DSC291: Machine Learning with Few Labels

Reinforcement Learning

Zhiting Hu Lecture 14, May 15, 2025



HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

• Reinforcement Learning

- Paper presentation:
 - Max Irby, Davit Abrahamyan: "Training Large Language Models to Reason in a Continuous Latent Space"
 - Yuyuan, Varun: "The Belief State Transformer"

Markov Decision Process

- At time step t=0, environment samples initial state $s_0 \sim p(s_0)$
- Then, for t=0 until done:
 - Agent selects action $a_t \sim \pi(a_t/S_t)$

sparse re

- 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 policy π is a function from S to A that specifies what action to take in each state

Agent

Environment

1220 t<T 1770 0

 $t \ge 0$

Action a

State s.

Reward r,

Next state s

- 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:

The optimal policy π^*

Following a policy produces sample *trajectories* (or paths) s_0 , a_0 , r_0 , s_1 , a_1 , r_1 , ...

We want to find optimal policy π^* that maximizes the sum of rewards $\sum_{n=1}^{\infty} 2^{n}$

Question: How do we handle the randomness (initial state, transition probability...)? Maximize the **expected sum of rewards!**

Formally:
$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t \ge 0} \gamma^t r_t | \pi \right]$$
 with $s_0 \sim p(s_0), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$

Definitions: Value function and Q-value function

Following a policy produces sample *trajectories* (or paths) s_0 , a_0 , r_0 , s_1 , a_1 , r_1 , ...

Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) s_0 , a_0 , r_0 , s_1 , a_1 , r_1 , ...

How good is a state?

The value function at state s, is the expected cumulative reward from following the policy from state s: $V^{\pi}(s) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, \pi
ight]$

 $V(\hat{s}_0) = \hat{s}_0 \cdot P(\hat{s}_0)$ Question: express the objective with the value function

Śo



Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) s_0 , a_0 , r_0 , s_1 , a_1 , r_1 , ...

How good is a state?

The value function at state s, is the expected cumulative reward from following the policy from state s: $V^{\pi}(s) = \mathbb{E}\left[\sum \gamma^{t} r_{t} | s_{0} = s, \pi\right] \quad Q \quad (s, a) = \int \sqrt{\pi} \int \sqrt{\pi}$

How good is a state-action pair

The **Q-value function** at state s and action a, is the expected cumulative reward from taking action a in state s and then following the policy:

$$Q^{\pi}(s,a) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi\right]$$

Question: express the value function with the Q-value function

Bellman equation

The optimal Q-value function Q* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s,a) = \max_{\pi} \mathbb{E}\left[\sum_{t \ge 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi
ight]$$

Bellman equation

The optimal Q-value function Q* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s,a) = \max_{\pi} \mathbb{E} \left[\sum_{t \ge 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

Q^{*} satisfies the following **Bellman equation**:

lowing Bellman equation:

$$Q^*(s,a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^*(s',a') | s, a \right]$$

Intuition: if the optimal state-action values for the next time-step Q*(s',a') are known, then the optimal strategy is to take the action that maximizes the expected value of

Sa -> C'

Bellman equation

The optimal Q-value function Q* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s,a) = \max_{\pi} \mathbb{E}\left[\sum_{t \ge 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi
ight]$$

Q* satisfies the following **Bellman equation**:

$$Q^*(s,a) = \mathbb{E}_{s'\sim\mathcal{E}}\left[r + \gamma \max_{a'} Q^*(s',a')|s,a\right] \quad \mathcal{U}$$

Intuition: if the optimal state-action values for the next time-step Q*(s',a') are known, then the optimal strategy is to take the action that maximizes the expected value of $r + \gamma Q^*(s',a')$ The optimal policy π^* corresponds to taking the best action in any state as specified by Q*

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s,a) = \mathbb{E}\left[r + \gamma \max_{a'} Q_i(s',a')|s,a\right]$$

 Q_i will converge to Q^* as i -> infinity

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s,a) = \mathbb{E}\left[r + \gamma \max_{a'} Q_i(s',a')|s,a\right]$$

 Q_i will converge to Q^* as i -> infinity

Question: What's the problem with this?

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s,a) = \mathbb{E}\left[r + \gamma \max_{a'} Q_i(s',a') | s, a'\right]$$

 Q_i will converge to Q^* as i -> infinity

Question: What's the problem with this?

Not scalable. Must compute Q(s.a) for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!



Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s,a) = \mathbb{E}\left[r + \gamma \max_{a'} Q_i(s',a')|s,a\right]$$

 Q_i will converge to Q^* as i -> infinity

Question: What's the problem with this?

Not scalable. Must compute Q(s,a) for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!

Question: How would you solve the issue?

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s,a) = \mathbb{E}\left[r + \gamma \max_{a'} Q_i(s',a')|s,a\right]$$

 Q_i will converge to Q^* as i -> infinity

Question: What's the problem with this?

Not scalable. Must compute Q(s,a) for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!

Question: How would you solve the issue?

Use a function approximator to estimate Q(s,a). E.g. a neural network!

 $a) \sim Q$

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s,a;\theta) \approx Q^*(s,a)$$

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s,a;\theta) \approx Q^*(s,a)$$

If the function approximator is a deep neural network => deep q-learning!

Alpha Go

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s,a;\theta) \approx Q^*(s,a)$$

function parameters (weights)

If the function approximator is a deep neural network => deep q-learning!



Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^{*}(s,a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^{*}(s',a') | s, a \right]$$

$$Q^{*}(S,a) = G_{S'} Q^{*}(s',a') = Q_{Q^{*}}(s',a') | s, a = Q_{Q^{*}}(s',a$$

Remember: want to find a Q-function that satisfies the Bellman Equation:

Forward Pass
Loss function:
$$L_{i}(\theta_{i}) = \mathbb{E}_{s,a \sim \rho(\cdot)} \left[(y_{i} - Q(s,a;\theta_{i}))^{2} \right]$$
where $y_{i} = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s',a';\theta_{i-1}) | s, a \right]$

Solving for the optimal policy: Q-learning Remember: want to find a Q-function that satisfies the Bellman Equation: $Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^*(s', a') | s, a \right]$

Forward Pass

Loss function:
$$L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot)} \left[(y_i - Q(s,a;\theta_i))^2 \right]$$

where
$$y_i \models \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a \right]$$

Backward Pass

Gradient update (with respect to Q-function parameters)

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s',a';\theta_{i-1}) - Q(s,a;\theta_i)) \nabla_{\theta_i} Q(s,a;\theta_i) \right]$$

$$\overline{}$$

tan

[Mnih et al. NIPS Workshop 2013; Nature 2015]

Case Study: Playing Atari Games



Objective: Complete the game with the highest score

State: Raw pixel inputs of the game stateAction: Game controls e.g. Left, Right, Up, DownReward: Score increase/decrease at each time step



Q(s,a; heta) : neural network with weights $\, heta$



Q(s,a; heta) : neural network with weights heta



Q(s,a; heta) : neural network with weights heta



Q(s,a; heta) : neural network with weights heta



Q(s,a; heta): neural network with weights heta

A single feedforward pass to compute Q-values for all actions from the current state => efficient!



Last FC layer has 4-d output (if 4 actions), corresponding to $Q(s_t, a_1)$, $Q(s_t, a_2)$, $Q(s_t, a_3)$, $Q(s_t, a_4)$

Number of actions between 4-18 depending on Atari game

Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:

- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:

- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

Address these problems using **experience replay**

- Continually update a replay memory table of transitions (s_t, a_t, r_t, s_{t+1}) as game (experience) episodes are played
- Train Q-network on <u>random minibatches of transitions</u> from the replay memory, instead of consecutive samples

Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:

- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

Address these problems using **experience replay**

- Continually update a replay memory table of transitions (s_t, a_t, r_t, s_{t+1}) as game (experience) episodes are played
- Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples

Each transition can also contribute to multiple weight updates => greater data efficiency

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N Initialize action-value function Q with random weights for episode = 1, M do Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$ for t = 1, T do With probability ϵ select a random action a_t otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ Execute action a_t in emulator and observe reward r_t and image x_{t+1} Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$ Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D} Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D} Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3 end for end for

Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory \mathcal{D} to capacity N Initialize replay memory, Q-network Initialize action-value function Q with random weights for episode = 1, M do Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$ for t = 1, T do With probability ϵ select a random action a_t otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ Execute action a_t in emulator and observe reward r_t and image x_{t+1} Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$ Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D} Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D} Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3 end for end for

Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory \mathcal{D} to capacity N Initialize action-value function Q with random weights —— Play M <u>episodes</u> (full games) for episode = 1, M do Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$ for t = 1, T do With probability ϵ select a random action a_t otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ Execute action a_t in emulator and observe reward r_t and image x_{t+1} Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$ Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D} Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D} Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3 end for end for

Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory \mathcal{D} to capacity N Initialize action-value function Q with random weights for episode = 1, M do Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$ Initialize state for t = 1, T do (starting game) With probability ϵ select a random action a_t screen pixels) at the otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ beginning of each Execute action a_t in emulator and observe reward r_t and image x_{t+1} episode Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$ Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D} Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D} Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3 end for end for



Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory \mathcal{D} to capacity N Initialize action-value function Q with random weights for episode = 1, M do Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$ for t = 1, T do With small probability, With probability ϵ select a random action a_t otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ select a random Execute action a_t in emulator and observe reward r_t and image x_{t+1} action (explore Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$ otherwise select Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D} greedy action from Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D} current policy Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3 end for end for

Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory \mathcal{D} to capacity N Initialize action-value function Q with random weights for episode = 1, M do Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$ for t = 1, T do With probability ϵ select a random action a_t otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ Execute action a_t in emulator and observe reward r_t and image x_{t+1} Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$ Take the action (a_t) , Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D} and observe the Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D} reward r, and next Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ state s_{t+1} Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3 end for end for

Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory \mathcal{D} to capacity N Initialize action-value function Q with random weights for episode = 1, M do Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$ for t = 1, T do With probability ϵ select a random action a_t otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ Execute action a_t in emulator and observe reward r_t and image x_{t+1} Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$ Store transition in Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D} replay memory Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D} Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3 end for end for

Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory \mathcal{D} to capacity N Initialize action-value function Q with random weights \rightarrow for episode = 1, M do Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$ for t = 1, T do With probability ϵ select a random action a_t otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ Execute action a_t in emulator and observe reward r_t and image x_{t+1} Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$ Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D} Experience Replay: Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D} Sample a random Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ minibatch of transitions from replay memory Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3 and perform a gradient end for descent step end for

Summary so far

- Q-learning:
 - Bellman equation
 - Value-based RL
 - Off-policy RL

$$\left[\begin{array}{ll} \text{Loss function:} & L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot)} \left[(y_i - Q(s,a;\theta_i))^2 \right] \\ \text{where} & y_i = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s',a';\theta_{i-1}) | s,a \right] \end{array} \right]$$

- Next: Policy gradient
 - Policy-based RL
 - On-policy RL

Questions?