration

# Policy Evaluation

- Estimate the expected future reward when following policy  $\pi$
- From the Bellman equation we know that
  - $v_\pi(s) = \sum_a \pi(a|s)(r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_\pi(s'))$
- Given a policy  $\pi$ , what if we searched for a function  $V$  that satisfies the Bellman equation?
  - Will we have successfully evaluated  $\pi$ ?
- If so, how should we search for  $V$ ?
- Use the Bellman equation as an update rule
  - $V(s) = \sum_a \pi(a|s)(r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s'))$

# Policy Evaluation

---

## Algorithm 1 Policy Evaluation

---

```
1: procedure POLICY EVALUATION( $\mathcal{S}, V, \pi, \gamma$ )
2:    $\Delta \leftarrow \text{inf}$ 
3:   while  $\Delta > 0$  do
4:      $\Delta \leftarrow 0$ 
5:     for  $s \in \mathcal{S}$  do
6:        $v \leftarrow V(s)$ 
7:        $V(s) \leftarrow \sum_a \pi(a|s)(r(s, a) + \gamma \sum_{s'} p(s'|s, a)V(s'))$ 
8:        $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
9:     end for
10:  end while
11:  return  $V$ 
12: end procedure
```

▷  $v_\pi$

---

Is this guaranteed to converge to  $v_\pi$  if the threshold is 0?



# Policy Evaluation: Convergence

$$v_{\pi}(s) = \sum_a \pi(a|s) (r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\pi}(s'))$$

$$v_{\pi}(s) = \sum_a \pi(a|s) r(s, a) + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) v_{\pi}(s')$$

$$\mathbf{v}^{\pi} = \mathbf{r}^{\pi} + \gamma \mathbf{P}^{\pi} \mathbf{v}^{\pi}$$

- Matrices:

- $\mathbf{v}^{\pi}$  is a vector of values for each state, size  $|\mathcal{S}|$
- $\mathbf{r}^{\pi}$  is a vector of rewards for each state, size  $|\mathcal{S}|$
- $\mathbf{P}^{\pi}$  is a matrix of transition probabilities for each pair of states, size  $|\mathcal{S}| \times |\mathcal{S}|$

# Policy Evaluation: Convergence

- Define  $T$  as the Bellman backup operator:

$$T(\mathbf{v}) := \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{v}$$

- We know that there exists a fixed point

$$T(\mathbf{v}^\pi) = \mathbf{v}^\pi$$

- The infinity norm:

$$\|\mathbf{x}\|_\infty := \max_i |x_i|$$

- If:

$$\|T(\mathbf{u}) - T(\mathbf{v})\|_\infty \leq \gamma \|\mathbf{u} - \mathbf{v}\|_\infty$$

for  $\gamma \in [0,1)$

- Then  $T$  is a **contraction mapping**
- Banach-Caccioppoli fixed point theorem proves that repeated applications of  $T$  will converge to a **unique** fixed point (i.e.  $\mathbf{v}^\pi$ ).

# Policy Evaluation: Convergence

$$\begin{aligned} \|T(\mathbf{u}) - T(\mathbf{v})\|_{\infty} &= \|\mathbf{r}^{\pi} + \gamma \mathbf{P}^{\pi} \mathbf{u} - \mathbf{r}^{\pi} + \gamma \mathbf{P}^{\pi} \mathbf{v}\|_{\infty} \\ &= \|\gamma \mathbf{P}^{\pi} \mathbf{u} + \gamma \mathbf{P}^{\pi} \mathbf{v}\|_{\infty} \\ &= \gamma \|\mathbf{P}^{\pi} (\mathbf{u} + \mathbf{v})\|_{\infty} \\ &\leq \gamma \left\| \left\| \mathbf{P}^{\pi} \|\mathbf{u} + \mathbf{v}\|_{\infty} \right\|_{\infty} \right\|_{\infty} // \mathbf{P}^{\pi} \text{ is a matrix of transition probabilities, sums to 1} \\ &\leq \gamma \|\mathbf{u} + \mathbf{v}\|_{\infty} \end{aligned}$$

# Policy Evaluation AI Farm

- Policy: Uniform random
- AI Farm took 4576 iterations to converge with  $\gamma = 1$ 
  - $\gamma = 0.99$  takes 1345 iterations
  - $\gamma = 0.9$  takes 191 iterations
- What if policy said to always go up?
  - $\gamma = 1$
  - $\gamma < 1$

-2274.06	-2246.84	-2197.41	-2135.24	-2069.99	-2009.64	-1959.49	-1918.20	-1857.38	-1820.67
-2297.28	-2265.07	-2206.14	-2134.31	-2061.11	-1995.42	-1946.64		-1829.27	-1779.97
-2348.70	-2306.01	-2223.76	-2130.77	-2040.71	-1960.30	-1893.92		-1742.02	-1685.95
		-2248.12	-2120.29	-2006.68	-1907.14	-1823.96		-1608.05	-1531.88
-2385.21	-2337.15		-2091.58	-1954.59	-1833.61		-1541.14		-1297.63
		-2170.76	-2025.55	-1882.47	-1731.44	-1550.44	-1268.95		-956.19
-2226.86	-2176.25	-2076.61	-1953.40	-1814.29	-1655.26				-610.45
-2124.66	-2082.58	-2002.04	-1893.14	-1762.04		-1410.82			
-2060.53	-2023.37	-1951.82	-1851.10	-1725.48	-1576.81	-1402.05	-1211.91	-1032.18	-985.69
-2029.55	-1994.57	-1926.78	-1829.95	-1707.97		-1404.66	-1243.45	-1109.79	-1049.74

# Outline

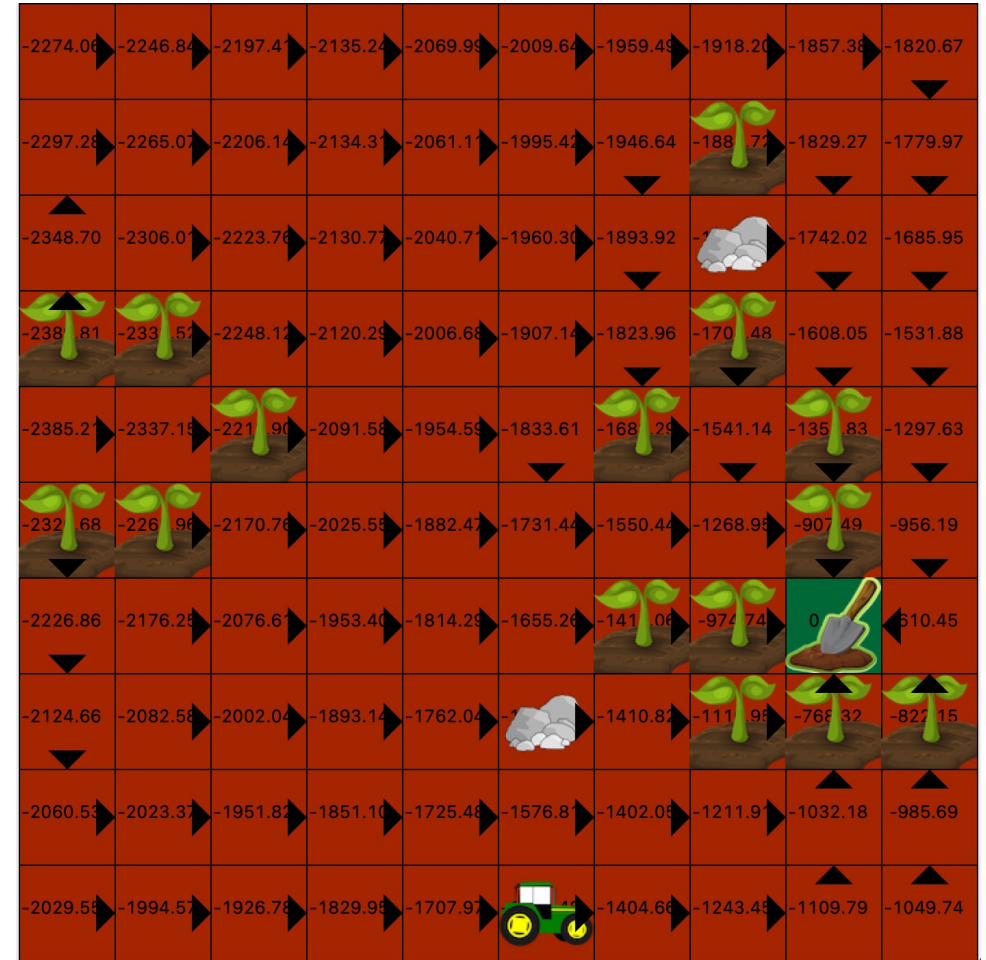
- Background
- Policy Evaluation
- Policy Improvement
- Policy Iteration
  - Modified Policy Iteration
  - Value Iteration
- Approximate value iteration

# Policy Improvement

- We have evaluated policy  $\pi$ , how can we find a better policy?
- $\pi' \geq \pi$  if and only if  $v_{\pi'}(s) \geq v_{\pi}(s)$  for all  $s \in \mathcal{S}$
- We set the policy to be greedy with respect to  $v_{\pi}$
- $\pi'(s) = \underset{a}{\operatorname{argmax}}(r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\pi}(s'))$
- This policy will always be the same or better than the previous one
  - Policy improvement theorem.

# Policy Improvement

- Better than original uniform random policy
- Still not optimal



When acting greedily with respect to  $v_{\pi}$

# Outline

- Background
- Policy Evaluation
- Policy Improvement
- **Policy Iteration**
  - Modified Policy Iteration
  - Value Iteration
- Approximate value iteration



# Policy Iteration

- Policy improvement improves a policy  $\pi$  and obtains a new policy  $\pi'$  such that,  $\pi' \geq \pi$
- If  $\pi' = \pi$ , then  $v_{\pi'} = v_{\pi}$
- Therefore,  $v_{\pi'}(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\pi'}(s'))$
- This is the same as the Bellman optimality equation!
- Does this mean  $\pi' = \pi_*$  and  $v_{\pi'} = v_*$ ?
- There is a proof showing that the Bellman optimality equation is a unique fixed point
  - Similar to that of the Bellman equation

# Policy Iteration

---

## Algorithm Policy Iteration

---

1: **Inputs:**

2:  $\pi$ : Initial Policy

3:  $V$ : Initial value function. Value of terminal state must be zero

4:  $\theta$ : Termination threshold

5:

6: **while** has not converged **do**

7:      $V = \text{Policy\_Evaluation}(\pi, V, \theta)$

8:      $\pi = \text{Policy\_Improvement}(V)$

9: **end while**

10: **return**  $\pi$

▷ Approximation of  $\pi^*$

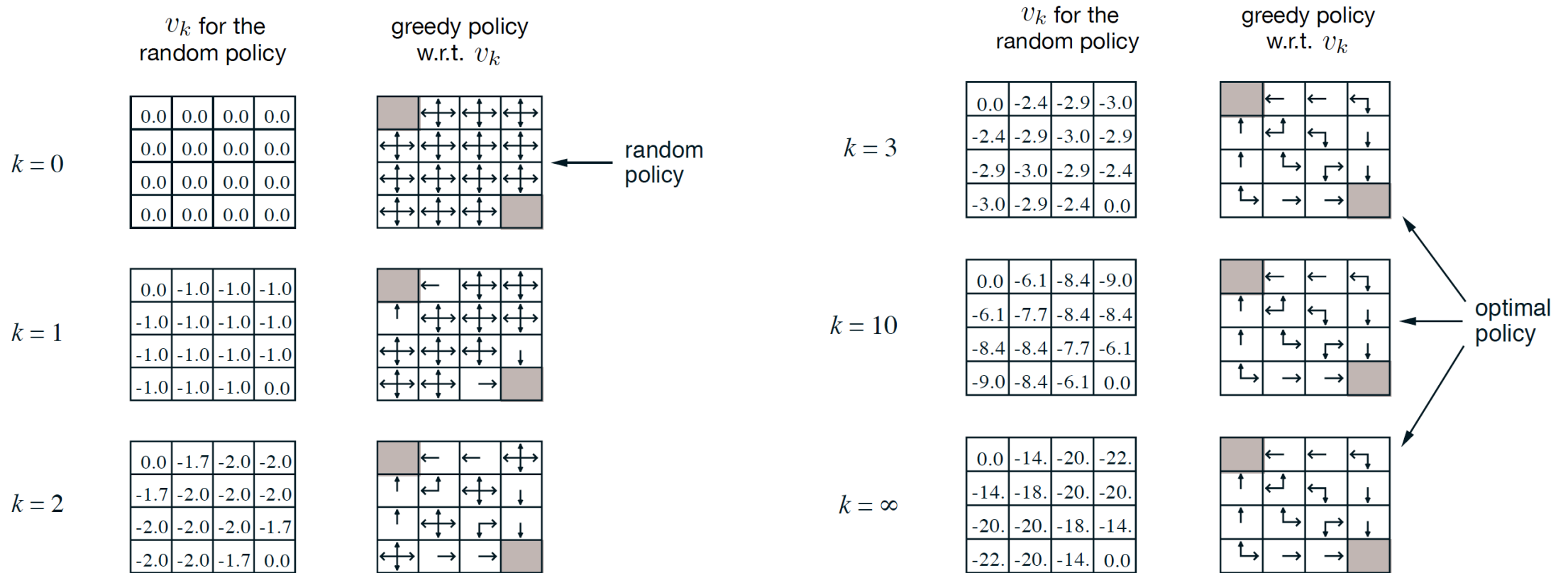
---

# Policy Iteration: AI Farm



```
(drl) forestagostinelli@Forests-MacBook-Pro Farm_Grid_World % python run_policy_iteration.py --map maps/map1.txt --wait 1.0 --wait_evaluation 0.0 --rand_right 0.0
```

# Do We Have to Wait For Policy Evaluation to Converge?

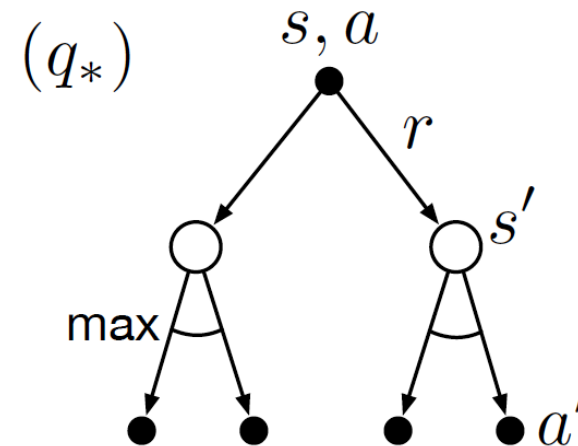
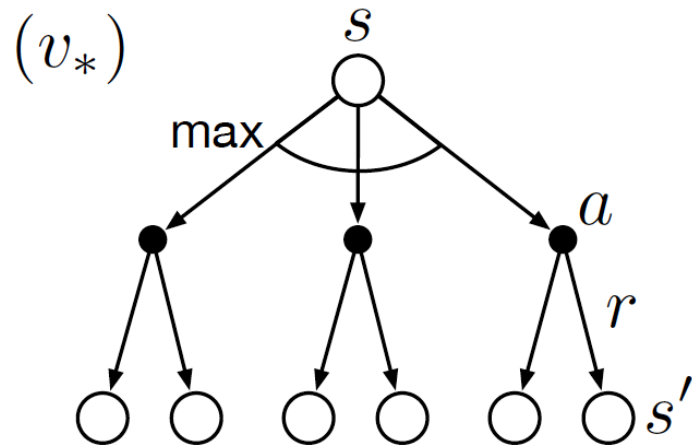


# Modified Policy Iteration

- We do not have to run policy evaluation to convergence.
- We can stop after reaching some threshold
- We can stop after  $k$  iterations
  - $k = 1$  is the same as **value iteration**

# Bellman Optimality Equation

- $v_*(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_*(s'))$
- $q_*(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_{a'} q_*(s', a')$



# Value Iteration

- Find the optimal value function
- Recall the Bellman optimality equation
  - $v_*(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_*(s'))$
- Use this as an update rule
  - $V(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s'))$
- Combines policy evaluation and policy improvement into one step
- There is a proof similar to that of policy evaluation to prove that repeatedly updating  $V$  will converge to a unique fixed point
  - Bellman optimality equation

# Value Iteration

---

## Algorithm Value Iteration

---

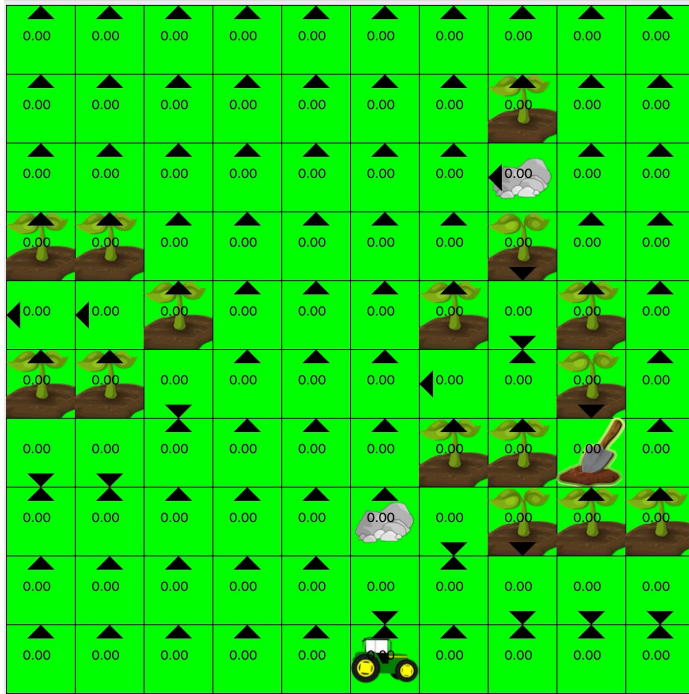
```
1: Inputs:  
2:    $V$ : Initial value function. Value of terminal state must be zero  
3:    $\theta$ : Termination threshold  
4:  
5: procedure VALUE_ITERATION( $V, \theta, \gamma$ )  
6:    $\Delta \leftarrow \text{inf}$   
7:   while  $\Delta > \theta$  do  
8:      $\Delta \leftarrow 0$   
9:     for  $s \in \mathcal{S}$  do  
10:       $v \leftarrow V(s)$   
11:       $V(s) \leftarrow \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a)V(s'))$   
12:       $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
13:     end for  
14:   end while  
15:   return  $V$  ▷ Approximation of  $v_*$   
16: end procedure
```

---



# Value Iteration: AI Farm

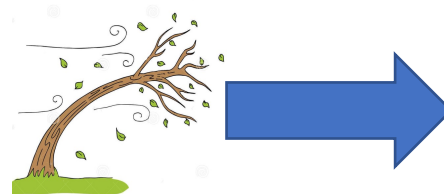
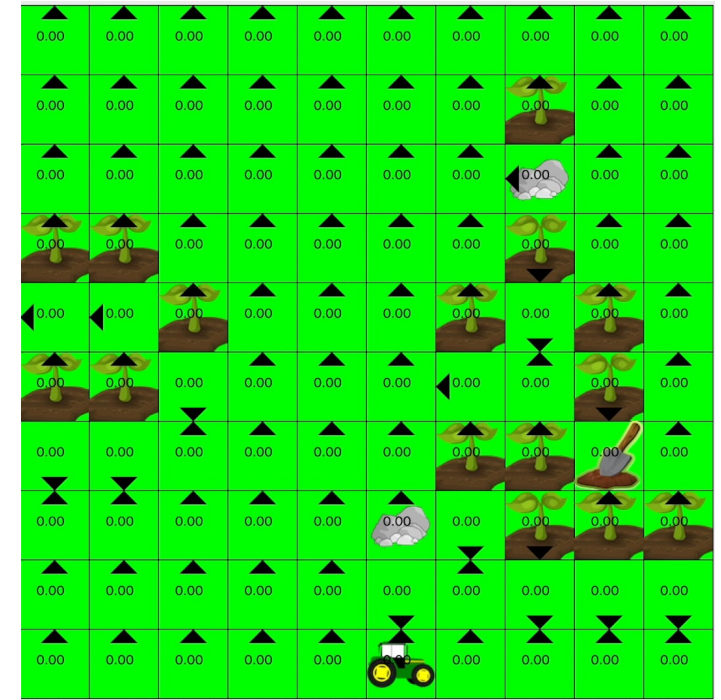
$p = 0$



$p = 0.1$



$p = 0.5$



# Asynchronous Policy Iteration

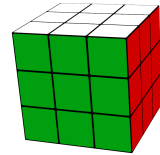
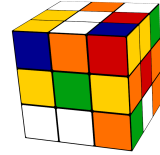
- We do not have to update all states at every iteration
- Prioritize states based on Bellman error
  - Difference between current and updated estimation of state value
- Prioritize states based on relevance to the agent's experience
- To converge, must continue to update all states, though it does not have to be at every iteration

# Outline

- Background
- Policy Evaluation
- Policy Improvement
- Policy Iteration
  - Modified Policy Iteration
  - Value Iteration
- **Approximate value iteration**

# Approximate Value Iteration

- Sometimes, we can accurately define the MDP, however, the state space is too large to store in a table
- Therefore, we will need to use an architecture to approximate the value
- The goal is to have an architecture whose number of parameters is far less than the size of the state space

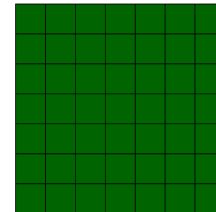
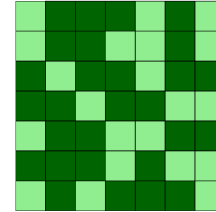


Rubik's cube

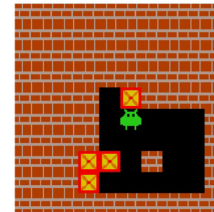
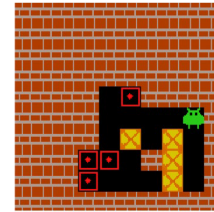
22	12	4	2	5
17	16	3	6	9
20	19	18	11	7
23	1		24	13
21	14	10	8	15

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	

24 puzzle

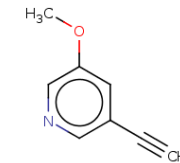
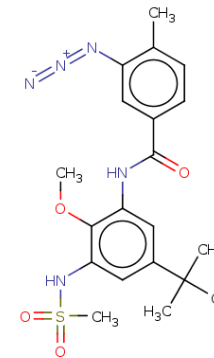


Lights Out (7x7)

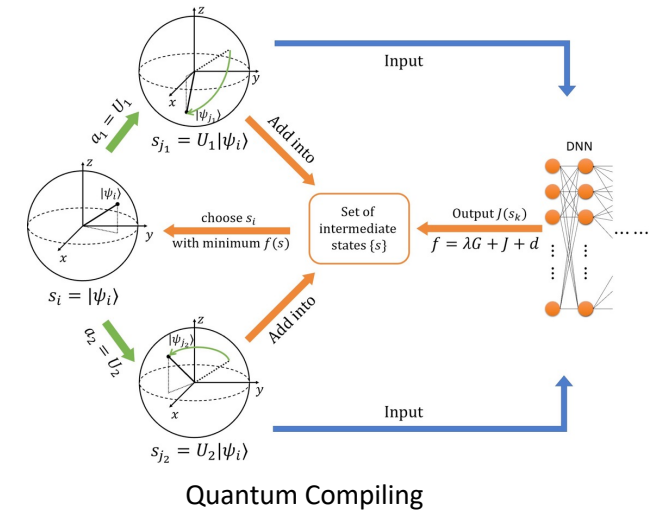


Sokoban

Puzzles



Chemical Synthesis



# Approximate Value Iteration

- Value Iteration

- $V(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s'))$

- Approximate Value Iteration

- $y = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) \hat{v}(s', \mathbf{w}))$

- $E(\mathbf{w}) = \frac{1}{2} (y - \hat{v}(s, \mathbf{w}))^2$

- $\nabla_{\mathbf{w}} E(\mathbf{w}) = (y - \hat{v}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w})$

- Even though  $y$  depends on  $\mathbf{w}$ , we do not differentiate  $y$  with respect to  $\mathbf{w}$ .

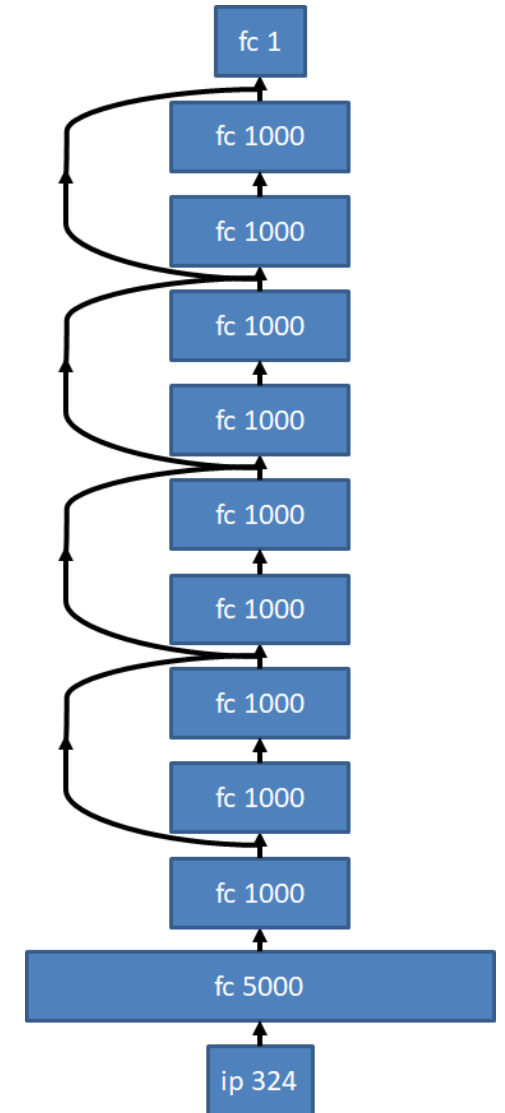
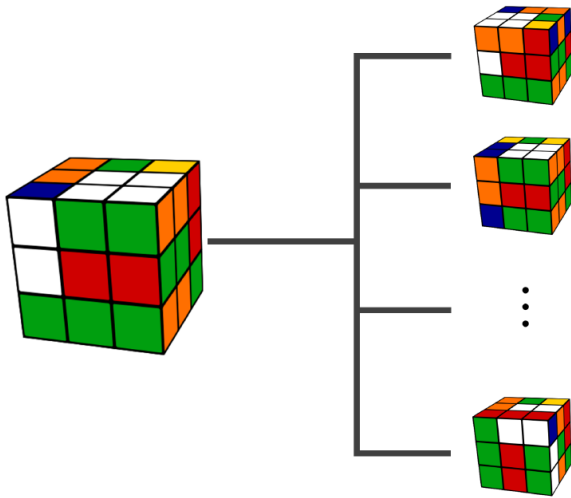
- Only supervision is that  $y = 0$  for terminal states

# Convergence

- Unlike in the tabular case, we cannot guarantee convergence in the case of non-linear function approximators
- However, there are many different methods that we can use to make deep reinforcement learning work in practice

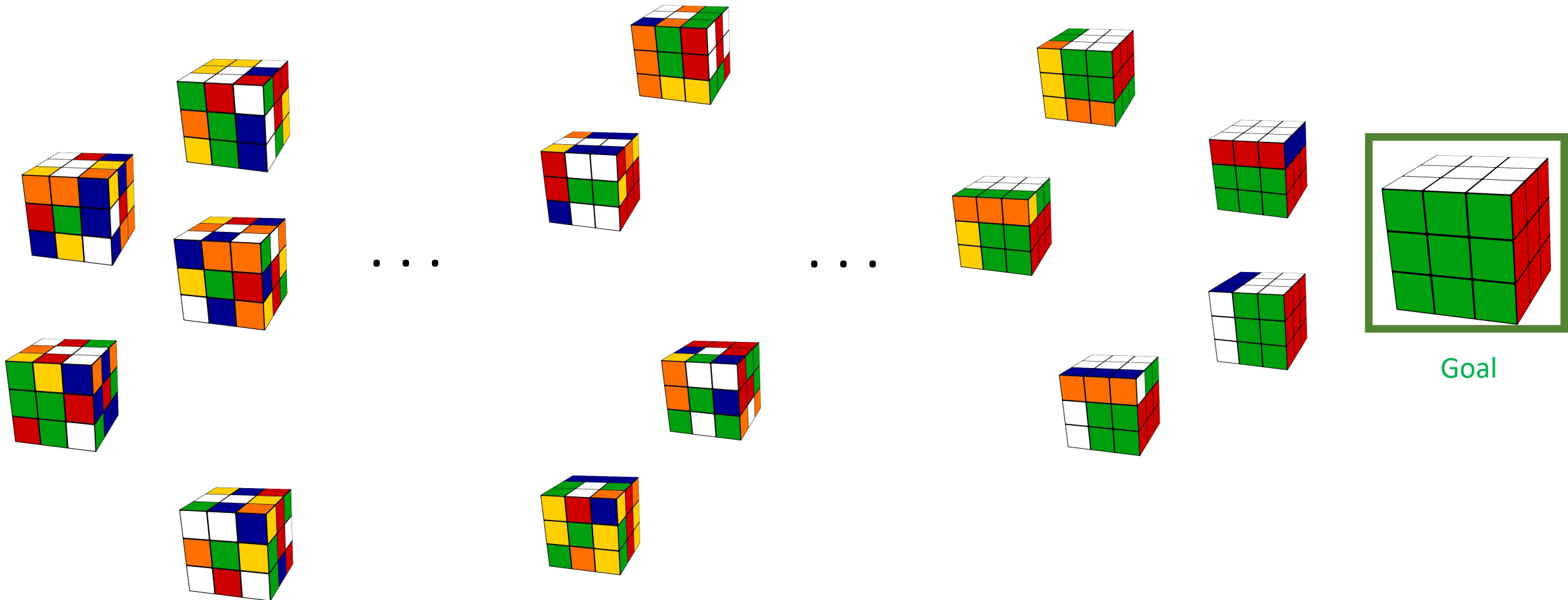
# Deep Approximate Value Iteration

- $y = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) \hat{v}(s', \mathbf{w}))$
- $E(\mathbf{w}) = \frac{1}{2} (y - \hat{v}(s, \mathbf{w}))^2$



# Prioritized Sweeping

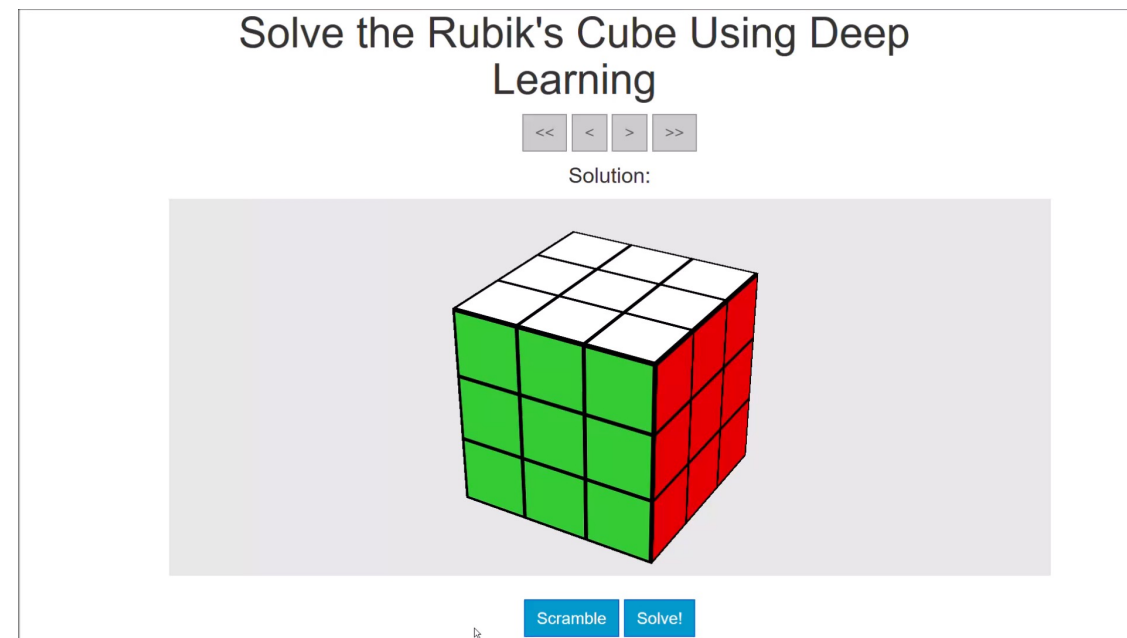
- Prioritized sweeping: Generate training data by taking moves in reverse from the **goal**



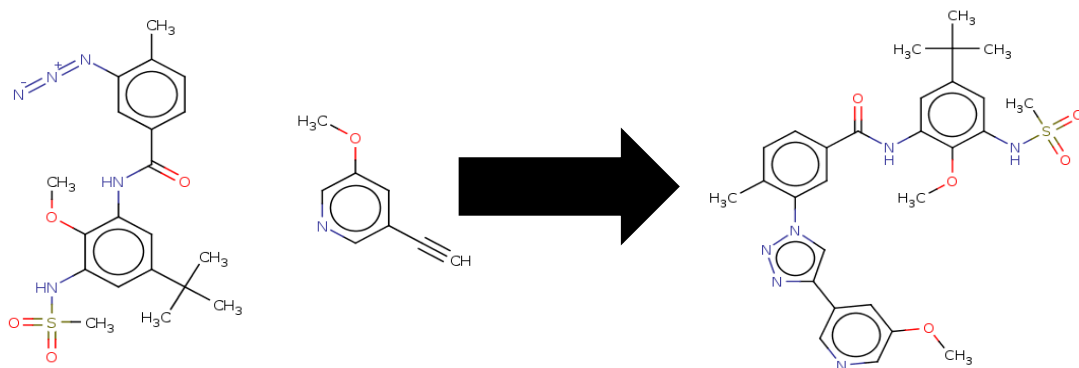


# Solving the Rubik's Cube

- Deep approximate value iteration and A\* search
- Solves the Rubik's cube and 6 other puzzles
- Can be applied to other areas in the natural sciences



<http://deepcube.igb.uci.edu/>



# Summary

- **Policy Evaluation:** Uses Bellman equation as an update rule

$$V(s) = \sum_a \pi(a|s)(r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s'))$$

- **Policy Improvement:** Behave greedily with respect to value function

$$\pi'(s) = \operatorname{argmax}_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s'))$$

- **Policy Iteration:** Iterate between policy evaluation and policy improvement until convergence

- **Value Iteration:** Uses Bellman optimality equation as an update rule

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