Reinforcement Learning: Overview

Sargur N. Srihari srihari@cedar.buffalo.edu

Topics in Reinforcement Learning

- 1. RL as a topic in Machine Learning
- 2. Tasks performed by reinforcement learning
- 3. Policies with exploration and exploitation
- 4. RL connected to a deep neural net
 - Deep RL for playing the Atari game

Srihari

Reinforcement Learning Srihari Association Strihari



Three Types of Machine Learning

- 1. Supervised Learning (Predictive)
 - Learn mapping $y(\mathbf{x})$ given dataset $D = \{(\mathbf{x}_i, y_i)\}$
 - E.g., MNIST classification
- 2. Unsupervised Learning (Descriptive)
 - Given only inputs $D = \{x_i\}$ find interesting patterns
 - E.g., Determine *k* cluster centers
- 3. Reinforcement Learning
 - How to act or behave when given occasional reward or punishment signals
 - E.g., how a robot learns to walk to a power outlet

Analogy of teaching a dog

- Consider teaching a dog a new trick:
 - You cannot tell it what to do, but you can reward/punish it if it does the right/wrong thing
 - It has to figure out what it did that made it get the reward/punishment, which is known as the credit assignment problem

5

Applications of Reinforcement Learning

- Use similar method to train computers for
- Game playing (backgammon, chess, GO)
- Scheduling jobs
- Robots (in a maze, controlling robot limbs)
- Multiple agents, Partial observability
- RL system based on deep learning (Deepmind)
 - Play Atari video games
 - Robotics
 - Reaching human level performance on many tasks

A robot game

- Goal of robot: get reward of diamond and avoid the hurdles (fire)
- Robot learns by trying all
- possible paths and choosing path
- which gives reward with the least
- hurdles



- Each right step will give robot a reward and each wrong step will subtract the reward
- Total reward is calculated when it reaches the final reward that is the diamond



https://www.youtube.com/watch?v=gn4nRCC9TwQ



Task of Reinforcement learning

- Autonomous agent must learn to perform a task by trial and error without any guidance from the human operator
- Reinforcement learning is the problem of getting an agent to act in the world so as to maximize its rewards

Srihari

Data Sets

- Unlike supervised and unsupervised learning, reinforcement learning does not just experience a fixed data set
- Reinforcement learning algorithms interact with an environment
 - Q-learning generates data exclusively from experience, without incorporation of the prior knowledge.
 - If we put all our history data into a table with *state, action, reward, next state* and then sample from it, it should be possible to train our agent that way, without the dataset.

- Agent (algorithm) interacts with its environment
- A feedback loop between agent (a system) and its experience (in the environment)



- A mobile robot has actions (move forward, turn).
 - Its task is to learn a control strategy or policy for choosing actions that achieve its goals

- E.g., goal of docking onto battery charger when battery is low

Srihari

The Learning Task

- Agent exists in environment with set of states S
 - It can perform any of a set of actions A
 - Performing action A_t in state S_t receives reward R_t
- Agent's task is to learn control policy $\pi: S \rightarrow A$
 - That maximizes expected sum of rewards
 - with future rewards discounted exponentially



Goal: Learn to choose actions that maximize

$$R_0 + \gamma R_1 + \gamma^2 R_2 + \dots, \text{ where } 0 \leq \gamma < 1$$

Summary of Terminology

- Action (A): possible moves that agent can take
- State (S): Current situation returned by environment
- Reward (*R*): Immediate return sent back from the environment to evaluate the last action
- Policy (π) : Strategy that agent employs to determine next action based on current state
- Value (V): Expected long-term return with discount, as opposed to short-term reward

Learning a Control Policy

- Target function to be learned is a control policy, $\pi: S \rightarrow A$, that outputs action *a* given state $s \in S$
- Determine as to what action to take in a particular situation, so as to maximize cumulative reward
- Problem is one of learning to control a sequential process
 - In manufacturing optimization
 - What sequence of manufacturing actions must be chosen
 - Reward to be maximized is value of goods produced minus cost

How RL differs from other ML

- 1. Delayed Reward
 - In other types of ML, training example is $<\!s,a=\pi(s)>$
 - $\ln RL$, trainer provides immediate reward for a
 - Which actions are to be credited for outcome
- 2. Exploration
 - Agent influences distribution: by action sequence chosen. So which experimentation produces best learning?
 - Exploration of unknown states and actions?
 - Or exploitation of states and actions already learned
- 3. Partially observable states
 - Entire state may not be observable
 - Need to combine previous observations with current sensor data when choosing actions

Multi-tasking

- Robot learning may involve learning several related tasks
- Mobile robot may need to:
 - Dock on its battery charger
 - Navigate through narrow corridors
 - How to pick up output from a laser printer

Exploration and Exploitation

- RL requires choice: exploring vs exploiting
- Exploitation
 - Refers to taking actions that come from the current best version of the learned policy
 - Actions that we know will achieve a high reward
- Exploration
 - Refers to taking actions specifically to obtain more training data

Policy with exploration/exploitation

- Given context \boldsymbol{x} , action a gives us a reward of 1
- We do not know if it is the best possible reward
- We may want to exploit our current policy and continue taking action *a* to be sure of obtaining reward of 1
- We may also want to explore by trying action a'
- We do not know what will happen if we try a'
- We hope to get a reward of $2,\,\mathrm{but}$ we run risk of getting a reward of 0
- Either way we get some knowledge

Implementation of Exploration

- Implemented in many ways
 - Occasionally taking random actions intended to cover the entire range of possible actions
 - Model-based approaches that compute a choice of action based on its expected reward and the model's uncertainty about that reward

Preference for exploration or exploitation

- Factors for preference
- If agent has only a short amount of time to accrue reward then we prefer exploitation
- If agent has a long time to accrue reward we begin with more exploraion
 - So that future actions can be planned more effectively with knowledge
 - As time progresses we move towards more exploitation

RL connected to deep neural net

- Task: Learning to navigate in complex environments without prior knowledge.
- RL agent infers from complex environments by punishment-reward system. It can model decision making process.
- Example Applications:
 - AlphaGo Zero beat the world champion (December 2017)
 - OpenAI bot won in Dota2 world championship (Aug 2018)

The Learning Task

- Formulate the problem of learning sequential control strategies more precisely
- Many ways to do so
 - Agents actions are deterministic or nondeterministic
 - Agent can predict next state due to an action or it cannot
 - Training by an expert may show optimal action sequences
- A general formulation is based on Markov
 Decision Processes

Markov Decision Process (MDP)

- Agent can perceive a set *S* of discrete states in its environment
- Has a set A of actions that it can perform
- At each discrete time step t the agent senses the current state s_t , chooses a current action a_t and performs it
- Environment responds by giving it a reward $r_t = r(s_t, a_t)$ and by producing state $s_{t+1} = \delta(s_t, a_t)$
 - The function δ and r are part of the environment and not necessarily know to the agent

Srihari

MDP dependency

• In an MDP the two functions

 $\delta(s_t, a_t) = s_{t+1}$ which maps the current state to the next state

 $r(s_t, a_t) = r_t$ which provides the reward for the action

depend only on the current state and action and not on earlier states and actions

- We consider only case where \boldsymbol{S} and \boldsymbol{A} are finite
- We first consider the deterministic case for δ and r

A simple grid world environment



Six grid squares represent six possible states for the agent Each arrow represents a possible action can take to move to another state The number with each arrow is immediate reward r(s,a) agent receives gives reward of 100 for actions entering the goal state G and zero otherwise

G is the goal state for which the only action available is to remain in the state

How to specify Policy?

- Task of agent is to learn a policy $\pi: S \rightarrow A$
 - For selecting its next action a_t based on current observed state s_t , *i.e.*, $\pi(s_t) = a_t$
- How to specify the policy?
 - One approach: Require policy that produces greatest cumulative reward over time
 - For this we need to define cumulative value $V^{\pi}(s_t)$ of a policy π starting from an initial state s_t as follows

$$\begin{split} V^{\pi}(s_t) &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots \\ &= \sum_{i=0}^{\infty} \gamma^i r_{i+1} \end{split}$$

Here $0 \le \gamma < 1$ is a constant that determines the relative value of delayed vs immediate rewards (see next slide)

• Where sequence of rewards r_{t+1} is generated by beginning at state s_t and by repeatedly using policy π

26

