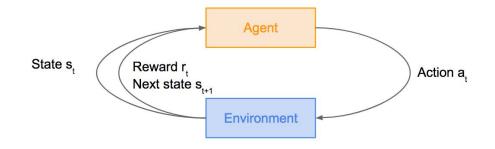
Lecture 17: Reinforcement Learning

Today: Reinforcement Learning

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

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





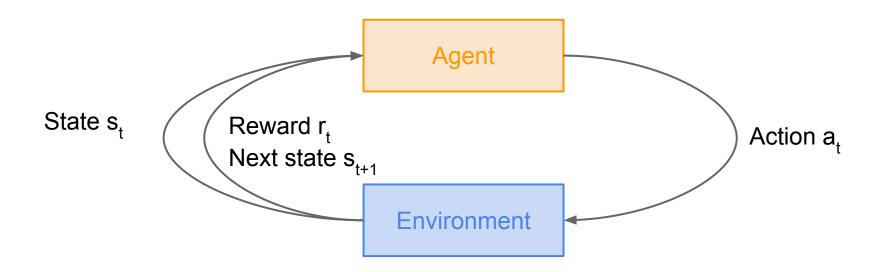
Atari games figure copyright Volodymyr Mnih et al., 2013. Reproduced with permission.

Overview

- What is Reinforcement Learning?
- Markov Decision Processes
- Q-Learning
- Policy Gradients
- SOTA Deep RL

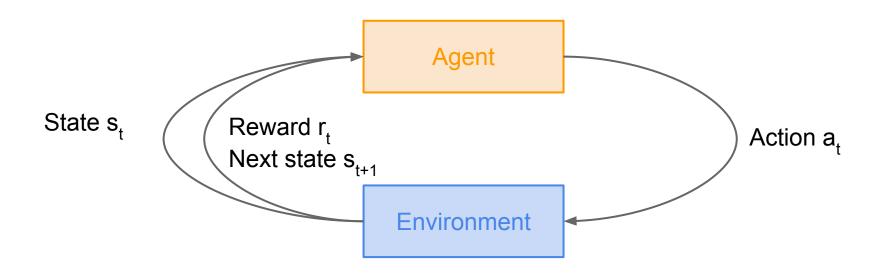
What is reinforcement learning?

Reinforcement Learning



Markov decision processes (MDPs)

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 world

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

 ${\cal S}\,$: set of possible states

 \mathcal{A} : set of possible actions

 \mathcal{R} : distribution of reward given (state, action) pair

P: transition probability i.e. distribution over next state given (state, action) pair

 γ : discount factor

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,
 - Environment samples reward $r_t \sim R(. | s_t, a_t)$
 - Environment samples next state $s_{t+1} \sim P(\cdot, |s_t, a_t)$
 - Agent receives reward r, and next state s,

- A policy π is a function from S to A that specifies what action to take in each state
- **Objective**: find policy $\mathbf{\pi}^*$ that maximizes cumulative discounted reward: $\sum \gamma^t r_t$

A simple MDP: Grid World

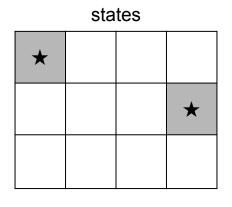
```
actions = {

1. right →

2. left →

3. up ↑

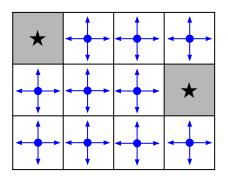
4. down ↑
```



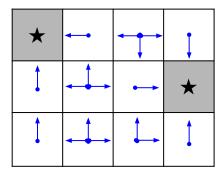
Set a negative "reward" for each transition (e.g. r = -1)

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

A simple MDP: Grid World



Random Policy



Optimal Policy

The optimal policy π^*

We want to find optimal policy π^* that maximizes the sum of rewards.

The optimal policy π^*

We want to find optimal policy π^* that maximizes the sum of rewards.

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

The optimal policy π^*

We want to find optimal policy π^* that maximizes the sum of rewards.

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

Formally:
$$\pi^* = \arg\max_{\pi} \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | \pi\right]$$
 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:

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

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:

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

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:

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

Q-learning

Bellman equation

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

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

Bellman equation

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

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

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

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

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 $r + \gamma Q^*(s',a')$

Bellman equation

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

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

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

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

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 $\,\pi^* = rg \max_a Q^*(s,a)\,$

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

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

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

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!

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

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!

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

$$Q(s,a; heta)pprox Q^*(s,a)$$
 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:

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

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

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

Backward Pass

Gradient update (with respect to Q-function parameters θ):

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

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

where
$$y_i = \mathbb{E}_{s' \sim \mathcal{E}}\left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a \right]$$
 Iteratively try to make the Q-value close to the target value (y_i) it

should have, if Q-function corresponds to optimal Q* (and optimal policy π^*)

Backward Pass

Gradient update (with respect to Q-function parameters θ):

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

Case Study: Playing Atari Games



Objective: Complete the game with the highest score

State: Raw pixel inputs of the game state

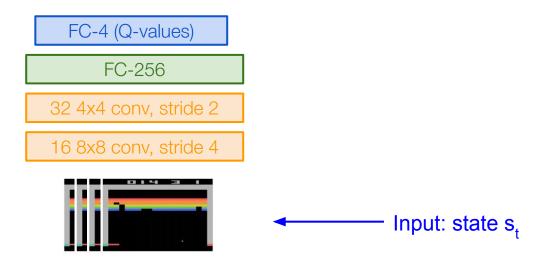
Action: Game controls e.g. Left, Right, Up, Down

Reward: Score increase/decrease at each time step

Figures copyright Volodymyr Mnih et al., 2013. Reproduced with permission.

Q-network Architecture

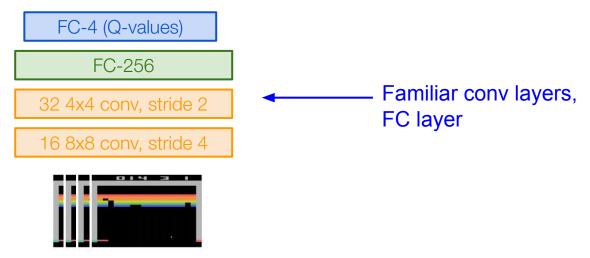
Q(s,a; heta) : neural network with weights heta



Current state s_t: 84x84x4 stack of last 4 frames (after RGB->grayscale conversion, downsampling, and cropping)

Q-network Architecture

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

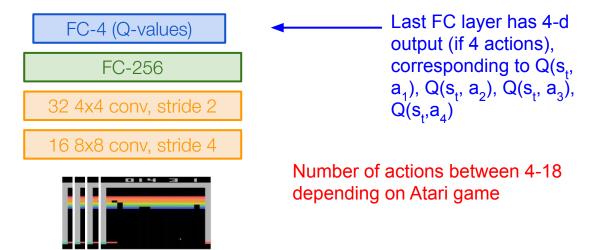


Current state s_t: 84x84x4 stack of last 4 frames (after RGB->grayscale conversion, downsampling, and cropping)

Q-network Architecture

we don't have softmax layer after the FC because here our goal is to directly predict our Q-value functions.

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

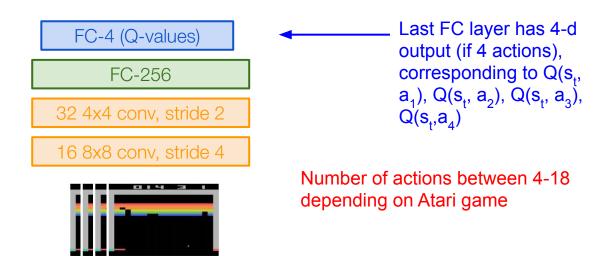


Current state s_t: 84x84x4 stack of last 4 frames (after RGB->grayscale conversion, downsampling, and cropping)

Q-network Architecture

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!



Current state s_t: 84x84x4 stack of last 4 frames (after RGB->grayscale conversion, downsampling, and cropping)

Training the Q-network: Loss function (from before)

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

where
$$y_i = \mathbb{E}_{s' \sim \mathcal{E}}\left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a \right]$$
 Iteratively try to make the Q-value close to the target value (y_i) it

Iteratively try to make the Q-value close to the target value (y_i) it should have, if Q-function corresponds to optimal Q* (and optimal policy π^*)

Problem: Nonstationary! The "target" for Q(s, a) depends on the current weights θ !

Solution: Use a slow-moving "target" Q-network that is delayed in parameter updates to generate target value labels for the Q-network.

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 random minibatches of transitions 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 \epsilon 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_i, a_i, r_i, \phi_{i+1}) from \mathcal{D}
            \text{Set } y_j = \left\{ \begin{array}{ll} 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{array} \right.
            Perform a gradient descent step on (y_i - Q(\phi_i, a_i; \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 \epsilon 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_i, a_i, r_i, \phi_{i+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_i - Q(\phi_i, a_i; \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 episodes (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 \epsilon 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_i, a_i, r_i, \phi_{i+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_i - Q(\phi_i, a_i; \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 \epsilon 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}
                                                                                                                             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_i, a_i, r_i, \phi_{i+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_i - Q(\phi_i, a_i; \theta))^2 according to equation 3
       end for
   end for
```

Initialize state (starting game screen pixels) at the beginning of each

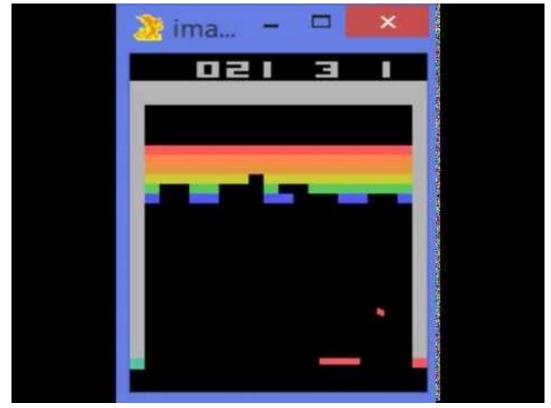
```
Algorithm 1 Deep Q-learning with Experience Replay
   Initialize replay memory \mathcal{D} to capacity N
   Initialize action-value function Q with random weights
   for episode = 1, M do
        Initialise sequence s_1 = \{x_1\} and preprocessed sequenced \phi_1 = \phi(s_1)
       for t = 1, T do
                                                                                                                               For each timestep t
            With probability \epsilon select a random action a_t
                                                                                                                               of the game
            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_i, a_i, r_i, \phi_{i+1}) from \mathcal{D}
            \text{Set } y_j = \left\{ \begin{array}{ll} 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{array} \right.
            Perform a gradient descent step on (y_i - Q(\phi_i, a_i; \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 NInitialize 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 With small probability, 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 $(\phi_i, a_i, r_i, \phi_{i+1})$ 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 $(y_i - Q(\phi_i, a_i; \theta))^2$ 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 \epsilon 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,),
                                                                                                                           and observe the
            Store transition (\phi_t, a_t, r_t, \phi_{t+1}) in \mathcal{D}
            Sample random minibatch of transitions (\phi_i, a_i, r_i, \phi_{i+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<sub>t+1</sub>
            Perform a gradient descent step on (y_i - Q(\phi_i, a_i; \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 \epsilon 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_i, a_i, r_i, \phi_{i+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_i - Q(\phi_i, a_i; \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 \epsilon 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_i, a_i, r_i, \phi_{i+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}
                                                                                                                  Sample a random
                                                                                                                  minibatch of transitions
                                                                                                                  from replay memory
           Perform a gradient descent step on (y_i - Q(\phi_i, a_i; \theta))^2 according to equation 3
                                                                                                                  and perform a gradient
       end for
                                                                                                                  descent step
   end for
```



https://www.youtube.com/watch?v=V1eYniJ0Rnk

Video by Károly Zsolnai-Fehér. Reproduced with permission.

What is a problem with Q-learning?

The Q-function can be very complicated!

Example: a robot grasping an object has a very high-dimensional state => hard to learn exact value of every (state, action) pair

What is a problem with Q-learning?

The Q-function can be very complicated!

all the joint positions, angles,...

Example: a robot grasping an object has **a very high-dimensional state** => hard to learn exact value of every (state, action) pair

But the policy can be much simpler: just close your hand Can we learn a policy directly, e.g. finding the best policy from a collection of policies?

Policy gradients

Policy Gradients

Formally, let's define a class of parameterized policies: $\Pi = \{\pi_{\theta}, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(heta) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | \pi_ heta
ight]$$

Policy Gradients

Formally, let's define a class of parameterized policies: $\Pi = \{\pi_{\theta}, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(heta) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | \pi_ heta
ight]$$

We want to find the optimal policy $\ heta^* = \arg\max_{ heta} J(heta)$

How can we do this?

Policy Gradients

Formally, let's define a class of parameterized policies: $\Pi = \{\pi_{\theta}, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(heta) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | \pi_ heta
ight]$$

We want to find the optimal policy $\ \theta^* = \arg\max_{\theta} J(\theta)$

How can we do this?

Gradient ascent on policy parameters!

Mathematically, we can write:

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)} [r(\tau)]$$
$$= \int_{\tau} r(\tau) p(\tau;\theta) d\tau$$

Where $\mathbf{r}(au)$ is the reward of a trajectory $\ au=(s_0,a_0,r_0,s_1,\ldots)$

And
$$p(\tau;\theta) = \prod_{t \geq 0} p(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t)$$

Expected reward: $J(heta) = \mathbb{E}_{ au \sim p(au; heta)}\left[r(au)
ight] \ = \int_{ au} r(au) p(au; heta) \mathrm{d} au$

Expected reward:
$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)}\left[r(\tau)\right]$$

$$= \int_{\tau} r(\tau)p(\tau;\theta)\mathrm{d}\tau$$

Now let's differentiate to calculate the gradients:

$$\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$$

Expected reward:
$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)}\left[r(\tau)\right]$$

$$= \int_{-}^{} r(\tau)p(\tau;\theta)\mathrm{d}\tau$$

Now let's differentiate to calculate the gradients:

$$\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$$

Intractable! Gradient of an expectation is problematic when p depends on θ

$$p(\tau;\theta) = \prod_{t>0} p(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t)$$

Expected reward:
$$J(\theta) = \mathbb{E}_{ au \sim p(au; heta)} \left[r(au)
ight]$$
 $= \int_{ au} r(au) p(au; heta) \mathrm{d} au$

Now let's differentiate to calculate the gradients:

$$abla_{ heta}J(heta) = \int_{ au} r(au)
abla_{ heta} p(au; heta) \mathrm{d} au$$

Intractable! Gradient of an expectation is problematic when p depends on θ

$$p(\tau;\theta) = \prod_{t>0} p(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t)$$

However, we can use a nice trick:

$$abla_{\theta} p(\tau; \theta) = p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta)$$

 $J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)} \left[r(\tau) \right]$ Expected reward: $=\int_{-}^{}r(\tau)p(\tau;\theta)\mathrm{d}\tau$

Now let's differentiate to calculate the gradients:

$$\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$$

Intractable! Gradient of an expectation is problematic when p depends on θ

$$p(\tau;\theta) = \prod_{t>0} p(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t)$$

However, we can use a nice trick:

$$\begin{split} \nabla_{\theta} p(\tau;\theta) &= p(\tau;\theta) \frac{\nabla_{\theta} p(\tau;\theta)}{p(\tau;\theta)} = p(\tau;\theta) \nabla_{\theta} \log p(\tau;\theta) \\ \nabla_{\theta} J(\theta) &= \int_{\tau} \left(r(\tau) \nabla_{\theta} \log p(\tau;\theta) \right) p(\tau;\theta) \mathrm{d}\tau \\ &= \mathbb{E}_{\tau \sim p(\tau;\theta)} \left[r(\tau) \nabla_{\theta} \log p(\tau;\theta) \right] \end{split} \quad \begin{array}{l} \text{Can estimate with Monte Carlo sampling!} \\ \text{Let's see how in the next slide} \end{split}$$

Let's see how in the next slide

Can we compute gradients of a trajectory without knowing the transition probabilities?

We have:
$$p(au; heta) = \prod_{t>0} p(s_{t+1}|s_t,a_t)\pi_{ heta}(a_t|s_t)$$

Can we compute gradients of a trajectory without knowing the transition probabilities?

We have:
$$p(au; heta) = \prod p(s_{t+1}|s_t,a_t)\pi_{ heta}(a_t|s_t)$$

Thus:
$$\log p(\tau;\theta) = \sum_{t>0}^{t\geq 0} \log p(s_{t+1}|s_t,a_t) + \log \pi_{\theta}(a_t|s_t)$$

Can we compute gradients of a trajectory without knowing the transition probabilities?

We have:
$$p(au; heta) = \prod p(s_{t+1}|s_t,a_t)\pi_{ heta}(a_t|s_t)$$

Thus:
$$\log p(\tau;\theta) = \sum_{t>0}^{t\geq 0} \log p(s_{t+1}|s_t,a_t) + \log \pi_{\theta}(a_t|s_t)$$

And when differentiating:
$$\nabla_{\theta} \log p(au; heta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Doesn't depend on transition probabilities!

$$\nabla_{\theta} J(\theta) = \int_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau$$
$$= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau) \nabla_{\theta} \log p(\tau; \theta)]$$

Can we compute gradients of a trajectory without knowing the transition probabilities?

We have:
$$p(au; heta) = \prod p(s_{t+1}|s_t, a_t) \pi_{ heta}(a_t|s_t)$$

Thus:
$$\log p(\tau; \theta) = \sum_{t>0}^{t\geq 0} \log p(s_{t+1}|s_t, a_t) + \log \pi_{\theta}(a_t|s_t)$$

And when differentiating:
$$\nabla_{\theta} \log p(\tau; \theta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$$
 Doesn't depend on transition probabilities!

Therefore when sampling a trajectory τ , we can estimate $J(\theta)$ with

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Intuition

Gradient estimator: $\nabla_{\theta} J(\theta) pprox \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen
- If $r(\tau)$ is low, push down the probabilities of the actions seen

Intuition

Gradient estimator: $\nabla_{\theta} J(\theta) pprox \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen
- If $r(\tau)$ is low, push down the probabilities of the actions seen

Might seem simplistic to say that if a trajectory is good then all its actions were good. But in expectation, it averages out!

However, this also suffers from high variance because credit assignment is really hard. Can we help the estimator?

Variance reduction

Gradient estimator: $\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

Variance reduction

Gradient estimator: $\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

First idea: Push up probabilities of an action seen, only by the cumulative future reward from that state

$$abla_{ heta} J(heta) pprox \sum_{t \geq 0} \left(\sum_{t' \geq t} r_{t'} \right)
abla_{ heta} \log \pi_{ heta}(a_t | s_t)$$

Variance reduction

Gradient estimator: $\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

First idea: Push up probabilities of an action seen, only by the cumulative future reward from that state

$$abla_{ heta} J(heta) pprox \sum_{t \geq 0} \left(\sum_{t' \geq t} r_{t'} \right)
abla_{ heta} \log \pi_{ heta}(a_t | s_t)$$

Second idea: Use discount factor γ to ignore delayed effects

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Variance reduction: Baseline

Problem: The raw value of a trajectory isn't necessarily meaningful. For example, if rewards are all positive, you keep pushing up probabilities of actions.

What is important then? Whether a reward is better or worse than what you expect to get

Idea: Introduce a baseline function dependent on the state. Concretely, estimator is now:

$$abla_{\theta} J(\theta) pprox \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_{t'} - b(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

Intuitively, we are happy with an action a_t in a state s_t if $Q^{\pi}(s_t, a_t) - V^{\pi}(s_t)$ is large. On the contrary, we are unhappy with an action if it's small.

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

Intuitively, we are happy with an action a_t in a state s_t if $Q^{\pi}(s_t, a_t) - V^{\pi}(s_t)$ is large. On the contrary, we are unhappy with an action if it's small.

Using this, we get the estimator: $\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} (Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t)) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

Actor-Critic Algorithm

Problem: we don't know Q and V. Can we learn them?

Yes, using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q-function).

- The actor decides which action to take, and the critic tells the actor how good its action was and how it should adjust
- Also alleviates the task of the critic as it only has to learn the values of (state, action) pairs generated by the policy
- Can also incorporate Q-learning tricks e.g. experience replay
- **Remark:** we can define by the **advantage function** how much an action was better than expected $A^{\pi}(s, x) = O^{\pi}(s, x)$

 $A^{\pi}(s,a) = Q^{\pi}(s,a) - V^{\pi}(s)$

Actor-Critic Algorithm

```
Initialize policy parameters \theta, critic parameters \phi
For iteration=1, 2 ... do
          Sample m trajectories under the current policy
          \Delta\theta \leftarrow 0
          For i=1, ..., m do
                    For t=1, ..., T do
                              A_t = \sum_{t' \ge t} \gamma^{t'-t} r_t^i - V_\phi(s_t^i)
                              \Delta \theta \leftarrow \Delta \theta + A_t \nabla_{\theta} \log(a_t^i | s_t^i)
         \Delta \phi \leftarrow \sum_{i} \sum_{t} \nabla_{\phi} ||A_{t}^{i}||^{2}\theta \leftarrow \alpha \Delta \theta
          \phi \leftarrow \beta \Delta \phi
```

End for

REINFORCE in action: Recurrent Attention Model (RAM)

Objective: Image Classification

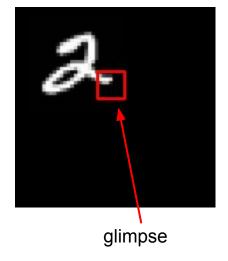
Take a sequence of "glimpses" selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image

State: Glimpses seen so far

Action: (x,y) coordinates (center of glimpse) of where to look next in image

Reward: 1 at the final timestep if image correctly classified, 0 otherwise



[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)

Objective: Image Classification

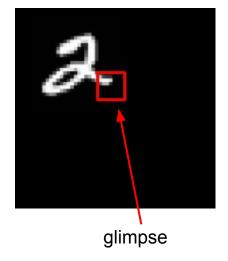
Take a sequence of "glimpses" selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image

State: Glimpses seen so far

Action: (x,y) coordinates (center of glimpse) of where to look next in image

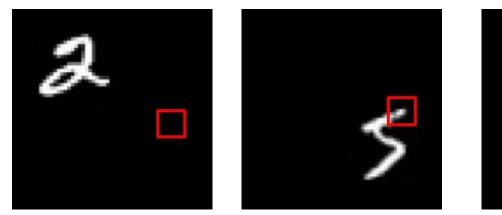
Reward: 1 at the final timestep if image correctly classified, 0 otherwise

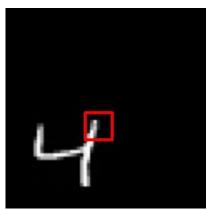


Glimpsing is a non-differentiable operation => learn policy for how to take glimpse actions using REINFORCE Given state of glimpses seen so far, use RNN to model the state and output next action

[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)





Has also been used in many other tasks including fine-grained image recognition, image captioning, and visual question-answering!

Figures copyright Daniel Levy, 2017. Reproduced with permission.

[Mnih et al. 2014]

SOTA Deep Reinforcement Learning

Recall Policy Gradient with baseline:

$$\max_{ heta} \mathbb{E}_t[\log \pi_{ heta}(a_t|s_t)\hat{A}_t]$$

Recall Policy Gradient with baseline:

$$\max_{ heta} \mathbb{E}_t[\log \pi_{ heta}(a_t|s_t)\hat{A}_t]$$

Objective from Conservative Policy Iteration [Kakade et al. 2002]:

$$egin{aligned} \max_{ heta} & \mathbb{E}_t ig[rac{\pi_{ heta}(a_t|s_t)}{\pi_{ ext{old}}(a_t|s_t)} \hat{A}_t ig] \ & \mathbb{E}_t [ext{KL}[\pi_{ ext{old}}(\cdot|s_t)||\pi_{ heta}(\cdot|s_t)]] \leq \delta \end{aligned}$$

Recall Policy Gradient with baseline:

$$\max_{ heta} \mathbb{E}_t[\log \pi_{ heta}(a_t|s_t)\hat{A}_t]$$

Objective from Conservative Policy Iteration [Kakade et al. 2002]:

$$egin{aligned} \max_{ heta} & \mathbb{E}_t ig[rac{\pi_{ heta}(a_t|s_t)}{\pi_{ ext{old}}(a_t|s_t)} \hat{A}_t ig] \ & \mathbb{E}_t [ext{KL}[\pi_{ ext{old}}(\cdot|s_t)||\pi_{ heta}(\cdot|s_t)]] \leq \delta \end{aligned}$$

Relaxation for PPO Objective:

$$\max_{ heta} \mathbb{E}_t [rac{\pi_{ heta}(a_t|s_t)}{\pi_{ ext{old}}(a_t|s_t)} \hat{A}_t - eta \cdot ext{KL}[\pi_{ ext{old}}(\cdot|s_t) || \pi_{ heta}(\cdot|s_t)]]$$

Relaxation for PPO Objective:

$$\max_{ heta} \mathbb{E}_t [rac{\pi_{ heta}(a_t|s_t)}{\pi_{ ext{old}}(a_t|s_t)} \hat{A}_t - eta \cdot ext{KL}[\pi_{ ext{old}}(\cdot|s_t) || \pi_{ heta}(\cdot|s_t)]]$$

Relaxation for PPO Objective:

$$\max_{ heta} \mathbb{E}_t [rac{\pi_{ heta}(a_t|s_t)}{\pi_{ ext{old}}(a_t|s_t)} \hat{A}_t - eta \cdot ext{KL}[\pi_{ ext{old}}(\cdot|s_t) || \pi_{ heta}(\cdot|s_t)]]$$

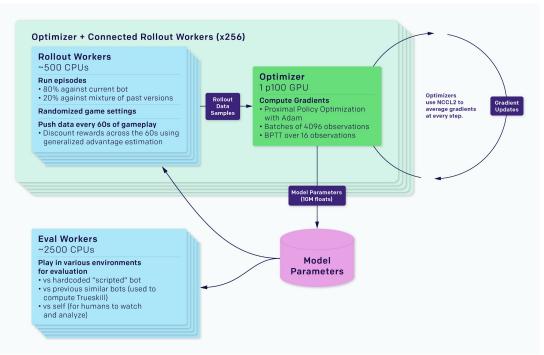
Key Idea: Tune KL penalty weight adaptively based on KL violation!

$$egin{aligned} d &= ext{KL}[\pi_{ ext{old}}(\cdot|s_t)||\pi_{ heta}(\cdot|s_t)] \ & \ d < d_{ ext{target}}/1.5: eta \leftarrow eta/2 \ & \ d < 1.5 \ d_{ ext{target}}: eta \leftarrow eta imes 2 \end{aligned}$$

PPO at scale: OpenAl Five and Dota 2

- 180 years of self-play per day
- 256 GPUs and 128,000
 CPU cores





["OpenAl Five", OpenAl 2018]

SOTA Deep RL: Soft Actor-Critic

Entropy-Regularized RL - reward policy for producing multi-modal actions

$$\mathcal{H}(\pi_{ heta}(\cdot|s)) = \mathbb{E}_{a \sim \pi_{ heta}(\cdot|s)}[-\log \pi_{ heta}(a|s)]$$

$$r'(s_t, a_t) = r(s_t, a_t) + lpha \mathcal{H}(\pi_{ heta}(\cdot|s))$$

["Soft Actor-Critic - Deep Reinforcement Learning with Real-World Robots", Haarnoja et al. 2018]

SOTA Deep RL: Soft Actor-Critic

Entropy-Regularized RL - reward policy for producing multi-modal actions

$$egin{aligned} \mathcal{H}(\pi_{ heta}(\cdot|s)) &= \mathbb{E}_{a \sim \pi_{ heta}(\cdot|s)}[-\log \pi_{ heta}(a|s)] \ r'(s_t, a_t) &= r(s_t, a_t) + lpha \mathcal{H}(\pi_{ heta}(\cdot|s)) \end{aligned}$$

Learn Soft Q-function with entropy regularization and the best policy under the Q-function.

$$\pi_{ heta} = rg\min_{ heta} \mathsf{KL}[\pi_{ heta}(\cdot|s)||rac{\exp(rac{1}{lpha}Q_{\psi}(s,\cdot))}{Z_{\psi}(s)}]$$

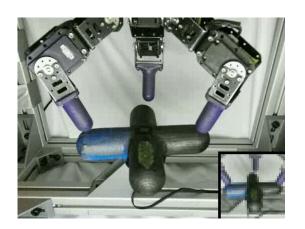
Best policy is given by a Boltzmann (softmax) distribution instead of maximum!

["Soft Actor-Critic - Deep Reinforcement Learning with Real-World Robots", Haarnoja et al. 2018]

SAC in action: Robotic Manipulation

Learn policies from scratch on a real robot in hours!





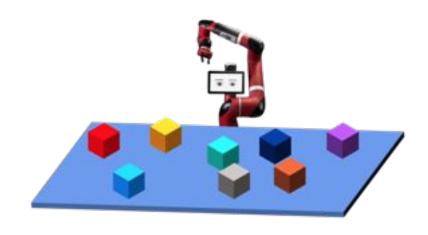


["Soft Actor-Critic - Deep Reinforcement Learning with Real-World Robots", Haarnoja et al. 2018]

Policy Learning in Robotics

Why is it hard?

- High-dimensional state space (e.g. visuomotor control)
- Continuous action space
- Hard exploration
- Sample efficiency

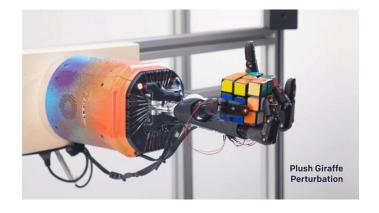


Addressing Sample Efficiency in Robotics

Sim2Real

- Train a policy in simulation, and then transfer it to real world.





["Solving Rubik's Cube with a Robotic Hand", OpenAI, 2019]

Addressing Sample Efficiency in Robotics

Imitation Learning

- Learn a policy from human demonstrations (provided via VR, phone, etc.)



["Deep Imitation Learning for Complex Manipulation Tasks from Virtual Reality Teleoperation", Zhang et al., 2017]



["Learning to Generalize Across Long-Horizon Tasks from Human Demonstrations", Mandlekar*, Xu* et al., 2020]

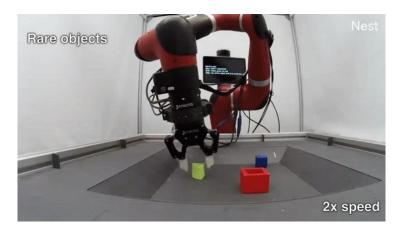
Addressing Sample Efficiency in Robotics

Batch Reinforcement Learning

Learn a policy from offline data collected by other policies (not necessarily expert).



["Implicit Reinforcement without Interaction at Scale for Learning Control from Offline Robot Manipulation Data", Mandlekar et al., 2020]



["Scaling data-driven robotics with reward sketching and batch reinforcement learning", Cabi et al., 2020]

Competing against humans in game play

AlphaGo [DeepMind, Nature 2016]:

- Required many engineering tricks
- Bootstrapped from human play
- Beat 18-time world champion Lee Sedol

AlphaGo Zero [Nature 2017]:

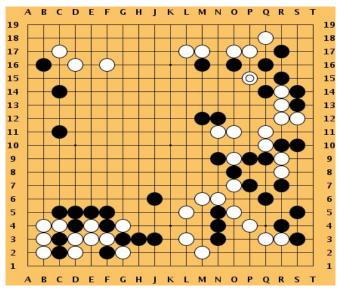
- Simplified and elegant version of AlphaGo
- No longer bootstrapped from human play
- Beat (at the time) #1 world ranked Ke Jie

Alpha Zero: Dec. 2017

 Generalized to beat world champion programs on chess and shogi as well

MuZero (November 2019)

Plans through a learned model of the game



This image is CC0 public domain

Silver et al, "Mastering the game of Go with deep neural networks and tree search", Nature 2016

Silver et al, "Mastering the game of Go without human knowledge", Nature 2017

Silver et al, "A general reinforcement learning algorithm that masters chess, shogi, and go through self-play", Science 2018

Schrittwieser et al, "Mastering Atari, Go, Chess and Shogi by Planning with a Learned Model", arXiv 2019

Competing against humans in game play

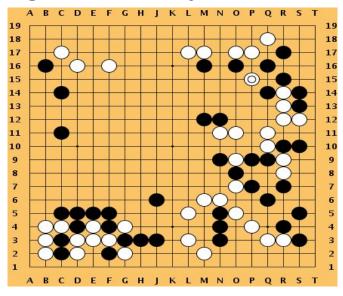
November 2019: Lee Sedol announces retirement



Quotes from <u>link</u>: <u>Image</u> of Lee Sedol is licensed under CC BY 2.0

"With the debut of AI in Go games, I've realized that I'm not at the top even if I become the number one through frantic efforts"

"Even if I become the number one, there is an entity that cannot be defeated"



This image is CC0 public domain

Silver et al, "Mastering the game of Go with deep neural networks and tree search", Nature 2016

Silver et al, "Mastering the game of Go without human knowledge", Nature 2017

Silver et al, "A general reinforcement learning algorithm that masters chess, shogi, and go through self-play", Science 2018

Schrittwieser et al, "Mastering Atari, Go, Chess and Shogi by Planning with a Learned Model", arXiv 2019

Summary

- Policy gradients: very general but suffer from high variance so requires a lot of samples. Challenge: sample-efficiency
- **Q-learning**: does not always work but when it works, usually more sample-efficient. **Challenge**: exploration
- Guarantees:
 - **Policy Gradients**: Converges to a local minima of $J(\theta)$, often good enough!
 - **Q-learning**: Zero guarantees since you are approximating Bellman equation with a complicated function approximator

Next Time: Scene Graphs



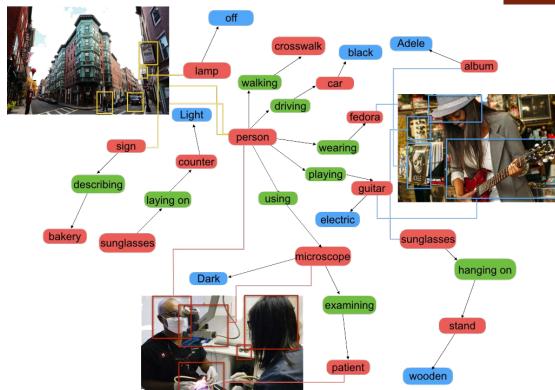


Figure copyright Krishna et al. 2017. Reproduced with permission.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen et al. "Visual genome: Connecting language and vision using crowdsourced dense image annotations." International Journal of Computer Vision 123, no. 1 (2017): 32-73.