DSC291: Machine Learning with Few Labels

Reinforcement Learning

Zhiting Hu Lecture 13, May 13, 2025



HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

• Reinforcement Learning

- Paper presentation:
 - David Lurie, Ben TenWolde: "Mixture of Agents with LLMs"
 - Qing Cheng, Zhaomei Geng: "ChatGPT in medicine: an overview of its applications, advantages, limitations, future prospects, and ethical considerations"



So far... Supervised Learning

Data: (x, y) x is data, y is label

Goal: Learn a *function* to map x -> y

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.





Classification

So far... Unsupervised Learning

Data: x no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



Today: Reinforcement Learning

Problems involving an **agent** interacting with an **environment**, State s, which provides numeric **reward** signals

Goal: Learn how to take actions in order to maximize reward

Superviser

angen 1



Agent

Environment

Action a

Reward r.

Next state s

Atari games figure copyright Volodymyr Mnih et al., 2013. Reproduced with permission.

Overview

- What is Reinforcement Learning? - Markov Decision Processes - Q-Learning - Value based RL & Af-polyc RL (- Policy Gradients - policy-based RL (on policy RL (-



Environment









Robot Locomotion



Objective: Make the robot move forward State: Angle and position of the joints Action: Torque applied on joints Reward: 1 at each time step upright + forward movement

Atari Games



Objective: Complete the game with the highest score

State: Raw pixel inputs of the game state **Action:** Game controls e.g. Left, Right, Up, Down **Reward:** Score increase/decrease at each time step

Go





Action: Where to put the next piece down Reward: 1 if win at the end of the game, 0 otherwise

-learny

dense renterf:

How can we mathematically formalize the RL problem?



Markov Decision Process

- Mathematical formulation of the RL problem
- Markov property: Current state completely characterises the state of the Stor Stor St

Defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

world

 $\rightarrow S$: set of possible states $\rightarrow \mathcal{A}$: set of possible actions $\rightarrow \mathcal{R}$: distribution of reward given (state, action) pair $\mathbb P$: transition probability i.e. distribution over next state given (state, action) pair γ : discount factor Styr M(Styl St, Ge)

Markov Decision Process

- At time step t=0, environment samples initial state $s_0 \sim p(s_0)$
- - Then, for t=0 until done: \rightarrow Agent selects action $a_t \oint a_t \sim \pi (a_t / S_t)$
 - Environment samples reward $r_t \sim R(. | s_t, a_t)$
 - Environment samples next state $s_{t+1} \sim P(. | s_t, a_t)$
 - Agent receives reward r_t and next state s_{t+1}

a~ R (a/S)

- A policy π is a function from S to A that specifies what action to take in each state
- Following a policy produces sample *trajectories* (or paths) $s_0, a_0, (r_0, s_1, a_1, (r_1), ...$

Objective: find policy π^* that maximizes cumulative discounted reward:



A simple MDP: Grid World

 A_{i} actions = {

1. right →

2. left ----

3. up
4. down





Objective: reach one of terminal states (greyed out) in least number of actions

 $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

A simple MDP: Grid World



Random Policy

 $\pi(a|s)$



Optimal Policy





Questions?