

# **TTIC 31230, Fundamentals of Deep Learning**

David McAllester, April 2017

## **Deep Reinforcement Learning**

# Definition of Reinforcement Learning

RL is defined by the following properties:

- An environment with **state**.
- State changes are influenced by **sequential decisions**.
- Reward (or loss) depends on **making decisions** that lead to **desirable states**.

# Reinforcement Learning Examples

- Board games (chess or go)
- Atari Games (pong)
- Robot control (driving)
- Dialog
- Life

# Policies

A policy is a way of behaving.

Formally, a (nondeterministic) policy maps a state to a probability distribution over actions.

$\pi(a_t|s_t)$  probability of action  $a_t$  in state  $s_t$

## Imitation Learning

Construct a dataset of state-action pairs  $(s, a)$  from experts.

Define stochastic policy  $\pi_{\Theta}(s)$ .

$$\Theta^* = \operatorname{argmin}_{\Theta} E_{(s,a)} - \ln \pi_{\Theta}(a \mid s)$$

This is just cross-entropy loss.

## Dangers of Imperfect Imitation Learning

Perfect imitation learning would reproduce expert behavior.

Imitation learning is **off-policy** — the state distribution in the training data is different from that defined by the policy being learned.

Imitating experts can be dangerous. “Don’t try this at home”.

On the flip side, imitating experts will never exceed expert performance (consider Alpha-Zero).

## Markov Decision Processes (MDPs)

For an RL problem we work with an action policy  $\pi$

$s_t$  is the state at time  $t$

$r_t$  is the reward at time  $t$

$a_t$  is the action taken at time  $t$ .

$r_t = R(s_t, a_t)$  reward at time  $t$

$T(s_{t+1}|s_t, a_t)$  state transition probability

The state space, action space,  $R$  and  $T$  define a **Markov Decision Process** or **MDP**.

## Optimizing Reward

In RL we maximize reward rather than minimize loss.

$$\pi^* = \operatorname{argmax}_{\pi} R(\pi)$$

$$R(\pi) = \mathbb{E} \left[ \sum_{t=0}^T r_t \right] \quad \text{episodic reward (go)}$$

$$\text{or } \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right] \quad \text{discounted reward}$$

$$\text{or } \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T r_t \quad \text{asymptotic average reward (driving)}$$



# The Value Function

For discounted reward:

$$V^\pi(s) = \mathbb{E} \left[ \sum_t \gamma^t r_t \mid \pi, s_0 = s \right]$$

$$V^*(s) = \max_{\pi} V^\pi(s)$$

$$\pi^*(a|s) = \operatorname{argmax}_a E_{s' \sim T(s'|s,a)} V^*(s')$$

$$V^*(s) = \max_a R(s, a) + \gamma E_{s' \sim T(\cdot|s,a)} V^*(s')$$

## Value Iteration

Suppose the state space and action space are finite.

In that case we can do value iteration.

$$V_0(s) = 0$$

$$V_{i+1}(s) = \max_a R(s, a) + \gamma E_{s' \sim T(\cdot|s,a)} V_i(s')$$

If all rewards are non-negative then

$$V_{i+1}(s) \geq V_i(s) \quad V_i(s) \leq V^*(s)$$

## Value Iteration

Theorem: **For discounted reward**

$$\lim_{i \rightarrow \infty} V_i(s) = V^*(s)$$

Proof:

$$\begin{aligned} \Delta &\doteq \max_s V^*(s) - V_\infty(s) \\ &= \max_s \left( \begin{aligned} &\max_a R(s, a) + E_{s'|a} \gamma V^*(s') \\ &- \max_a R(s, a) + E_{s'|a} \gamma V_\infty(s') \end{aligned} \right) \\ &\leq \gamma \Delta \end{aligned}$$

## The $Q$ Function

For discounted reward:

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

$$Q^*(s, a) = \max_{\pi} Q^\pi(s, a)$$

$$\pi^*(a|s) = \operatorname{argmax}_a Q^*(s, a)$$

$$Q^*(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim T(\cdot|s, a)} \max_{a'} Q^*(s', a')$$

## $Q$ -Learning

We will assume a parameterized  $Q$  function  $Q_{\Theta}(s, a)$ .

Define the **Bellman error**

$$\left( Q_{\Theta}(s, a) - \left( R(s, a) + \gamma E_{s' \sim S(\cdot|s, a)} \max_{a'} Q_{\Theta}(s', a') \right) \right)^2$$

**Algorithm:** run the policy  $\operatorname{argmax}_a Q_{\Theta}(s_t, a)$  and repeat

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} (Q_{\Theta}(s_t, a_t) - (r_t + \gamma \max_a Q_{\Theta}(s_{t+1}, a)))^2$$

## The Theorem

**Theorem: For discounted reward**, if the Bellman error is zero then the induced policy is optimal.

This  $Q$ -learning theorem is similar in strength to the contrastive divergence theorem (gradient is zero at solution).

Weaker than the pseudo-likelihood theorem (optimization problem is correct assuming universal expressiveness).

## Issues with $Q$ -Learning

Problem 1: Nearby states in the same run are highly correlated.

Problem 2: SGD on Bellman error tends to be unstable. (Weakness of the  $Q$ -learning theorem?)

To address these problems we can use a **replay buffer**.

## Using a Replay Buffer

We use a replay buffer of tuples  $(s_t, a_t, r_t, s_{t+1})$ .

Repeat:

1. Run the policy  $\operatorname{argmax}_a Q_{\Theta}(s, a)$  to add tuples to the replay buffer. Remove oldest tuples to maintain a maximum buffer size.
2.  $\Psi = \Theta$
3. for  $N$  times select a random element of the replay buffer and do

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} (Q_{\Theta}(s_t, a_t) - (r_t + \gamma \operatorname{argmax}_a Q_{\Psi}(s_{t+1}, a)))^2$$



## Replay is Off-Policy

Note that the replay buffer is (slightly) **off-policy**. This seems to be important for stability.

Seems related to the issue of stochastic vs. deterministic policies. More on this later.

## Multi-Step $Q$ -learning

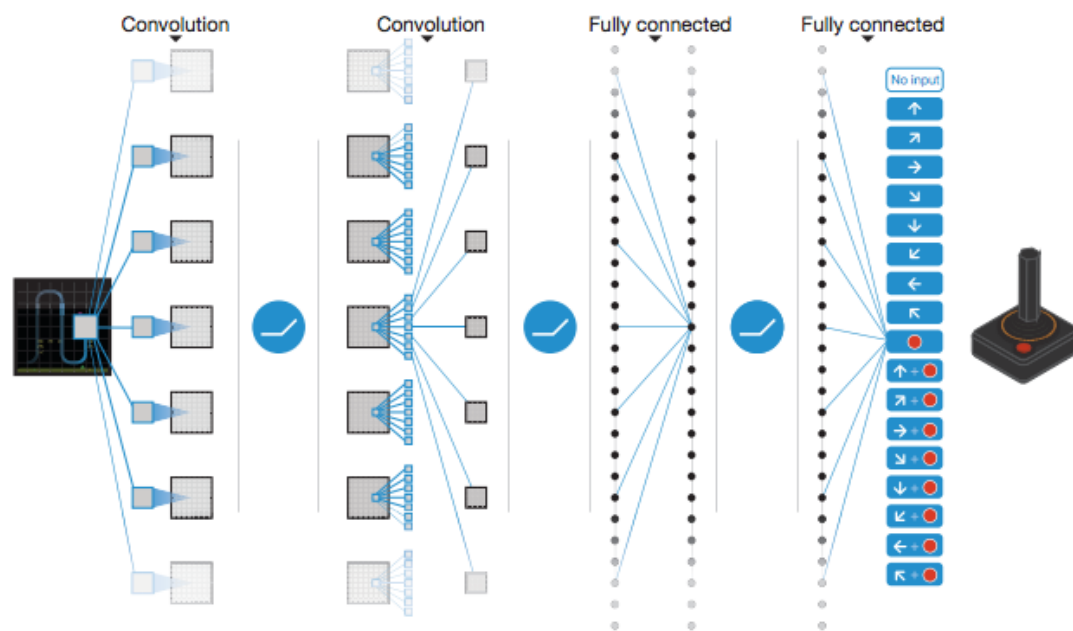
$$\Theta \leftarrow \sum_t \nabla_{\Theta} \left( Q_{\Theta}(s_t, a_t) - \sum_{\delta=0}^K \gamma^{\delta} r_{(t+\delta)} \right)^2$$

## Deep $Q$ -Learning (DQN)

Human-level control through deep reinforcement learning, Mnih et al., Nature, 2015. (Deep Mind)

# Deep Q-Networks

We consider a **Deep Q-network** — a deep network with parameters  $\Theta$  and computing a Q-value  $Q_{\Theta}(s, a)$ .



## Watch The Video

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

## Asynchronous Q-Learning (Simplified)

No replay buffer. Many asynchronous threads each repeating:

$$\tilde{\Theta} = \Theta \text{ (retrieve global } \Theta \text{)}$$

using policy  $\operatorname{argmax}_a Q_{\tilde{\Theta}}(s, a)$  compute

$$s_t, a_t, r_t, \dots, s_{t+K}, a_{t+K}, r_{t+K}$$

$$R_i = \sum_{\delta=0}^{t+K-i} \gamma^{i+\delta} r_{(i+\delta)}$$

Update global  $\Theta$ :

$$\Theta \leftarrow \Theta + \eta \sum_{i=t}^{t+K} \nabla_{\tilde{\Theta}} (Q_{\tilde{\Theta}}(s_i, a_i) - R_i)^2$$

# Policy Gradient

We assume a parameterized policy  $\pi_{\Phi}(a|s)$ .

$\pi_{\Phi}(a|s)$  is normalized while  $Q_{\Theta}(s, a)$  is not.

$$\Phi \ += \ \eta \nabla_{\Phi} R(\Phi)$$

## Policy Gradient Theorem (Episodic Case)

$$\mathbb{E}[R \mid \Phi] = \sum_{s_0, a_0, s_1, a_1, \dots, s_T, a_T} P(s_0, a_0, s_1, a_1, \dots, s_T, a_T) R$$

$$\begin{aligned} \nabla_{\Phi} P(\dots)R &= P(S_0) \nabla_{\Phi} \pi(a_0) P(s_1) \pi(a_1) \cdots P(s_T) \pi(a_T) R \\ &\quad + P(S_0) \pi(a_0) P(s_1) \nabla_{\Phi} \pi(a_1) \cdots P(s_T) \pi(a_T) R \\ &\quad \vdots \\ &\quad + P(S_0) \pi(a_0) P(s_1) \pi(a_1) \cdots P(s_T) \nabla_{\Phi} \pi(a_T) R \\ &= P(\dots) \left( \sum_i \frac{\nabla_{\Phi} \pi_{\Phi}(a_i)}{\pi_{\Phi}(a_i)} \right) R \end{aligned}$$



## Policy Gradient Theorem (Episodic Case)

$$\begin{aligned}
 \nabla_{\Phi} P(\dots)R &= P(S_0)\nabla_{\Phi} \pi(a_0)P(s_1)\pi(a_1)\cdots P(s_T)\pi(a_T) R \\
 &\quad + P(S_0)\pi(a_0)P(s_1)\nabla_{\Phi} \pi(a_1)\cdots P(s_T)\pi(a_T) R \\
 &\quad \vdots \\
 &\quad + P(S_0)\pi(a_0)P(s_1)\pi(a_1)\cdots P(s_T)\nabla_{\Phi} \pi(a_T) R \\
 &= P(\dots) \left( \sum_i \frac{\nabla_{\Phi} \pi_{\Phi}(a_i)}{\pi_{\Phi}(a_i)} \right) R \\
 \nabla_{\Phi} \mathbb{E}[R \mid \Phi] &= \mathbb{E} \left[ R \sum_t \nabla_{\Phi} \ln \pi_{\Phi}(a_t|s_t) \right]
 \end{aligned}$$

## Policy Gradient Theorem

$$\begin{aligned} & \nabla_{\Phi} \mathbb{E}[R \mid \Phi] \\ &= \mathbb{E} \left[ R \sum_t \nabla_{\Phi} \ln \pi_{\Phi}(a_t | s_t) \right] \\ &= E \sum_{t_1, t_2} r_{t_1} \nabla_{\Phi} \ln \pi_{\Phi}(a_{t_2} | s_{t_2}) \\ &= E \sum_{t_1 < t_2} r_{t_1} \nabla_{\Phi} \ln \pi_{\Phi}(a_{t_2} | s_{t_2}) + E \sum_{t_1 \geq t_2} r_{t_1} \nabla_{\Phi} \ln \pi_{\Phi}(a_{t_1} | s_{t_2}) \end{aligned}$$

# Policy Gradient Theorem

$$\nabla_{\Phi} E[R \mid \Phi]$$

$$= E \sum_{t_1 < t_2} r_{t_1} \nabla_{\Phi} \ln \pi_{\Phi}(a_{t_2} | s_{t_2}) + E \sum_{t_1 \geq t_2} r_{t_1} \nabla_{\Phi} \ln \pi_{\Phi}(a_{t_1} | s_{t_2})$$

$$= \sum_{t_1 < t_2} E_{s_{t_1}} (E r_{t_1} | s_{t_1}) E_{s_{t_2} | s_{t_1}} \sum_{a_{t_2}} \pi_{\Phi}(a_{t_2} | s_{t_2}) \nabla_{\Phi} \ln \pi_{\Phi}(a_{t_2} | s_{t_2}) + \dots$$

$$\sum_{t_1 < t_2} E_{s_{t_1}} (E r_{t_1} | s_{t_1}) E_{s_{t_2} | s_{t_1}} \sum_{a_{t_2}} \nabla_{\Phi} \pi_{\Phi}(a_{t_2} | s_{t_2}) = 0$$

## Policy Gradient Theorem

$$\begin{aligned} & \nabla_{\Phi} E[R \mid \Phi] \\ &= E \sum_{t_1 \leq t_2} \nabla_{\Phi} \ln \pi_{\Phi}(a_{t_1} | s_1) r_{t_2} \\ &= E \sum_t \nabla_{\Phi} \ln \pi_{\Phi}(a_t | s_t) R_{\geq t} \end{aligned}$$

Sampling runs and computing the above sum over  $t$  is Williams' REINFORCE algorithm.

# Optimizing Discrete Decisions

The REINFORCE algorithm is used generally for non-differentiable loss functions.

For example error rate and BLEU score are non-differentiable — they are defined on the result of discrete decisions.

## Policy Gradient Theorem

$$\begin{aligned}\nabla_{\Phi} E[R \mid \Phi] &= E \sum_t \nabla_{\Phi} \ln \pi_{\Phi}(a_t|s_t) R_{\geq t} \\ &= \sum_s \rho(s) \sum_a (\nabla_{\Phi} \pi_{\Phi}(a|s)) Q^{\pi_{\Phi}}(s, a)\end{aligned}$$

$\rho(s)$  is the expected number of occurrences of  $s$

$$Q^{\pi}(s, a) = E_{\pi} \sum_t r_t \mid s_0 = s, a_0 = a$$

## The Actor-Critic Algorithm

$$\begin{aligned}\nabla_{\Phi} E[R \mid \Phi] &= E \sum_t \nabla_{\Phi} \ln \pi_{\Phi}(a_t|s_t) Q^{\pi_{\Phi}}(s_t, a_t) \\ &\approx E \sum_t \nabla_{\Phi} \ln \pi_{\Phi}(a_t|s_t) Q_{\Psi}(s_t, a_t)\end{aligned}$$

We can reduce variance by using an estimator  $Q_{\Psi}(s, a)$  for  $Q^{\pi_{\Phi}}(s, a)$ .

$\pi_{\Phi}$  is the “actor” and  $Q_{\Psi}$  is the “critic”.

## The Actor-Critic Algorithm

$$\nabla_{\Phi} E[R \mid \Phi] \approx E \sum_t \nabla_{\Phi} \ln \pi_{\Phi}(a_t | s_t) Q_{\Psi}(s_t, a_t)$$

$$\begin{aligned} \Phi & += \eta_1 \sum_t (\nabla_{\tilde{\Phi}} \ln \pi_{\tilde{\Phi}}(a_i | s_i)) Q_{\Psi}(s_t, a_t) \\ \Psi & -= \eta_2 \sum_t \nabla_{\tilde{\Psi}} (Q_{\Psi}(s_t, a_t) - R_{\geq t})^2 \end{aligned}$$

$\pi_{\Phi}$  is the “actor” and  $Q_{\Psi}$  is the “critic”.



## Advantage-Actor-Critic Theorem

$$\nabla_{\Phi} R(\Phi) = \sum_s \rho(s) \sum_a (\nabla_{\Phi} \pi_{\Phi}(a|s)) (Q^{\pi_{\Phi}}(s, a) - V^{\pi_{\Phi}}(s))$$

$$V^{\pi}(s) = E_{a \sim \pi(\cdot|s)} Q(s, a)$$

## **Asynchronous Advantage Actor-Critic (A3C)**

Asynchronous Methods for Deep Reinforcement Learning, Mnih et al., Arxiv, 2016 (Deep Mind)

## Asynchronous Advantage Actor-Critic (A3C)

$\tilde{\Phi} = \Phi; \tilde{\Psi} = \Psi$  (retrieve global  $\Phi$  and  $\Psi$ )

using policy  $\pi_{\tilde{\Phi}}$  compute  $s_t, a_t, r_t, \dots, s_{t+K}, a_{t+K}, r_{t+K}$

$$R_i = \sum_{\delta=0}^{t+K-i} \gamma^{i+\delta} r_{(i+\delta)}$$

Update global  $\Phi$  and  $\Psi$

$$\begin{aligned} \Phi & += \eta \sum_{i=t}^{t+K} \left( \nabla_{\tilde{\Phi}} \ln \pi_{\tilde{\Phi}}(a_i | s_i) \right) (R_i - V_{\tilde{\Psi}}(s_i)) \\ \Psi & -= \eta \sum_{i=t}^{t+K} \nabla_{\tilde{\Psi}} (V_{\tilde{\Psi}}(s_i) - R_i)^2 \end{aligned}$$

## Issue: Policies must be Exploratory

The optimal policy is deterministic —  $a(s) = \operatorname{argmax}_a Q(s, a)$ .

However, a deterministic policy never samples alternative actions.

Typically one forces a random action some small fraction of the time.

## **Issue: Discounted Reward**

DQN and A3C use discounted reward on episodic or long term problems.

Presumably this is because actions have near term consequences.

This should be properly handled in the mathematics.

## Observation: Continuous Actions are Differentiable

In problems like controlling an inverted pendulum, or robot control generally, a continuous loss can be defined and the gradient of loss of with respect to a deterministic policy exists.

## More Videos

<https://www.youtube.com/watch?v=g59nSURxYgk>

<https://www.youtube.com/watch?v=rAai4QzcYbs>

**END**