
Practical Machine Learning

Reinforcement Learning

<https://bit.ly/2USYqjZ>

Today we cover

Reinforcement learning

Open AI Gym

Building an agent

Reinforcement Learning

Learning to act in complex environments

Reinforcement learning (RL) is an area of machine learning concerned with how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward.

internal state



reward

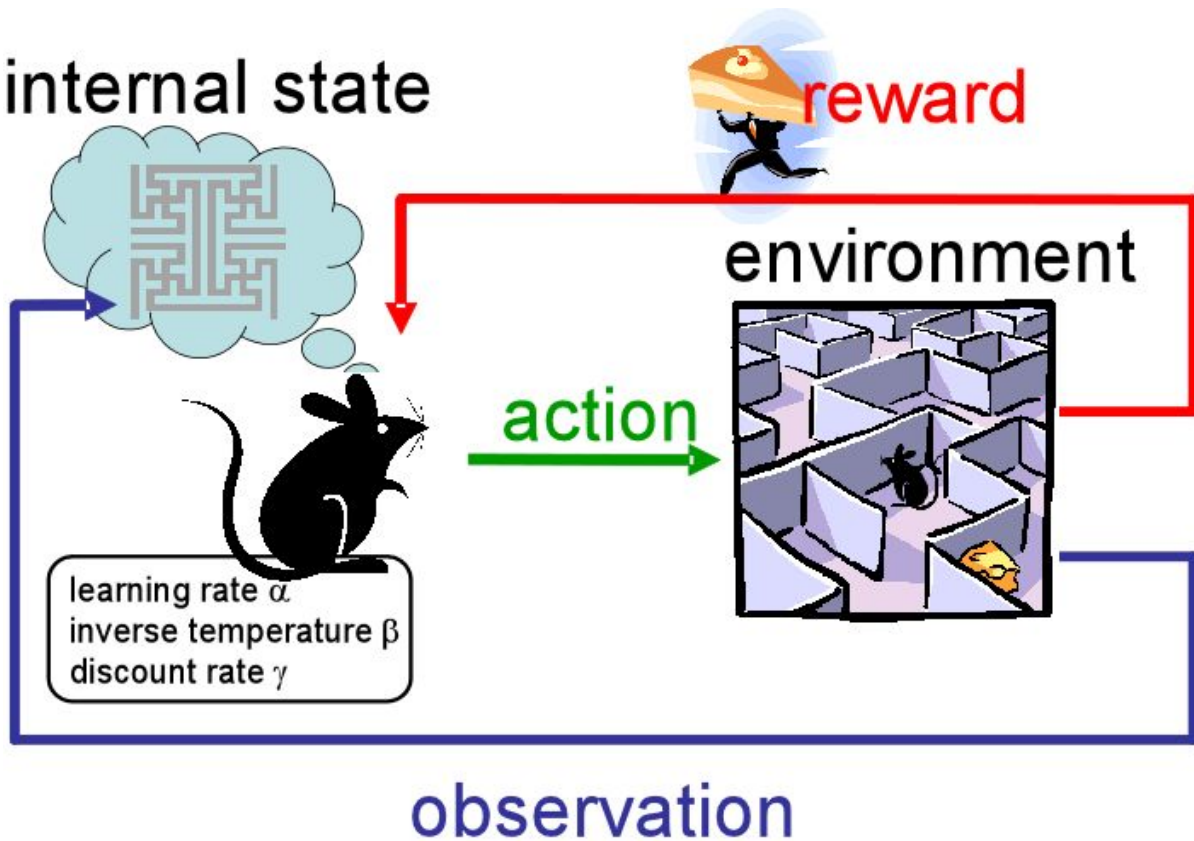
environment

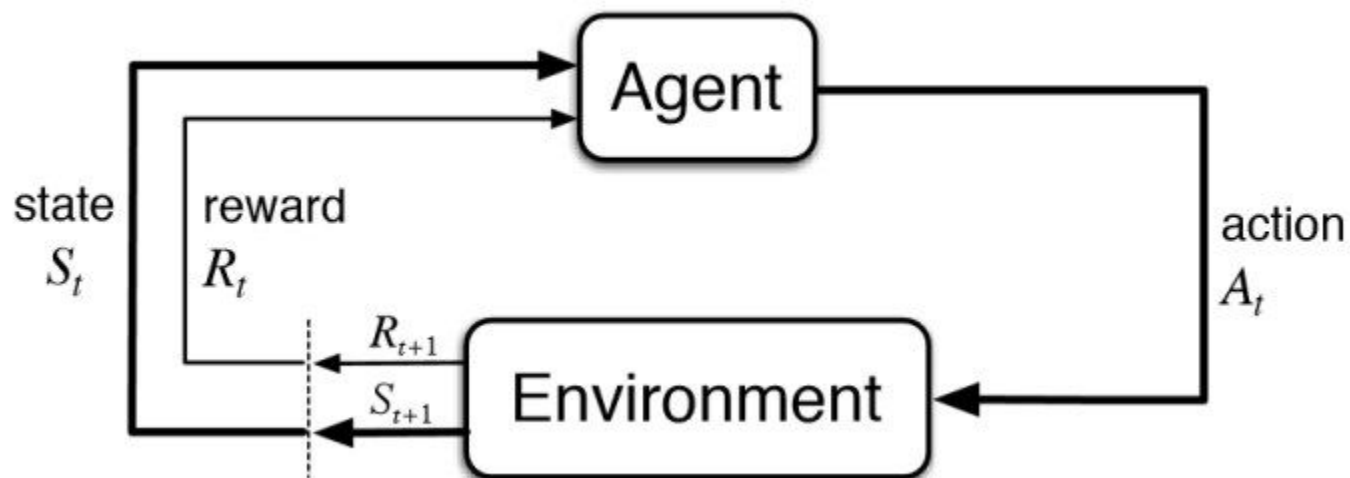


action

learning rate α
inverse temperature β
discount rate γ

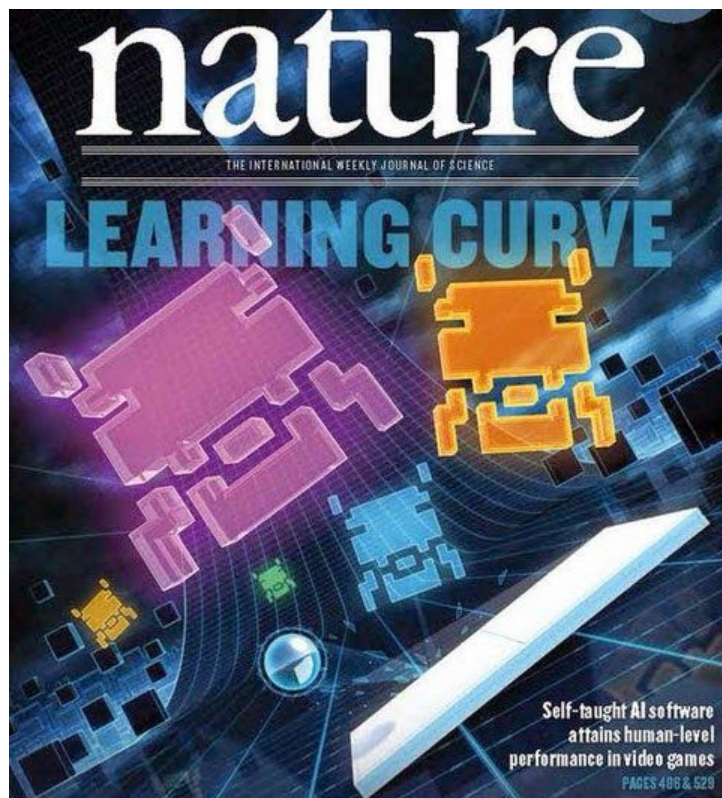
observation





Why is it difficult?

- No immediate label
 - Reward could be from a combination of actions
 - I need to explore the space
 - I need to find what was useful from my actions
-





nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

At last – a computer program that
can beat a champion Go player **PAGE 484**

ALL SYSTEMS GO



RL policy network

Value network



v_θ



Self-play positions

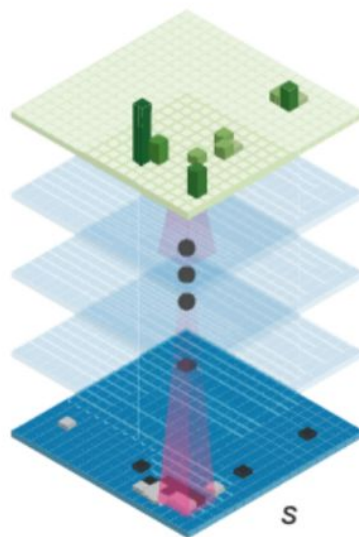
Neural network

Data

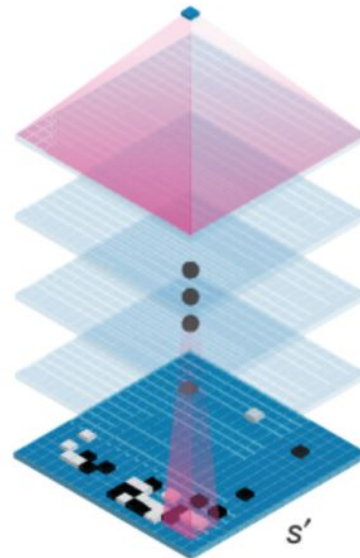
Policy network

Value network

$p_{\sigma|p}(a|s)$



$v_\theta(s')$






Key concepts

- Environment
 - Agent
 - State / observation
 - Action
 - Reward
 - Discount factor
 - Policy
-


Environment

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

Agent


	F	F	F
F	H	F	H
F	F	F	H
H	F	F	G

Action


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

Left / Right / Up / Down

State

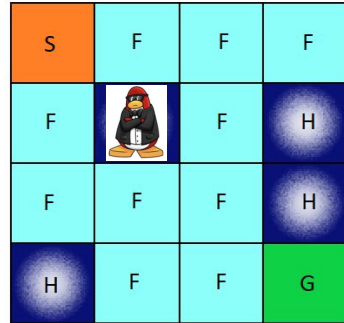
S	F		F
F	H	F	H
F	F	F	H
H	F	F	G

State = 2

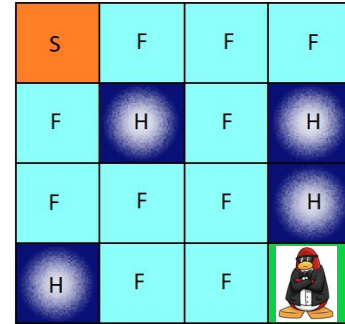
S	F	F	F
F	H	F	H
	F	F	H
H	F	F	G

State = 8

Reward



0 if falls



1 if goal reached

Return

$$R = \sum_{t=0}^{\infty} \gamma^t r_t,$$

Reward
at time t

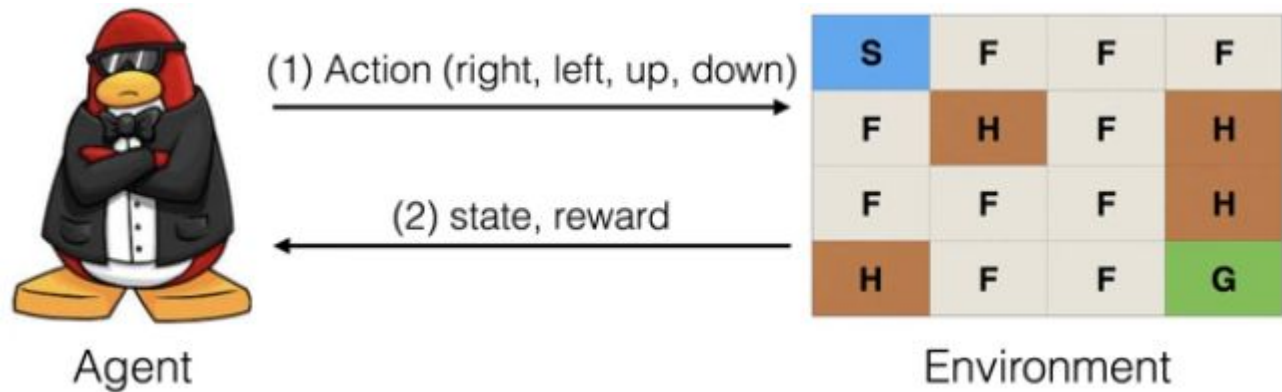
$$\gamma \in [0, 1]$$

Discount rate
modulated by time

Discount rate expresses my preference for immediate rewards.



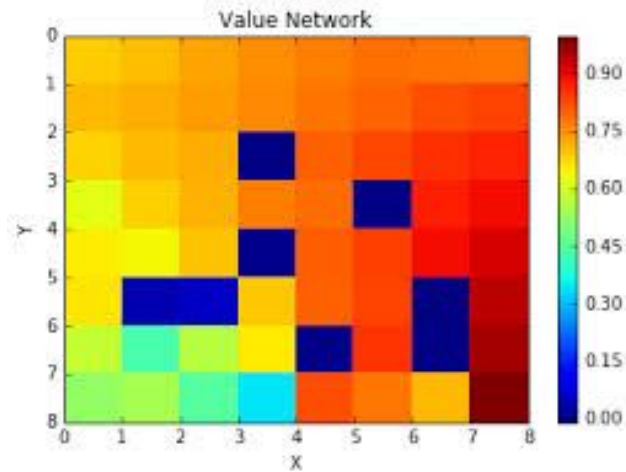
Policy



$$\pi : S \times A \rightarrow [0, 1]$$

$$\pi(a|s) = Pr(a_t = a | s_t = s)$$

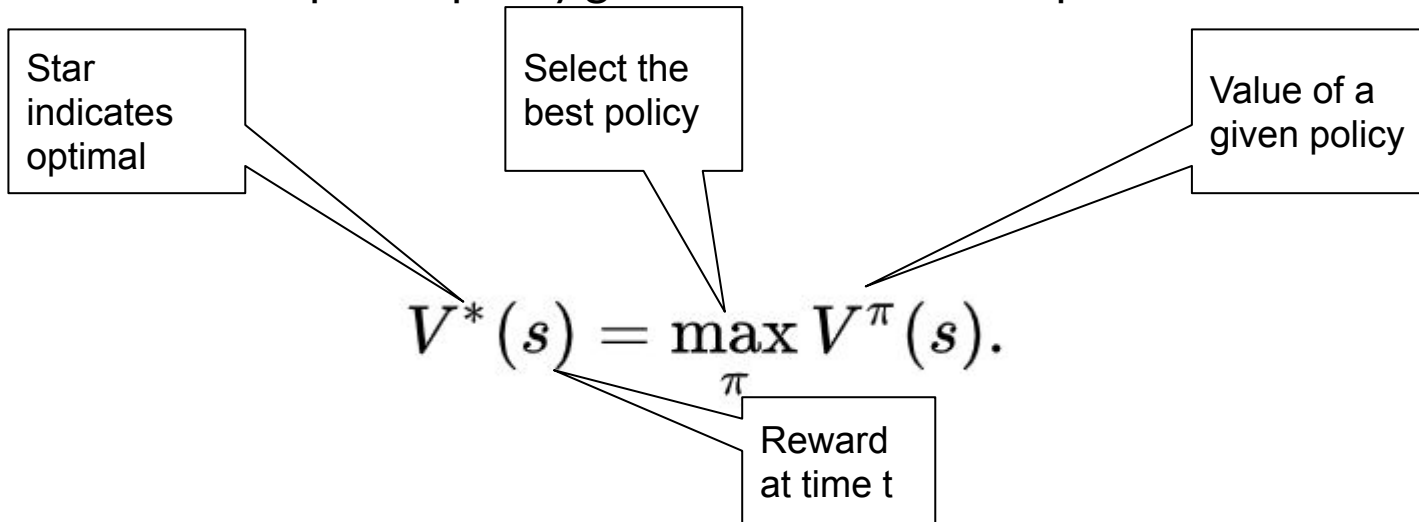
Value of a policy



$$V_{\pi}(s) = E[R] = E[\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s],$$

Optimal policy

The optimal policy guarantees the best expected returns

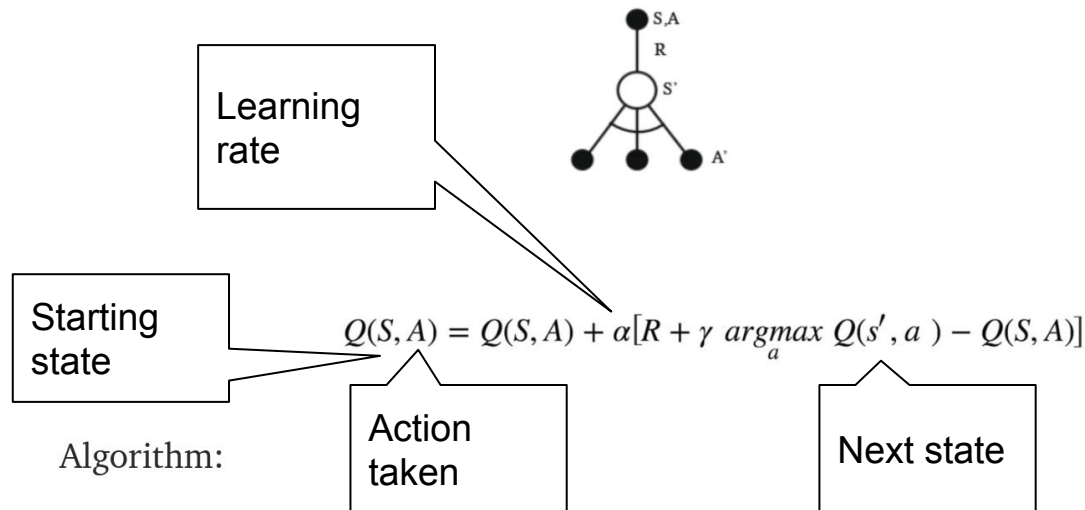


Q-Value

$$Q : S \times A \rightarrow \mathbb{R} .$$

$$Q^{\pi}(s, a) = E[R|s, a, \pi],$$

Q-learning (SARSA Max)



Initialize $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize S

Repeat (for each step of episode):

Choose A from S using policy derived from Q (e.g., ϵ -greedy)

Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$;

until S is terminal

Code

- Gym API
 - Env
 - Step
 - Observation
 - Reward
 - IsDone
-

```
[ ] !pip install gym
```

```
[ ] import gym
```

```
▶ env = gym.make('FrozenLake-v0')  
  obs = env.reset()  
  for _ in range(10):  
      env.render()  
      obs, reward, is_done, info = env.step(0)  
  env.close()
```
