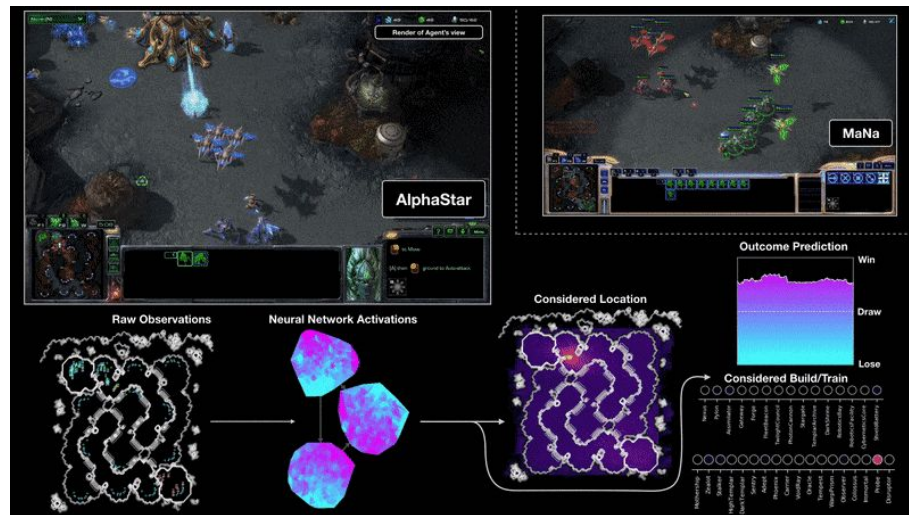
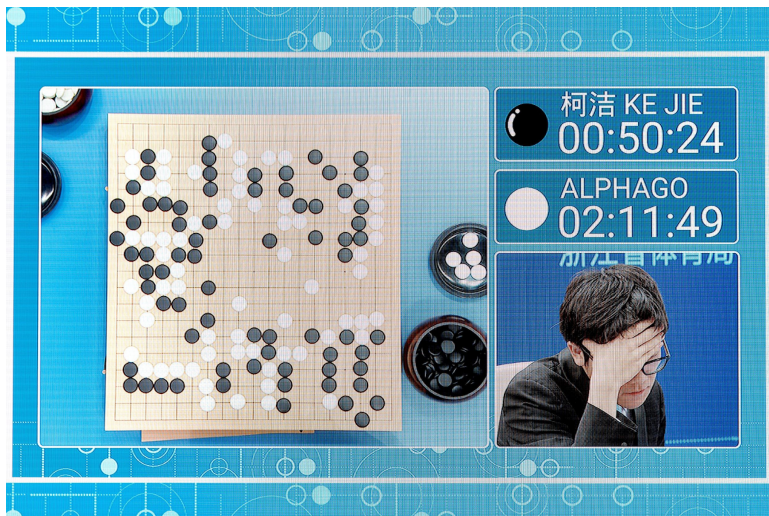




Introduction to Reinforcement Learning

Simon Fan
Ashish Gaurav
UWDSC W19





Brief timeline

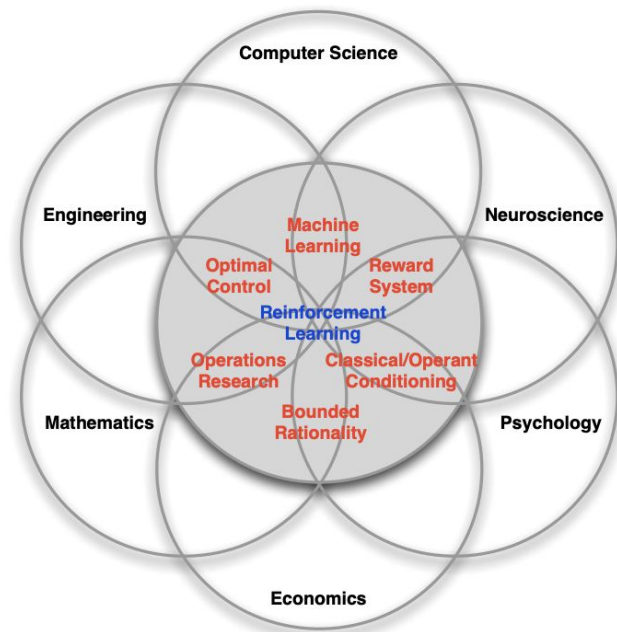
- 1963: RL plays tic-tac-toe
- 1970: AI winter
- 1989: Q-learning
- 2010: Deep Learning
- 2016: AlphaGo (supervised learning + deep RL + monte carlo + tricks)
- 2017: AlphaGo Zero (deep RL + monte carlo + more tricks)
- 2019: AlphaGo Star (supervised learning + deep RL + monte carlo + more and more tricks)

What is RL?

Train an **agent**, who interacts with an **environment** by performing **actions** which result in **rewards**.

A **policy** (which action to perform and when) is learnt by the agent and we can compute the cumulative reward the agent receives by following the policy.

The goal of RL is for the agent to learn the **optimal policy** to follow, in order to maximize the cumulative reward (value).





Contrast with other forms of learning

Suppose there is a distribution X with true labels y , we approximate the label mapping function with:

- Supervised learning
 - Given m samples of $\langle x, y \rangle$, assume correct
 - find $f: x \rightarrow y$ w.r.t. some objective function
- Unsupervised learning
 - Given m samples of x
 - find $f: x \rightarrow y$
- Reinforcement learning
 - Given m samples of x and some reward function z , assume correct
 - learn $f: x \rightarrow y$ w.r.t. the reward function z , i.e. maximize the cumulative reward
 - z looks a lot like y in supervised learning

Parallel with Psychology

Positive reinforcement (positive reward)

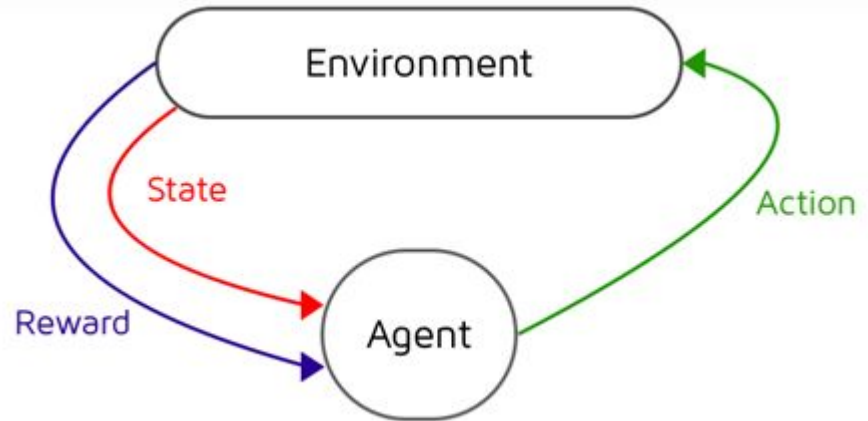


Negative reinforcement (negative reward)



RL definitions

- Agent: the one who interacts with the environment
- State (s): unique, observed status of the environment
- Possible actions (A): state to actions mapping
- Action ($a=A(s)$): action that the agent selected
- Reward (R): $\langle \text{state, action} \rangle$ to reward mapping
- Policy (π): state to action mapping

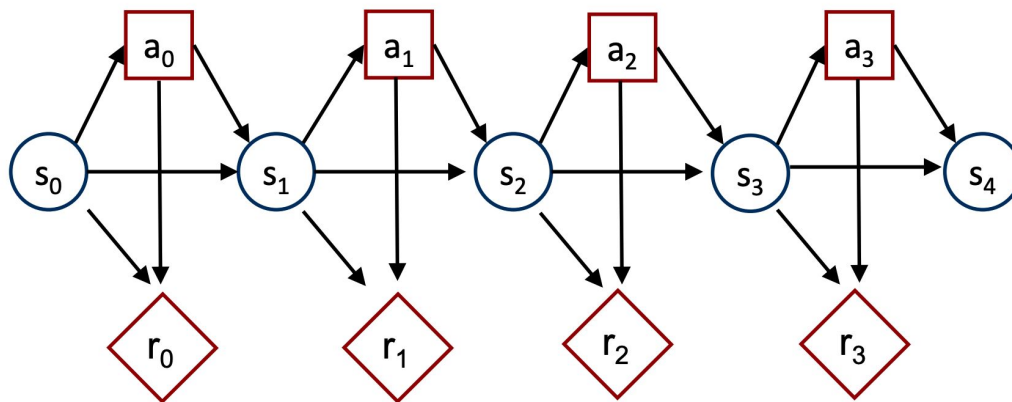


RL definitions

Formally it's a MDP

Markov Decision Process

- Markov process augmented with...
 - Actions e.g., a_t
 - Rewards e.g., r_t



An example: walking on a lake

- States: Each index is a state
 - S: start state
 - F: walkable surface
 - H: game ends
 - G: goal state
- Reward: 1 when goal is reached, 0 otherwise
- Agent: Tries to get from S to G such that reward is maximized
- Actions: Up, Down, Left, Right

What is the optimal policy? Are there more than one?

S	F	F	F
F	H	F	H
F	F	F	H
H	F	F	G

An example: Frozen Lake

- States: Each index is a state
 - S: start state
 - F: walkable surface, *but slippery*
 - H: game ends
 - G: goal state
- Reward: 1 when goal is reached, 0 otherwise
- Agent: Tries to get from S to G such that reward is maximized
- Actions: Up, Down, Left, Right

A lot harder to write a non RL solver

DEMO IN OPENAI GYM! <https://gym.openai.com/>

S	F	F	F
F	H	F	H
F	F	F	H
H	F	F	G

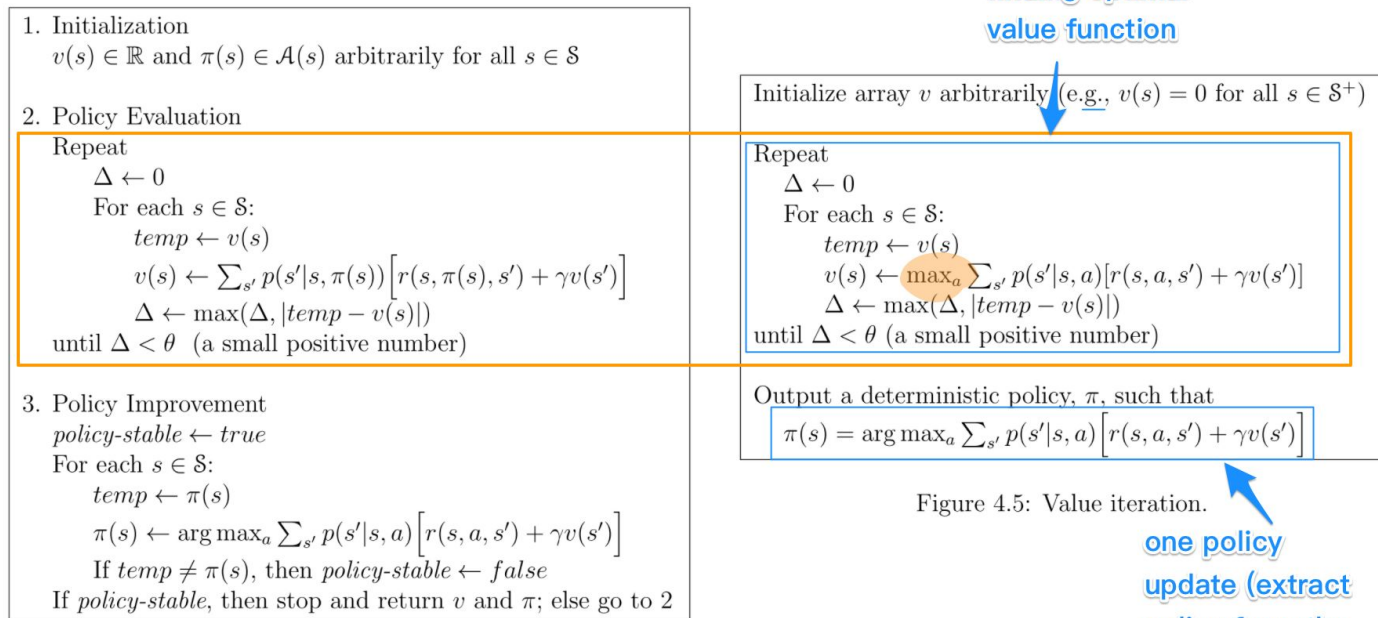


Figure 4.5: Value iteration.

Figure 4.3: Policy iteration (using iterative policy evaluation) for v_* . This algorithm has a subtle bug, in that it may never terminate if the policy continually switches between two or more policies that are equally good. The bug can be fixed by adding additional flags, but it makes the pseudocode so ugly that it is not worth it. :-)



Uncovered, but important topics

- Continuous state spaces
- Continuous action spaces
- Model vs Model-free
- On-policy vs off-policy
- Exploration vs exploitation tradeoff
 - Superstitions
 - “Solution”: epsilon-greedy, upper confidence bounds, ...
- Reward tuning
 - Horizon
 - Credit assignment
 - “Solutions”: shaping rewards, discount rate, ... (careful environment design)