



Reinforcement Learning

Computer Engineering Department Sharif University of Technology

Mohammad Hossein Rohban, Ph.D.

Spring 2025

Courtesy: Some slides are adopted from CS 285 Berkeley, and CS 234 Stanford, and Pieter Abbeel's compact series on RL.

Function Approximation

Function approximation and deep RL

- The **policy**, **value function**, **model**, and **agent state update** are all functions
- We want to learn these from experience (data).
- If there are too many states, we need to approximate.
- This is often called **deep reinforcement learning**, when using **neural networks** to represent these functions.

Large-Scale Reinforcement Learning

- Reinforcement learning can be used to solve **large** problems, e.g.
 - Backgammon: 10^{20} states
 - Go: 10^{170} states
 - Helicopter: continuous state space
 - Robots: real world
- How can we apply our methods for **prediction** and **control**?

Value Function Approximation

Value Function Approximation

- So far we mostly considered **lookup tables**
 - Every state s has an entry $v(s)$
 - Or every state-action pair s, a has an entry $q(s, a)$
- Problem with large MDPs:
 - There are too many states and/or actions to store in memory
 - It is too slow to learn the value of each state individually
 - Individual environment states are often **not fully observable**

Value Function Approximation

- Solution for large MDPs:
 - Estimate value function with **function approximation**

$$\begin{array}{ll} v_{\mathbf{w}}(s) \approx v_{\pi}(s) & \text{(or } v_*(s)\text{)} \\ q_{\mathbf{w}}(s, a) \approx q_{\pi}(s, a) & \text{(or } q_*(s, a)\text{)} \end{array}$$

- Update parameter w (e.g., using MC or TD learning)
- Generalize to unseen states

Key Requirements

Learning in new setup

- In principle, **any** function approximator can be used, but RL has specific properties:
 - Experience is **not iid** — successive time-steps are correlated
 - Agent's policy affects the data it receives
 - Regression targets can be **non-stationary**
 - ...because of changing policies (which can change the target and the data!)
 - ...because of bootstrapping
 - ...because of non-stationary dynamics (e.g., other learning agents)
 - ...because the world is large (never quite in the same state)

Gradient-based Algorithms

Approximate Values By Stochastic Gradient Descent

- Goal: find w that minimize the difference between $v_w(s)$ and $v_\pi(s)$

$$J(\mathbf{w}) = \mathbb{E}_{S \sim d}[(v_\pi(S) - v_w(S))^2]$$

- Gradient descent:

$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w}) = \alpha \mathbb{E}_d (v_\pi(S) - v_w(S)) \nabla_{\mathbf{w}} v_w(S)$$

- Stochastic gradient descent (SGD), sample the gradient:

$$\Delta \mathbf{w} = \alpha (G_t - v_w(S_t)) \nabla_{\mathbf{w}} v_w(S_t)$$

Note: **Monte Carlo return** G_t is a sample for $v_\pi(s_t)$

Training Loss

- We can't update towards the true value function $v_{\pi}(s)$
- We substitute a **target** for $v_{\pi}(s)$
 - For MC, the target is the return G_t

$$\Delta \mathbf{w}_t = \alpha(\mathbf{G}_t - v_{\mathbf{w}}(s)) \nabla_{\mathbf{w}} v_{\mathbf{w}}(s)$$

- For TD, the target is the TD target $R_{t+1} + \gamma v_{\mathbf{w}}(S_{t+1})$

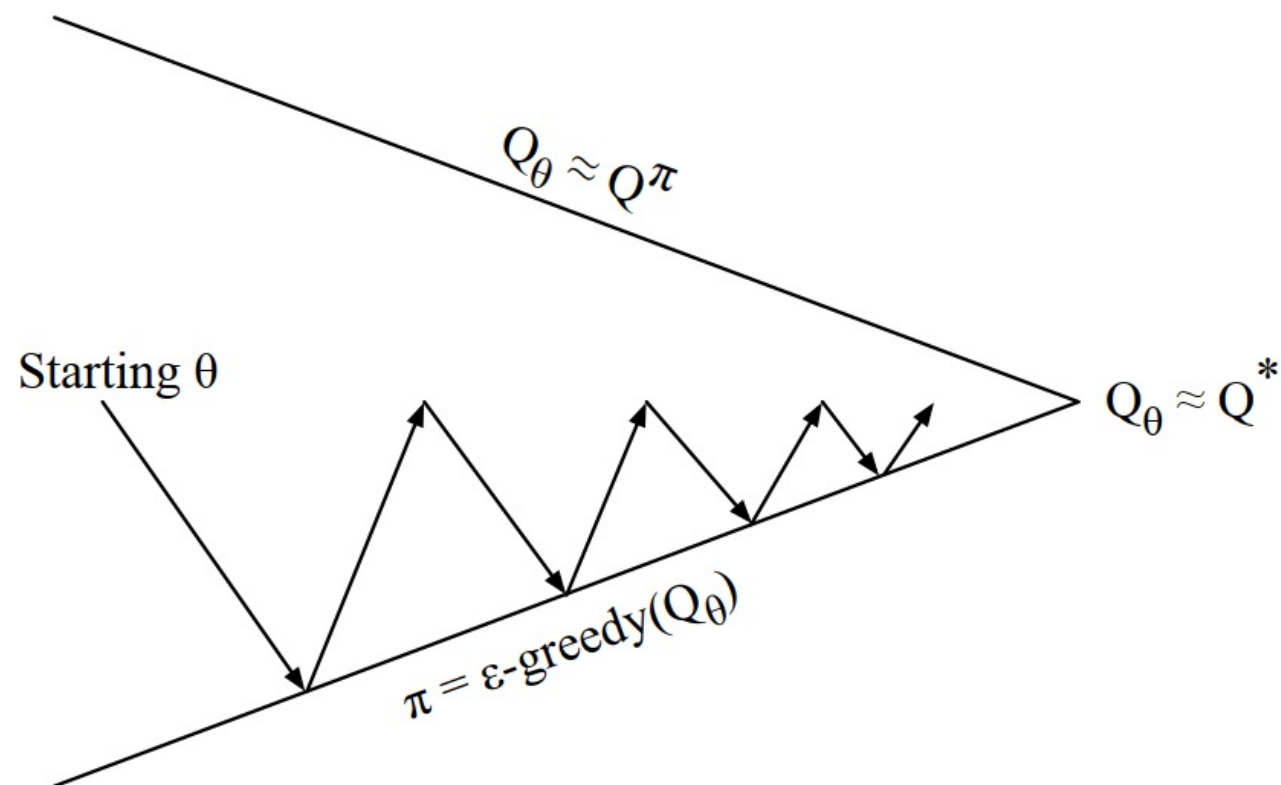
$$\Delta \mathbf{w}_t = \alpha(\mathbf{R}_{t+1} + \gamma v_{\mathbf{w}}(S_{t+1}) - v_{\mathbf{w}}(S_t)) \nabla_{\mathbf{w}} v_{\mathbf{w}}(S_t)$$

Monte-Carlo with Value Function Approximation

- The return G_t is an **unbiased** sample of $v_\pi(s)$
- Can therefore apply “supervised learning” to (online) “training data”:
$$\{(S_0, G_0), \dots, (S_t, G_t)\}$$

Control with value-function approximation

Control with Value Function Approximation

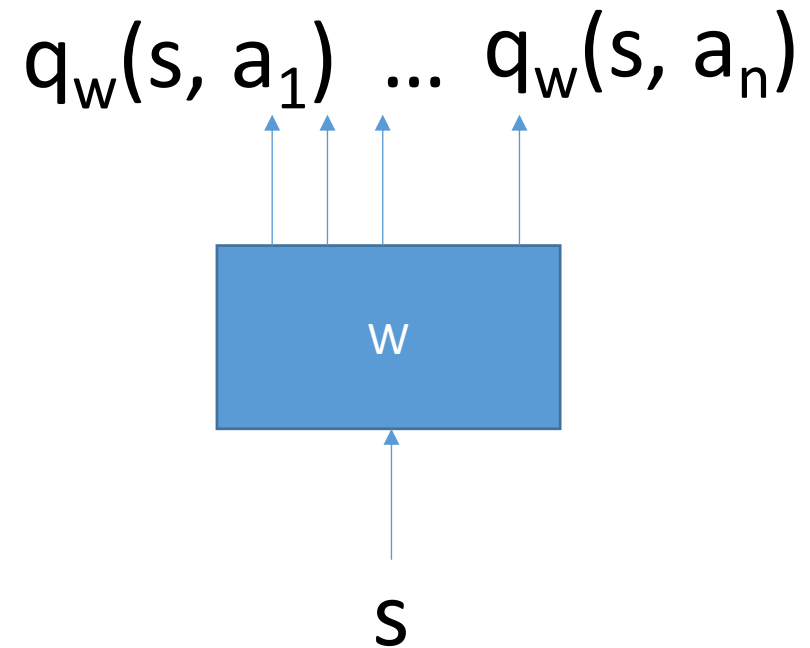
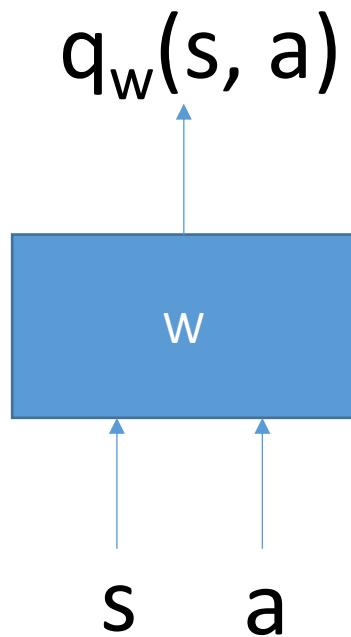


Policy evaluation **Approximate** policy evaluation, $q_w \approx q_\pi$

Policy improvement E.g., ϵ -greedy policy improvement

Action-Value Function Approximation

- Approximate the action-value function $q_w(s, a) \approx q_\pi(s, a)$



Action-Value Function Approximation

- Should we use action-in, or action-out?
 - Action in: $q_w(s, a) = w^T x(s, a)$
 - Action out: $q_w(s) = Wx(s)$ such that $q_w(s, a) = q_w(s)[a]$
- One **incorporates a** in feature learning, the other uses the **same features for all a's**
- Unclear which is better in general
- If we want to use **continuous actions**, **action-in** is easier
- For **(small) discrete action** spaces, **action-out** is common (e.g., DQN)

Deep reinforcement learning

Recall: Model free control

- Similar to policy evaluation, true state-action value function for a state is unknown and so substitute a target value
- In Monte Carlo methods, use a return G_t as a substitute target

$$\Delta \mathbf{w} = \alpha(G_t - \hat{Q}(s_t, a_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t; \mathbf{w})$$

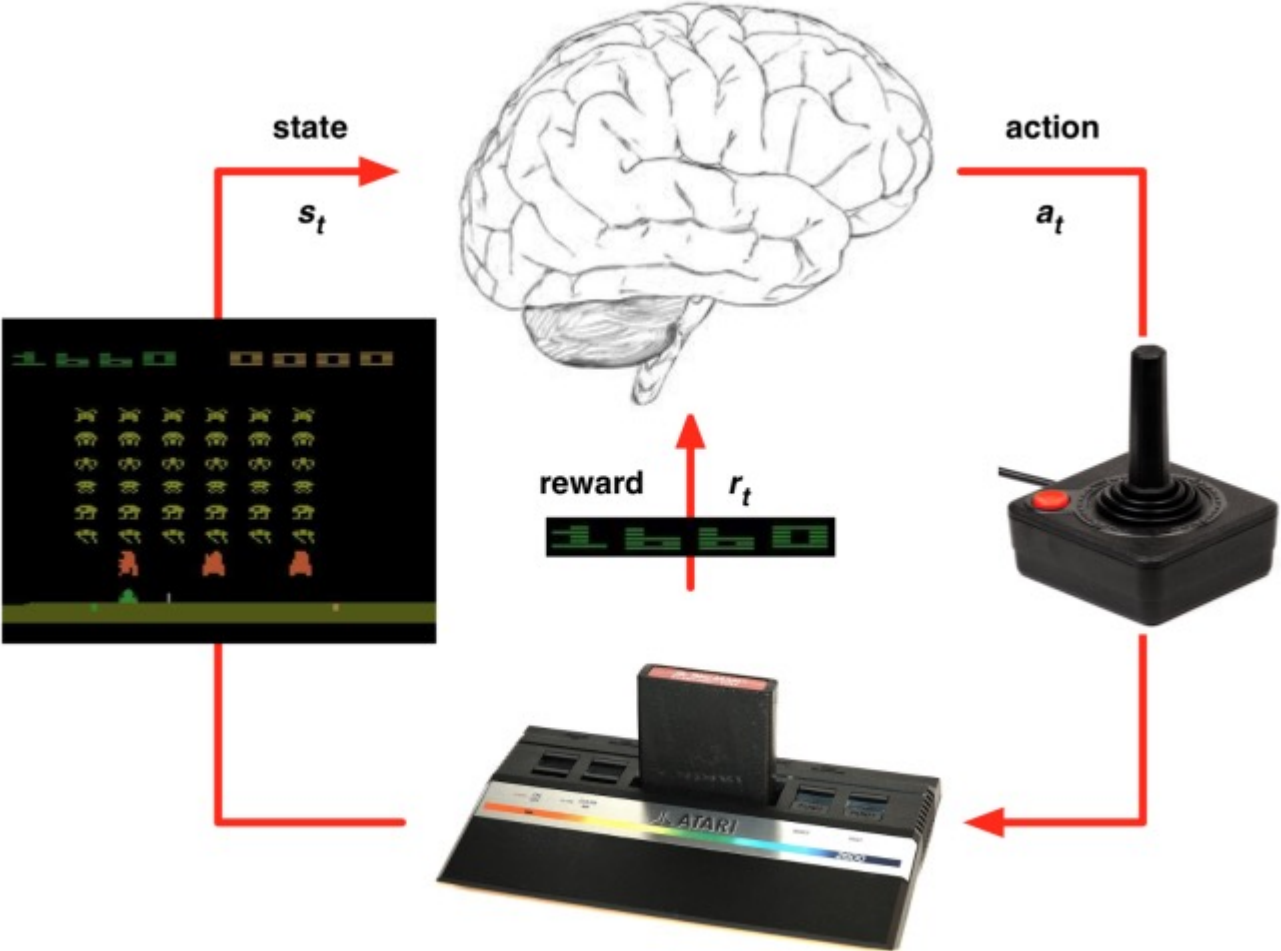
- For SARSA instead use a TD target

$$\Delta \mathbf{w} = \alpha(r + \gamma \hat{Q}(s_{t+1}, a_{t+1}; \mathbf{w}) - \hat{Q}(s_t, a_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t; \mathbf{w})$$

- For Q-learning

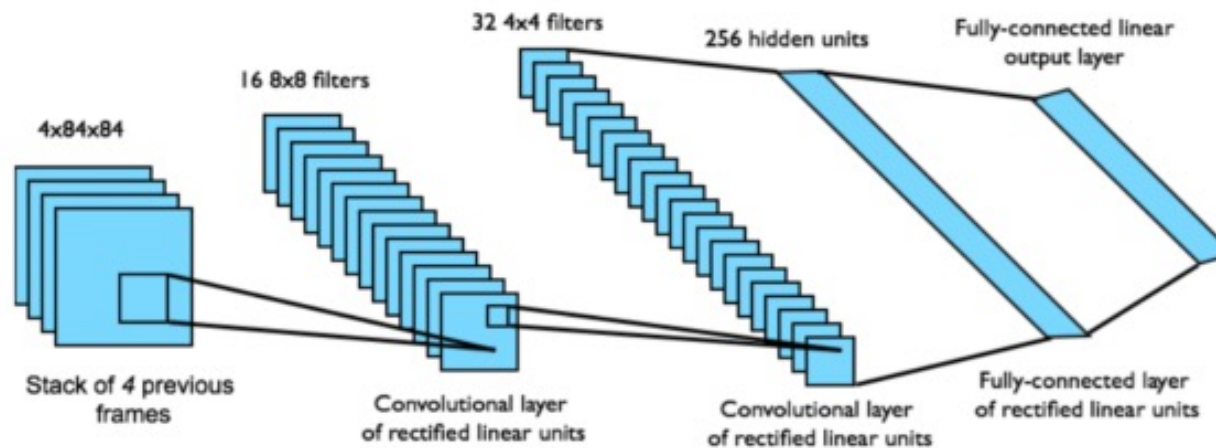
$$\Delta \mathbf{w} = \alpha(r + \gamma \max_a \hat{Q}(s_{t+1}, a; \mathbf{w}) - \hat{Q}(s_t, a_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t; \mathbf{w})$$

Doing deep RL in Atari



DQNs in Atari

- End-to-end learning of values $Q(s, a)$ from pixels s
- Input state s is stack of raw pixels from **last 4 frames**
- Output is $Q(s, a)$ for **18 joystick/button positions**
- Reward is **change in score** for that step
- Network architecture and hyperparameters fixed across all games

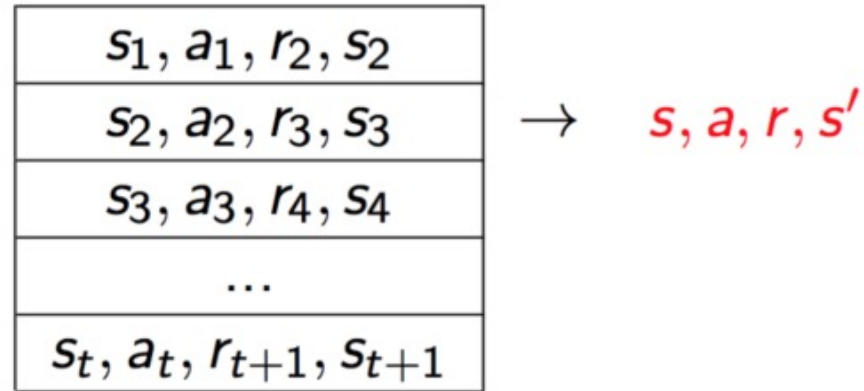


DQNs in Atari

- Q-learning **converges** to the optimal $Q^*(s, a)$ using table lookup representation.
- But Q-learning with Value Function Approximation can diverge
- Two of the issues causing problems:
 - **Correlations** between samples (non-iid training)
 - **Non-stationary** targets
- Deep Q-learning (DQN) addresses these challenges by
 - **Experience replay**
 - **Fixed Q-targets**

DQNs: Experience Replay

- To help remove correlations, store dataset (called a replay buffer) D from prior experience



- To perform experience replay, repeat the following:
 - $(s, a, r, s') \sim D$: sample an experience tuple from the dataset
 - Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; w)$
 - Use stochastic gradient descent to update the network weights

$$\Delta \mathbf{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

Problem

Can treat the target as a constant scalar, but the **weights will get updated** on the next round, **changing the target value**

DQNs: Fixed Q-Targets

- To help improve **stability**, **fix** the **target weights** used in the target calculation for multiple updates
- Target network uses a **different set of weights** than the weights being updated
- Let parameters w^- be the set of weights used in the target, and w be the weights that are being updated
- Slight change to computation of target value:
 - $(s, a, r, s') \sim D$: sample an experience tuple from the dataset
 - Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; w^-)$
 - Use stochastic gradient descent to update the network weights

$$\Delta \mathbf{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}^-) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

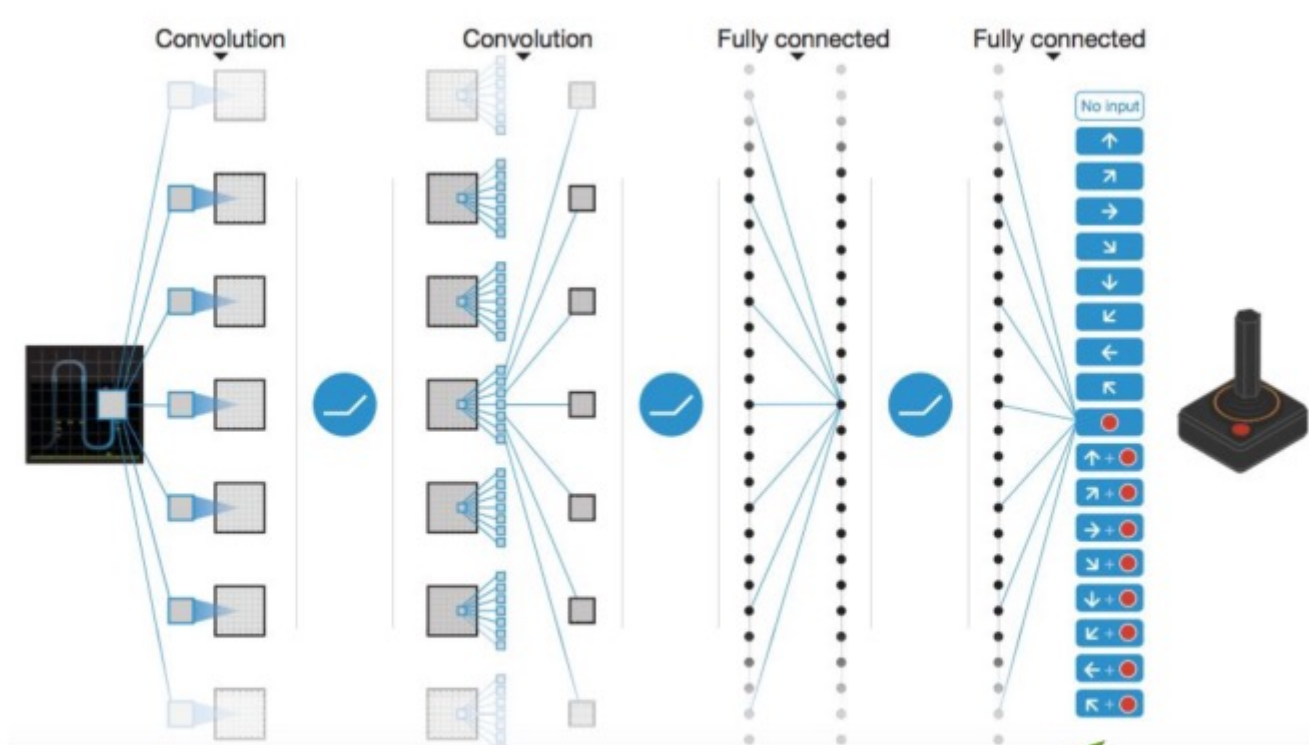
DQN Algorithm

```
1: Input  $C, \alpha, D = \{\}$ , Initialize  $\mathbf{w}, \mathbf{w}^- = \mathbf{w}, t = 0$ 
2: Get initial state  $s_0$ 
3: loop
4:   Sample action  $a_t$  given  $\epsilon$ -greedy policy for current  $\hat{Q}(s_t, a; \mathbf{w})$ 
5:   Observe reward  $r_t$  and next state  $s_{t+1}$ 
6:   Store transition  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $D$ 
7:   Sample random minibatch of tuples  $(s_i, a_i, r_i, s_{i+1})$  from  $D$ 
8:   for  $j$  in minibatch do
9:     if episode terminated at step  $i + 1$  then
10:       $y_i = r_i$ 
11:     else
12:       $y_i = r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a'; \mathbf{w}^-)$ 
13:     end if
14:     Do gradient descent step on  $(y_i - \hat{Q}(s_i, a_i; \mathbf{w}))^2$  for parameters  $\mathbf{w}$ :  $\Delta \mathbf{w} = \alpha(y_i - \hat{Q}(s_i, a_i; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_i, a_i; \mathbf{w})$ 
15:   end for
16:    $t = t + 1$ 
17:   if  $\text{mod}(t, C) == 0$  then
18:      $\mathbf{w}^- \leftarrow \mathbf{w}$ 
19:   end if
20: end loop
```

DQN Summary

- DQN uses experience replay and fixed Q-targets
- Store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory D
- Sample random mini-batch of transitions (s, a, r, s') from D
- Compute Q-learning targets w.r.t. old, fixed parameters w^-
- Optimizes MSE between Q-network and Q-learning targets
- Uses stochastic gradient descent

DQN



1 network, outputs Q value for each action

Results

Game	Linear	Deep Network	DQN w/ fixed Q	DQN w/ replay	DQN w/replay and fixed Q
Breakout	3	3	10	241	317
Enduro	62	29	141	831	1006
River Raid	2345	1453	2868	4102	7447
Seaquest	656	275	1003	823	2894
Space Invaders	301	302	373	826	1089

Policy Gradient

The Goal of Reinforcement Learning

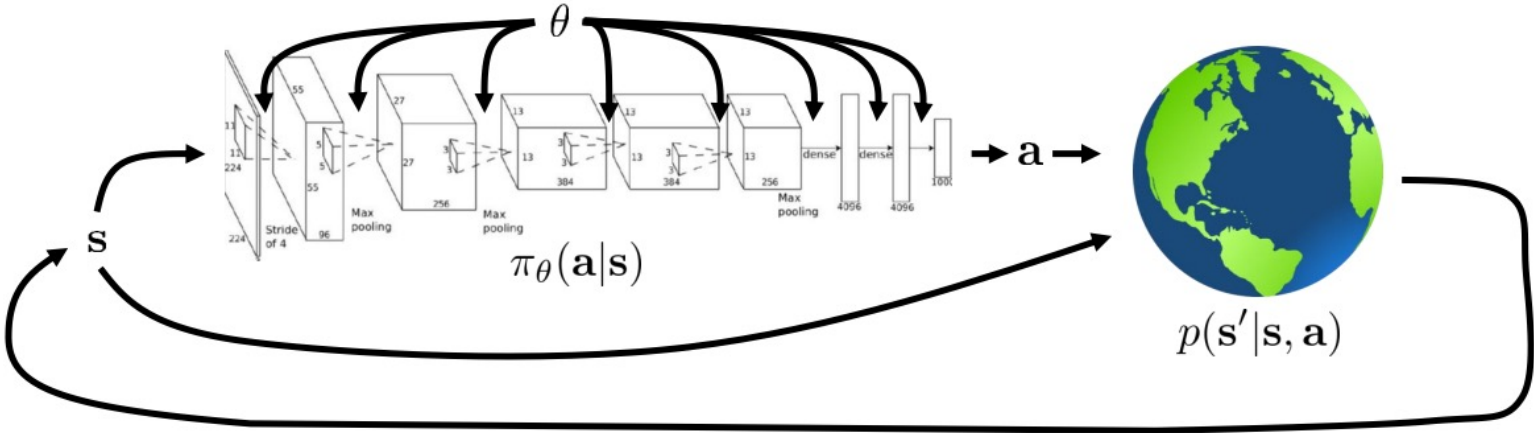
$$J(\theta) = E \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$$s_1 \sim p(s_1)$$

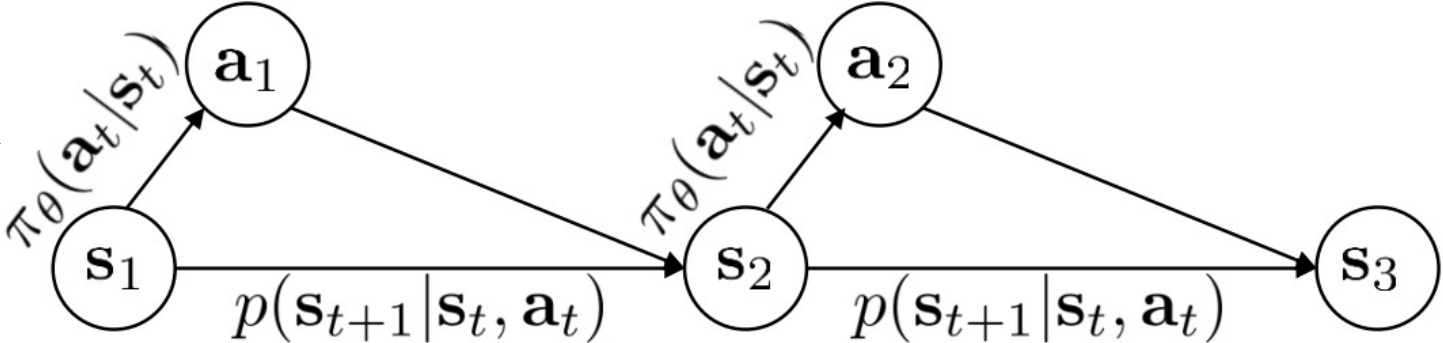
$$a_t \sim \pi(\cdot | s_t)$$

$$s_{t+1} \sim p(\cdot | s_t, a_t)$$

Trajectory Probability

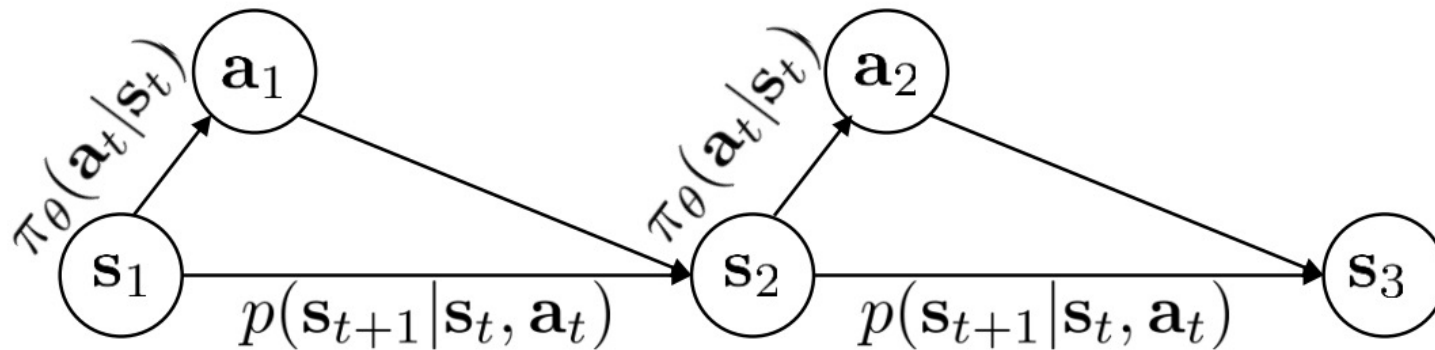


stochastic policy



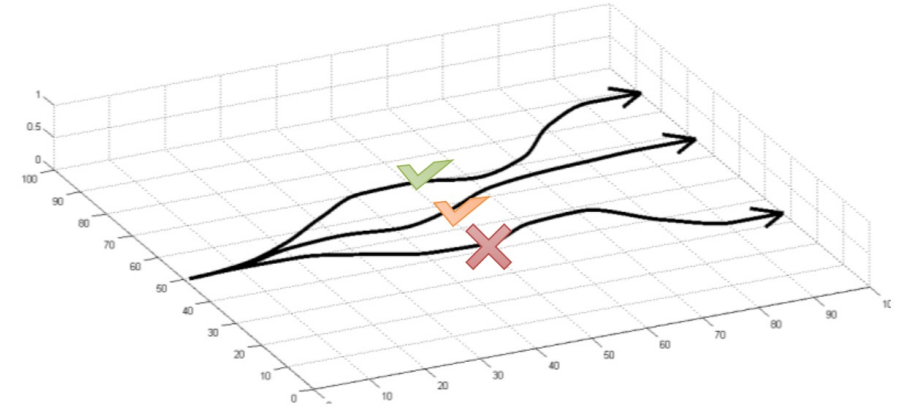
Trajectory Probability

$$\underbrace{p_{\theta}(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T)}_{p_{\theta}(\tau)} = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$



Evaluating the Objective

$$\theta^* = \arg \max_{\theta} \underbrace{E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]}_{J(\theta)}$$



Monte Carlo estimation for objective function:

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

sum over samples from π_{θ}

Direct Policy Differentiation

$$\theta^* = \arg \max_{\theta} \underbrace{E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]}_{J(\theta)}$$

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\underbrace{r(\tau)}_{\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t)} \right] = \int p_{\theta}(\tau) r(\tau) d\tau$$

$$\nabla_{\theta} J(\theta) = \int \underbrace{\nabla_{\theta} p_{\theta}(\tau)} r(\tau) d\tau = \int \underbrace{p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)} r(\tau) d\tau = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$$

a convenient identity

$$\underline{p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)} = p_{\theta}(\tau) \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} = \underline{\nabla_{\theta} p_{\theta}(\tau)}$$

Direct Policy Differentiation

$$\theta^* = \arg \max_{\theta} J(\theta)$$

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [r(\tau)]$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$$

log of both sides

$$p_{\theta}(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\log p_{\theta}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\nabla_{\theta} \left[\cancel{\log p(\mathbf{s}_1)} + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \cancel{\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} \right]$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

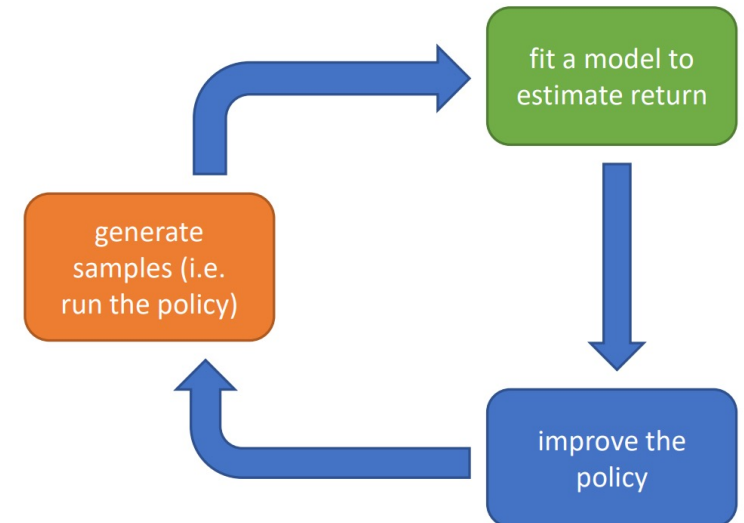
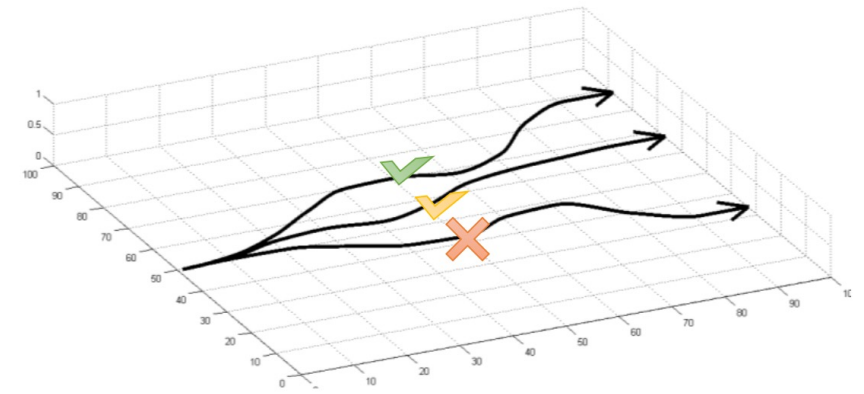
Evaluating the Policy Gradient

$$\text{recall: } J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

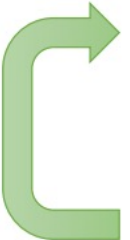
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

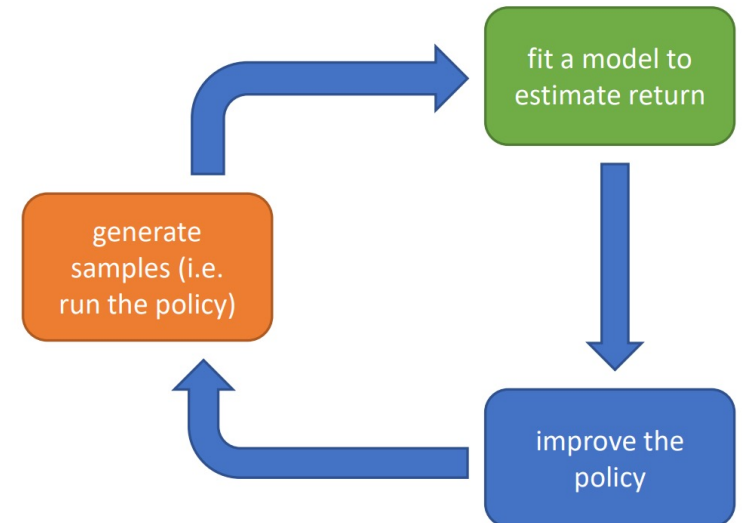
$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$



REINFORCE Algorithm

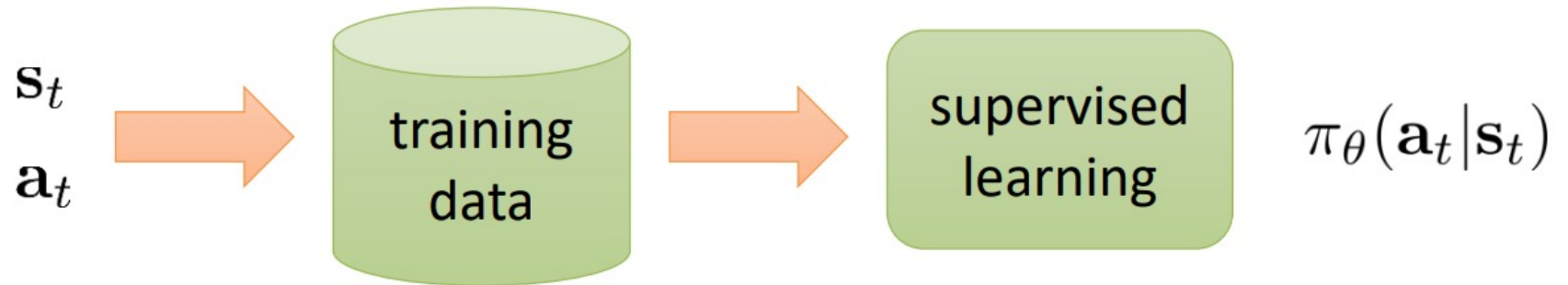
REINFORCE algorithm:

- 
1. sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ (run the policy)
 2. $\nabla_\theta J(\theta) \approx \sum_i (\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i|\mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i))$
 3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$



Imitation Learning: Behavioral Cloning

Directly learns a policy by **using supervised learning** on observation-action pairs **from expert demonstrations**.



maximum likelihood:
$$\nabla_{\theta} J_{\text{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right)$$

Policy Gradient Vs Maximum Likelihood

policy gradient:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

maximum likelihood:
$$\nabla_{\theta} J_{\text{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right)$$

