

---

# RL Algorithms

---

Kaige Yang  
University College London  
Kaige.yang.11@ucl.ac.uk

## Contents

<b>1</b>	<b>Overview of Algorithms</b>	<b>3</b>
<b>2</b>	<b>RL Problem</b>	<b>3</b>
<b>3</b>	<b>How to solve RL Problem?</b>	<b>4</b>
<b>4</b>	<b>Advantage Actor-Critic</b>	<b>6</b>
<b>5</b>	<b>Value-based and Policy-based Models</b>	<b>7</b>
<b>6</b>	<b>Model-free and Model-based Models</b>	<b>7</b>
<b>7</b>	<b>Algorithms</b>	<b>7</b>
7.1	TD-Learning . . . . .	7
7.2	Q-Learning . . . . .	8
7.3	DQN . . . . .	9
7.4	{N-step, Double, Noisy, Dueling, Categorical}-DQN . . . . .	10
7.5	REINFORCE . . . . .	11
7.6	Actor-Critic . . . . .	12
7.7	A3C . . . . .	13
7.8	A2C . . . . .	14
7.9	ACER . . . . .	14
7.10	TRPO . . . . .	14
7.11	PPO . . . . .	16
7.12	DPG . . . . .	18
7.13	DDPG . . . . .	19
7.14	TD3 . . . . .	21
7.15	Soft Q-Learning . . . . .	22
7.16	SAC . . . . .	24
7.17	State-of-The-Art Algorithms . . . . .	26



## Abstract

This report summary the detail of state-of-the-art RL algorithms. The goal is to provide a holistic view on the recent progress in RL. The content of this report is a summary of many learning resources.

## 1 Overview of Algorithms

name	value	policy	off-policy	on-policy	dis-state	cont-state	dis-action	con-action
Q-learning	×		×		×		×	
REINFORCE		×		×	×		×	
Actor-Critic		×		×	×	×	×	×
DQN	×		×		×	×	×	
DDPG		×	×		×	×		×
PPO		×		×	×	×	×	
SAC		×	×		×	×		×
TD3		×	×		×	×		×
A3C		×		×	×	×	×	×

Table 1: RL Algorithms

## 2 RL Problem

**Markov Decision Processes** (MDPs) are a tool for modeling sequential decision-making problems where an agent interacts with an environment in a sequential fashion. An MDP, denoted by  $\mathcal{M}$ , is defined as  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$ . The state space  $\mathcal{S}$ , the action space  $\mathcal{A}$ , the transition probability kernel  $\mathcal{P}$  where  $\mathcal{P}(s'|s, a)$  gives the probability that the next state  $s'$  given the current state  $s$  and taken action  $a$ . The reward function  $r(s, a)$  gives the immediate reward upon taking action  $a$  in state  $s$ .  $\gamma \in [0, 1)$  is the discounted factor.

Given an MDP  $\mathcal{M}$ , the interaction between the agent and environment happens as follows: Let  $t \in \mathbb{N}$  denote the current time, let  $s_t \in \mathcal{S}$  be the current state and  $a_t$  be the action chosen by the agent at time  $t$ . Upon the action is chosen, the state of the environment moves to  $s'$  following  $\mathcal{P}(s'|s, a)$ . The agent receives a reward  $r(s, a)$ . Then, the agent observes the next state  $s'$  and chooses next action  $a'$  and the process repeated. A trajectory  $\{(s_0, a_0, r_1), \dots, (s_{t-1}, a_{t-1}, r_t)\}$  summaries the interaction history up to time  $t$ .

The **goal** of the agent is to find a policy  $\pi$  so as to maximize the expected total discounted reward (expected discounted cumulative reward). At each time  $t$ ,  $\pi(a|s_t)$  gives the probability of action  $a$  being selecting given the current state  $s_t$ . The expected discounted cumulative reward is defined as

$$R = \sum_{t=0}^{\infty} \gamma^t r_{t+1} \quad (1)$$

The goal is to find optimal policy which maximizes  $R$ . Namely,

$$\pi^* = \arg \max_{\pi} \sum_{t=0}^{\infty} \gamma^t r_{t+1} \quad (2)$$

The state-value function,  $V : \mathcal{S} \rightarrow \mathbb{R}$ , is defined by

$$V(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s \right], s \in \mathcal{S}. \quad (3)$$

The action-value function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ .

$$Q(s, a) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s, a_0 = a \right], s \in \mathcal{S}, a \in \mathcal{A}. \quad (4)$$

It is known that both  $V(s)$  and  $Q(s, a)$  satisfy **Bellman consistency equation**:

$$V(s) = \mathbb{E} \left[ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s'|s, a) V(s') \right], \quad s, s' \in \mathcal{S} \quad (5)$$

$$Q(s, a) = \mathbb{E} \left[ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s'|s, a) Q(s', a') \right], \quad (s, a) \in (\mathcal{S}, \mathcal{A}) \quad (6)$$

The optimal state value  $V^*(s)$  is the maximum of the expected return under the constraints that the process starts at state  $s$ . Similarly, the optimal action value  $Q^*(s, a)$  is the maximum expected return under the constraints that the process starts at state  $s$  and the first action chosen is  $a$ .

The relationship between  $V^*$  and  $Q^*$  is

$$V^*(s) = \max_{a \in \mathcal{A}} Q^*(s, a), \quad s \in \mathcal{S} \quad (7)$$

$$Q^*(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s'|s, a) V^*(s'), \quad s \in \mathcal{S}, a \in \mathcal{A}. \quad (8)$$

$V^*(s)$  and  $Q^*(s, a)$  satisfy **Bellman optimality equations**:

$$V^*(s) = \max_{a \in \mathcal{A}} \left[ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s'|s, a) V^*(s') \right], \quad s \in \mathcal{S} \quad (9)$$

$$Q^*(s, a) = \sum_{s' \in \mathcal{S}} \mathcal{P}(s'|s, a) \left[ r(s, a) + \gamma \max_{a' \in \mathcal{A}} Q^*(s', a') \right], \quad (s, a) \in (\mathcal{S}, \mathcal{A}) \quad (10)$$

### 3 How to solve RL Problem?

There are three ways to find the optimal policy  $\pi^*$ .

1. The knowledge of  $Q^*$  alone is sufficient for finding  $\pi^*$ .
2. Knowing  $V^*$ ,  $r$  and  $\mathcal{P}$  is also sufficient to find  $\pi^*$ .
3. Find  $\pi^*$  directly without relying on  $Q^*$  and  $V^*$ .

Algorithms follows (1), (2) are called value-based, while algorithm follows (3) is called policy-based.

Basically all algorithms in RL try to answer one of the following questions:

1. How to learn  $V^*$ ?
2. How to learn  $Q^*$ ?
3. How to learn  $\pi^*$ ?

The fundamental mechanism of value-based algorithms is **Policy Iteration**: policy evaluation and policy improvement

- Policy evaluation

Based on the Bellman equation, the state-value under a specific policy  $\pi$  can be estimated.

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \left[ r + \gamma V_\pi(s') \right] \quad (11)$$

Using the Bellman equation as an update rule:

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \left[ r + \gamma V_k(s') \right] \quad (12)$$

In the limit,  $k \rightarrow \infty$ ,  $V_k(s) \rightarrow V_\pi(s)$ ,  $\forall s \in \mathcal{S}$ .

- Policy improvement

An improved policy  $\pi'$  can be obtained by acting greedy with respect to the current policy  $\pi$ .

$$\pi'(s) = \arg \max_a \sum_{s'} p(s'|s, a) \left[ r + \gamma V_\pi(s') \right] \quad (13)$$

#### Value Iteration:

The value iteration update is obtained simply by turning the Bellman optimality equation into an update rule. Also note how the value iteration update is identical to the policy evaluation update except that it requires the maximum to be taken over all actions. In fact, the policy evaluation step of policy iteration can be truncated in several ways without losing the convergence guarantees of policy iteration. One important special case is when policy evaluation is stopped after just one sweep (one update of each state).

$$\begin{aligned} V_{k+1}(s) &= \max_a \mathbb{E} \left[ r_{t+1} + \gamma V_k(s') | S_t = s, A_t = a \right] \\ &= \max_a \sum_{s', r} p(s'|s, a) \left[ r + \gamma V_k(s') \right] \end{aligned} \quad (14)$$

For arbitrary  $V_0$ , the sequence  $\{V_k\}$  can be shown to converge to  $V^*$  under the same conditions that guarantee the existence of  $V^*$ .

**Generalized Policy Iteration** refers to the general idea of letting policy evaluation and policy improvement processes interact independent of the granularity and other details of the two processes. Almost all reinforcement learning methods are well described as GPI. That is, all have identifiable policies and value functions, with the policy always being improved with respect to the value function and the value function always being driven toward to the value function for the policy.

**Policy Gradient Theorem** is the fundamental basis of policy-based algorithm.

The policy  $\pi(a|s, \theta) = Pr(A_t = a | S_t = s, \theta_t = \theta)$  gives the probability that action  $a$  given that the environment is in state  $s$  and the parameter  $\theta$ . Policy-based algorithms consider learning the policy parameter  $\theta$  based on the gradient of some performance measure  $J(\pi_\theta)$  with respect to the policy parameter. These methods seek to maximize performance, so their update approximate gradient ascent in  $J$ :

$$\theta_{t+1} = \theta_t + \alpha \nabla J(\pi_\theta) \quad (15)$$

where  $\nabla J(\pi_\theta)$  is a stochastic estimate whose expectation approximation approximates the gradient of the performance measure with respect to its argument  $\theta_t$ .

In episodic case, the performance measure  $J$  is the value of the start state of the episode.

$$J(\pi_\theta) = V^{\pi_\theta}(s_0), \quad s_0 \sim \rho_0 \quad (16)$$

where  $V^{\pi_\theta}(s_0)$  is the true value function for  $\pi_\theta$ .

With function approximation, it may seem challenging to change the policy parameter in a way that ensure improvement. The problem is that performance depends on both the action selections and the distribution of states in which these selections are made, and that both these are affected by the policy parameter. Given a state, the effect of the policy parameter on the action, and thus on reward, can be computed in a relatively straightforward way from knowledge of the parameterization. But the effect of the policy on the state distribution is a function of the environment and is typically unknown.

How can we estimate the performance gradient with respect to the policy parameter when the gradient depends on the unknown effect of policy changes on the state distribution.

The **Policy Gradient Theorem** answers this question, which provides us an analytic expression for the gradient of performance with respect to the policy parameter that does not involve the derivative of the state distribution. The theorem establishes that

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla_\theta \pi(a|s, \theta) \quad (17)$$

In continuing case, the performance measure is defined as the average reward rate.

$$J(\theta) = r(\pi_\theta) = \sum_a \mu_\pi(s) \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) r \quad (18)$$

Another objective for continuous case is average value

$$J(\theta) = \sum_s d_\pi(s) V_\pi(s) \quad (19)$$

where  $d_\pi(s)$  is the stationary distribution of Markov chain for policy  $\pi$ .

#### 4 Advantage Actor-Critic

Consider a stochastic parameterized policy  $\pi_\theta$ . We aim to maximize the expected return of the policy  $J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$  where  $\tau$  is the trajectory generated by policy  $\pi_\theta$ .

We would like to optimize the policy by gradient ascent

$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\pi_\theta) \quad (20)$$

To actually use this algorithm we need

- The analytical gradient of policy performance.
- Form a sample estimate of the gradient.

**The analytical gradient** The joint probability of the trajectory  $\tau = (s_0, a_0, \dots, S_{T+1})$  is

$$P(\tau|\theta) = \rho_0(s_0) \prod_{t=1}^T P(s_{t+1}|s_t, a_t) \pi_\theta(a_t|s_t) \quad (21)$$

**The Log-Probability of trajectory**

$$\log P(\tau|\theta) = \log \rho_0(s_0) + \sum_{t=0}^T (\log P(s_{t+1}|s_t, a_t) + \log \pi_\theta(a_t|s_t)) \quad (22)$$

**The log-derivative Trick:** The log-derivative trick is based on a simple rule from calculus: the derivative of  $\log x$  with respect to  $x$  to  $1/x$ .

$$\nabla_\theta P(\tau|\theta) = P(\tau|\theta) \frac{\nabla_\theta P(\tau|\theta)}{P(\tau|\theta)} = P(\tau|\theta) \nabla_\theta \log P(\tau|\theta) \quad (23)$$

**Gradient of Log-Prob**

$$\nabla_\theta \log P(\tau|\theta) = \nabla_\theta \log \rho_0(s_0) + \sum_{t=0}^T (\nabla_\theta \log P(s_{t+1}|s_t, a_t) + \nabla_\theta \log \pi_\theta(a_t|s_t)) \quad (24)$$

As  $\rho_0(s_0)$  and transition model  $P(s_{t+1}|s_t, a_t)$  are independent to  $\theta$ ,

$$\nabla_\theta \log P(\tau|\theta) = \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t|s_t) \quad (25)$$

**The derivation of policy gradient**

$$\begin{aligned} \nabla_\theta J(\pi_\theta) &= \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)] \\ &= \nabla_\theta \int_{\tau} P(\tau|\theta) R(\tau) \\ &= \int_{\tau} \nabla_\theta P(\tau|\theta) R(\tau) \\ &= \int_{\tau} P(\tau|\theta) \nabla_\theta \log P(\tau|\theta) R(\tau) \\ &= \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log P(\tau|\theta) R(\tau)] \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t|s_t) R(\tau) \right] \end{aligned} \quad (26)$$

**The sample estimate of the gradient:** If we collect a set of trajectories  $\mathcal{D} = \{\pi_i\}_{i=1,\dots,N}$ , the policy gradient can be estimated with

$$\hat{g} = \frac{1}{|\mathcal{D}|} \sum_{\tau \sim \mathcal{D}} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) R(\tau) \quad (27)$$

**An improved version**

$$\begin{aligned} \nabla_{\theta} J(\pi_{\theta}) &= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) R(\tau) \right] \\ &= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}) \right] \end{aligned} \quad (28)$$

where

$$\hat{R}_t = \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}) \quad (29)$$

## 5 Value-based and Policy-based Models

There are two main approaches to train agents.

**Policy-optimization:** Approaches in the family represent a policy explicitly  $\pi_{\theta}(a|s)$  and optimize the parameter  $\theta$  by gradient ascent on the performance objective  $J(\pi_{\theta})$ . This optimization is usually conducted on policy where each update only uses data collected from the most recent policy.

**Value-based Model:** Approaches in the family learn an approximator of action-value function  $Q_{\theta}(s, a)$ . The objective function is defined by Bellman equation. The optimization is conducted off-policy where each update use data from any training step. The corresponding policy is obtained by acting greedy with respect to the Q-value  $a(s) = \arg \max_a Q_{\theta}(s, a)$ .

Policy-optimization is more principled and more stable while less data efficient, while Value-based Model is more data efficient and less stable.

## 6 Model-free and Model-based Models

By a model of the environment, we mean a function which predicts state transition and reward function. The agent can be classified as model-free or model-based according to whether the agent has access to (or learns ) a model of the environments.

## 7 Algorithms

### 7.1 TD-Learning

**TD-Learning** is a policy evaluation algorithm, which is designed based on **Bellman equation**.

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \left[ r + \gamma V_{\pi}(s') \right] \quad (30)$$

It uses the Bellman equation as an update rule. It updates  $V(s_t)$  toward estimated return  $r_{t+1} + \gamma V(s_{t+1})$ .

$$V(s_t) \leftarrow V(s_t) + \alpha(r_{t+1} + \gamma V(s_{t+1}) - V(s_t)) \quad (31)$$

where  $r_{t+1} + \gamma \hat{V}(s_{t+1})$  is called TD-target,  $\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$  is the TD-error and  $\alpha$  is the learning step.

To update the policy based on  $V(s)$ , the model of MDP is required. i.e, the reward function  $\mathbf{r}$  and transition matrix  $\mathbf{P}$ .

$$\pi'(a|s_t) = \arg \max_{a \in \mathcal{A}} r(s, a) + P(s_{t+1}|s_t, a) V(s_{t+1}) \quad (32)$$

However, typically the model of MDP is unknown.

With function approximation, suppose the state-value function  $\hat{V}(s, \mathbf{w})$  is parameter by  $\mathbf{w}$ . The goal is to find parameter  $\mathbf{w}$  minimizing mean-squared error between approximate value  $\hat{V}(s, \mathbf{w})$  and true value  $V_\pi(s)$ .

$$J(\mathbf{w}) = \mathbb{E}_\pi[(V_\pi(s) - \mathbf{x}(s)^T \mathbf{w})^2] \quad (33)$$

Gradient descent finds a local minimum

$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w}) = \alpha \mathbb{E}_\pi[(V_\pi(s) - \hat{V}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{V}(s, \mathbf{w})] \quad (34)$$

Stochastic gradient descent samples the gradient

$$\Delta \mathbf{w} = \alpha (V_\pi(s) - \hat{V}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{V}(s, \mathbf{w}) \quad (35)$$

The expected update is equal to full gradient update.

In case of linear function approximator, suppose state-value function can be represented as a linear combination of state feature  $\mathbf{x}(s)$  and coefficient vector  $\mathbf{w}$

$$\hat{v}(s) = \mathbf{x}(s)^T \mathbf{w} \quad (36)$$

Following stochastic gradient descent,  $J(\mathbf{w})$  converges to global optimum. The update rule is particularly simple

$$\nabla_{\mathbf{w}} \hat{V}(s) = \mathbf{x}(s) \quad (37)$$

$$\Delta \mathbf{w} = \alpha (V_\pi(s) - \hat{V}(s)) \mathbf{x}(s) \quad (38)$$

where the true value  $V_\pi(s)$  is unknown and be replaced by TD-target

$$V_\pi(s_t) = r_{t+1} + \gamma \hat{V}_\pi(s_{t+1}) \quad (39)$$

The update step is

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha (r_{t+1} + \gamma \hat{V}_\pi(s_{t+1}) - \hat{V}_\pi(s_t)) \mathbf{x}(s) \quad (40)$$

## 7.2 Q-Learning

**Q-Learning** is an off-policy control algorithm. Off-policy means that evaluate a target policy  $\pi(a|s)$  to compute  $V_\pi(s)$  or  $Q_\pi(s, a)$ , while following behaviour policy  $\mu(a|s)$

$$\{s_1, a_1, r_2, \dots, s_t\} \sim \mu \quad (41)$$

Why this important ?

- Learning from observing humans or other agents.
- Reuse experience generated from old policies  $\pi_1, \pi_2, \dots, \pi_{t-1}$ .
- Learn about optimal policy while following exploratory policy.
- Learn about multiple policies while following on policy.

TD-learning can be used as off-policy with important sampling.

$$V(s_t) \leftarrow V(s_t) + \alpha \left( \frac{\pi(a_t|s_t)}{\mu(a_t|s_t)} (r_{t+1} + \gamma V(s_{t+1})) - V(s_t) \right) \quad (42)$$

Off-policy learning (evaluation) of action-value  $Q(s, a)$ .

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (r_{t+1} + \gamma Q(s_{t+1}, a') - Q(s_t, a_t)) \quad (43)$$

where the behaviour policy gives the next action  $a_{t+1} \sim \mu(a|s_t)$ , the target policy gives an alternative action  $a' \sim \pi(a|s)$  where  $a' \neq a_{t+1}$ .

**KG:** This is off-policy action value evaluation.

Combine the above with policy improvement step yields **Q-Learning**. The target policy is greedy w.r.t  $Q(s, a)$

$$\pi(s_t) = \arg \max_a Q(s_t, a) \quad (44)$$

The behaviour policy  $\mu$  is  $\epsilon$ -greedy w.r.t  $Q(s, a)$ .

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right) \quad (45)$$

**KG:** This update rule merges the policy evaluation and policy improvement in one place.



### 7.3 DQN

DQN is basically Q-Learning with function approximator in the form of neural network. In addition, to deal with highly correlated sample, replay buffer and delay-updated target net are introduced.

As the neural network is updated by SGD, but the fundamental requirements for SGD optimization is that the training data is i.i.d. In our case, the data samples are not independent and they will be close to each other as they belong to the same episode. In addition, the distribution of our training data won't be identical to samples provided by the optimal policy that we want to learn.

To deal with this issue, we usually need to use a large buffer of our past experience and sample training data from it, instead of using our latest experience. This is called replay buffer.

Another issue is the correlation between steps. To make the training more stable, there is a trick called target network, when we keep a copy of our network and use it for  $Q(s', a')$  is the bellman equation, This network is synchronized with our main network only periodically.

*Reinforcement learning is known to be unstable or even to diverge when a nonlinear function approximator such as a neural network is used to represent the Q function. This instability has several causes: the correlations present in the sequence of observations, the fact that small updates to Q may significantly change the policy and therefore change the data distribution, and the correlations between the action-values (Q) and the target values  $r + \gamma \max_a Q(s, a)$ . We address these instabilities with a novel variant of Q-learning, which uses two key ideas. First, we used a biologically inspired mechanism termed experience replay that randomizes over the data, thereby removing correlations in the observation sequence and smoothing over changes in the data distribution (see below for details). Second, we used an iterative update that adjusts the action-values (Q) towards target values that are only periodically updated, thereby reducing correlations with the target. [12]*

KG: Update rule of DQN

#### Algorithm 1: deep Q-learning with experience replay.

```

Initialize replay memory  $D$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights  $\theta$ 
Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$ 
For episode = 1,  $M$  do
    Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$ 
    For  $t = 1, T$  do
        With probability  $\epsilon$  select a random action  $a_t$ 
        otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$ 
        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  $D$ 
        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$ 
        Set  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$ 
        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the network parameters  $\theta$ 
        Every  $C$  steps reset  $\hat{Q} = Q$ 
    End For
End For

```

Figure 1: DQN Algorithm

## 7.4 {N-step, Double, Noisy, Dueling, Categorical}-DQN

In this section, we discuss variants of **DQN**.

**N-step DQN**: The bellman equation can be unrolled into many steps in the fashion as

$$Q(s_t, a_t) = r_t + \gamma \max_a Q(s_{t+1}, a) = r_t + \gamma r_{t+1} + \gamma^2 \max_a Q(s_{t+2}, a) \quad (46)$$

In the case that reward only received at the end of an episode, unroll the update rule many steps can improve the propagation speed of values which improves the convergence. **KG**: This is similar to the idea of eligibility trace. However, how many steps to unroll is a matter of hyper-parameter tuning.

If too many steps are unrolled, DQN will fail to converge as in the intermediate steps, the max operations are omitted which means that we assume the corresponding selected actions are optimal. However, we have no guarantee to assume this is true. Indeed, at the beginning of the training, the agent acts randomly. In this case, the calculated  $(Q_s, a_t)$  maybe smaller than the optimal value of the state-action. Therefore, more unrolling steps leads to more inaccurate estimate. In addition, the large replay buffer makes it even worse, as it increases the chance of getting transitions obtained from old bad policies. This will lead to a wrong update of the current Q-value, as it can easily break our training progress. Nevertheless, the convergence speeding up makes it a worth doing, though the number of unrolling steps is a matter of hyper-parameter tuning.

**KG**: Why off-policy works? Off-policy methods do not need freshness of data. For example, a simple DQN is off-policy, which means that we can use very old data sampled from the environment several million steps old and this data will still be useful for learning. That is because we are just updating the value of the action value  $Q(s_t, a_t)$  with immediate reward  $r_t$ , plus discounted current approximation of the best action's value. Even if the action was sampled randomly, it does not matter because for this particular action  $a_t$ , in the state  $s_t$ , our update will be correct. That is why in off-policy methods, we can use a very large experience replay buffer to make our data closer to being i.i.d.

On the other hand, on-policy methods heavily depend on the training data to be sampled according to the current policy we are updating. That happens because on-policy methods are trying to improve the current policy indirectly or directly.

**Double DQN**: It is noticed that **DQN** tends to overestimate Q-value which may be harmful to training performance and sometimes can lead to sub-optimal policies. The reason is that the max operation in the Bellman equation. In **DQN**, the update rule is

$$Q(s_t, a_t) = r_t + \gamma \max_a Q'(s_{t+1}, a) \quad (47)$$

$Q'(s_{t+1}, a)$  is calculated by target net. The author proposed choosing actions for the next state using the trained network instead of the target net, but calculate the Q-value using the target net.

$$Q(s_t, a_t) = r_t + \gamma Q'(s_t, \arg \max_a Q(s_{t+1}, a)) \quad (48)$$

The authors proved that this simple tweak fixed overestimation completely and they called this new architecture **Double DQN**.

**KG**: Why **DQN** overestimates the Q-value, how to proof?

**Noisy networks**: To balance the trade-off of exploration and exploitation, **DQN** employs  $\epsilon$ -greedy strategy. The author proposed to add noise to the weight of fully-connected layer and adjust the parameters of this noise during training using back-propagation. The authors proposed two ways of adding noises.

- Independent Gaussian noise: For every weight in the fully-connected layer, we have a random value that we draw from the normal distribution.
- Factorized Gaussian noise: To minimize the amount of random values to be samples, the author proposed keeping only two random vectors, one with the size of input and another with the size of the output of the layer. Then, a random matrix for the layer is created by calculating the outer product of the vectors.

**Prioritized replay buffer**: This method tries to improve the efficiency of samples in the replay buffer by prioritized those samples according to the training loss. The author of the paper questioned this uniform random sample policy and proved that by assigning priorities to buffer samples, according to

training loss and sampling the buffer proportional to those priorities, we can significantly improve convergence and the policy quality of the DQN. The method can be seen as train on data that surprises you. The tricky point here is to keep the balance of training on an unusual sample and training on the rest of the buffer.

**Dueling DQN:** The paper noticed that Q-value can be divided into state value and advantage  $Q(s, a) = V(s) + A(s, a)$ . The advantage represents how much extra reward some particular action than in average in the state.

The paper learns the state value and advantage value separately. The networks has two heads one outputs state value the other output advantage value. The combination gives the Q-value. However, we need to make sure the average advantage is zero. This is achieved by  $Q(s, a) = V(s) + A(s, a) - 1/N \sum_k A(s, k)$ .

**Categorical DQN:** DQN output the Q-value of actions. This paper challenges that the distribution of Q-values of actions are more informative in particular in complex environment.

**KG:** Why distribution of Q-value of each action could help?

## 7.5 REINFORCE

**REINFORCE** is a Monte-Carlo Policy gradient algorithm. Before show the detail of **REINFORCE**, we first lay more foundations of policy gradient algorithms.

As discussed in section 2. The policy  $\pi(s, \theta)$  is parameterized by  $\theta$ . Policy-based algorithm is an optimization problem: find  $\theta$  maximizing objective  $J(\theta)$ .

In episodic case, the performance measure  $J$  is the value of the start state of the episode.

$$J(\theta) = V_{\pi_{\theta}}(s_0) \quad (49)$$

In continuing case, the performance measure is defined as the average reward rate.

$$J(\theta) = r(\pi_{\theta}) = \sum_a \mu_{\pi}(s) \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) r \quad (50)$$

Another objective for continuous case is average value

$$J(\theta) = \sum_s d_{\pi}(s) V_{\pi}(s) \quad (51)$$

where  $d_{\pi}(s)$  is the stationary distribution of Markov chain for policy  $\pi$ .

The optimization problem can be solved by gradient descent to search local optimum. What we need is the gradient of  $\theta$  w.r.t the objective function.

$$\Delta \theta = \alpha \nabla_{\theta} J(\theta) \quad (52)$$

where  $\nabla_{\theta} J(\theta)$  is the policy gradient.

**KG:** How to compute the policy gradient analytically?

Assume the policy  $\pi_{\theta}$  is differentiable whenever it is non-zero and we know the gradient  $\nabla_{\theta} \pi_{\theta}(s, a)$ . The likelihood ratios is defined as

$$\nabla_{\theta} \pi_{\theta}(s, a) = \pi_{\theta}(s, a) \frac{\nabla_{\theta} \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)} = \pi_{\theta} \nabla_{\theta} \log \pi_{\theta}(s, a) \quad (53)$$

The score function is defined as  $\nabla_{\theta} \log \pi_{\theta}(s, a)$ .

Remember that what we want is policy gradient  $\nabla_{\theta} J(\theta)$ . **Policy Gradient Theorem** links policy gradient and score function.

**Theorem 1. (Policy Gradient Theorem)** For any differentiable policy  $\pi_{\theta}(s, a)$ , for any of the policy objective functions defined above, the policy gradient is

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[ Q_{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a) \right] \quad (54)$$

Now, we turn our attention to **REINFORCE** exploits Theorem 1 directly. However, there is an issue that  $Q_{\pi_{\theta}}(s, a)$  is unknown. In **REINFORCE**,  $Q_{\pi_{\theta}}(s, a)$  is estimated at the end of each episode where

$$\hat{Q}(s_t, a_t) = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{T-1} r_T \quad (55)$$

At the end of each episode, the update rule is applied to each state-action pair along the episode

$$\theta \leftarrow \theta + \alpha \hat{Q}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) \quad (56)$$

## 7.6 Actor-Critic

In **REINFORCE**,  $Q(s, a)$  is estimated in Monte-Carlo way. This results in high variance. However,  $Q(s, a)$  can be approximated by another function  $Q_{\mathbf{w}}(s, a)$ . The rest are the same as **REINFORCE**, this leads to **Actor-Critic** where  $\pi_{\theta}(s, a)$  is called Actor and  $Q_{\mathbf{w}}(s, a)$  is called Critic. In such way,  $Q(s, a)$  can be estimated at any step instead of waiting until the end of episode. Specifically, After each transition step  $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})$ , the following update rules apply

$$\delta = r_{t+1} + \gamma \hat{Q}_{\mathbf{w}}(s_{t+1}, a_{t+1}) - \hat{Q}_{\mathbf{w}}(s_t, a_t) \quad (57)$$

$$\mathbf{w} \leftarrow \mathbf{w} + \beta \delta \nabla_{\mathbf{w}} \hat{Q}_{\mathbf{w}}(s_t, a_t) \quad (58)$$

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) \hat{Q}_{\mathbf{w}}(s_t, a_t, \mathbf{w}) \quad (59)$$

where  $\alpha, \beta$  are learning step size. In linear case,  $\nabla_{\mathbf{w}} \hat{Q}_{\mathbf{w}}(s_t, a_t) = \mathbf{x}(s_t, a_t)$ .

High variance can be reduced using a baseline  $B(s)$ . It requires that  $B(s)$  reduces variance without changing expectation. A good baseline is the state value function  $B(s) = V(s)$ . Thus the policy gradient can be written using advantage function

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) (Q_{\pi_{\theta}}(s, a) - V_{\pi_{\theta}}(s))] = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A_{\pi_{\theta}}(s, a)] \quad (60)$$

where  $A_{\pi_{\theta}}(s, a) = Q_{\pi_{\theta}}(s, a) - V_{\pi_{\theta}}(s)$ .

Note that TD-error is an unbiased estimate of advantage.

$$\delta = r + \gamma V(s_{t+1}) - V(s_t) \quad (61)$$

$$\mathbb{E}[\delta] = \mathbb{E}[r + \gamma V(s_{t+1}) - V(s_t)] = Q(s_t, a_t) - V(s_t) \quad (62)$$

Thus, in practice, **Advantage Actor-Critic** consists two function approximators:  $\pi_{\theta}$  and  $V_{\mathbf{w}}$ .

**KG: Why does baseline reduce variance? any proof?**

A potential issue of **Actor-Critic** is the bias introduced when approximating the policy gradient.

**KG: Why it is biased?**

The policy gradient is approximated as

$$\nabla_{\theta} J(\theta) \approx \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q_{\mathbf{w}}(s, a)] \quad (63)$$

This is approximation, as  $Q_{\mathbf{w}}(s, a) \approx Q(s, a)$ .

**KG: Does REINFORCE suffer the same issue?** No, as Monte-Carlo estimation  $\hat{Q}(s, a)$  is an unbiased estimator.

Luckily, if we choose value function approximation carefully, we can avoid introducing any bias. The following gives some conditions.

**Theorem 2.** (Compatible function approximation Theorem) If the following two conditions are satisfies:

- Value function is compatible to the policy:

$$\nabla_{\mathbf{w}} Q_{\mathbf{w}}(s, a) = \nabla_{\theta} \log \pi_{\theta}(s, a) \quad (64)$$

- Value function parameters  $\mathbf{w}$  minimize the mean-squared error:

$$\mathbb{E}_{\pi_{\theta}} [(Q_{\pi_{\theta}}(s, a) - Q_{\mathbf{w}}(s, a))^2] \quad (65)$$

Then, the policy gradient is exact.

**KG: AC for continuous action space**

## 7.7 A3C

**A3C:** Asynchronous Advantage Actor-Critic, an on-policy policy-gradient algorithm, gathers transitions (samples) from several parallel environments, all of them exploiting the current policy. The motivation of this setting is to improve the stability of policy gradient method.

The reason behind is the correlation between samples, which breaks the i.i.d assumption of SGD. The negative consequence of correlated samples is very high variance in gradients, which means our training batch contains very similar examples, all of them pushing our network in the same direction. However, this direction may be totally the wrong direction in the global sense, all those examples could be from sub-optimal policies.

**KG:** Can we use replay buffer as in DQN? No, as **Actor-critic** is an on-policy algorithm. Samples from other policies are not useful.

**KG:** Why on-policy algorithm can only use samples from the current policy?

In A3C, samples are collected from several parallel environment while following the same current policy. This breaks the correlation within one single episode, as we now train on several episodes obtained from different environments. As the same time, we are still using the same current policy. A disadvantage is sample inefficiency, as we throw away all experiences that we have just got after one single training. **KG:** This issue will be dealt with by **DPG**, which transform policy-gradient method to be off-policy.

In terms of implementation, **KG:** How samples are gathered together from different environments? **KG:** How the actor net and critic net are updated?

We propose a conceptually simple and lightweight framework for deep reinforcement learning that uses asynchronous gradient descent for optimization of deep neural network controller. We present asynchronous variants of four standard reinforcement learning algorithms and show that parallel actor-learners have a stabilizing effect on training allowing all four methods to successfully training neural network controller. The best performing methods, **A3C**, surpasses the current state-of-the-art on the Atari domain.

**KG:** The Asynchronous RL Framework

We now present multi-threaded asynchronous variants of on-step SARSA, one-step Q-learning, n-step Q-learning and advantage actor-critic. The aim in designing these methods was to find RL algorithms that can train deep neural network policies reliably and without large resource requirements. While the underlying RL methods are quite different, with actor-critic being an on-policy method and Q-learning being an off-policy value-based method, we use two main ideas to make all four algorithms practical given our design goal.

First, we use asynchronous actor-critics, where we use multiple CPU threads on a single machine. Keeping the learners on a single machine removes the communication costs of sending gradients and parameters and enables us to use Hogwild! style updates for training. **KG:** What is Hogwild!?

Second, we make the observation that multiple actors learners running in parallel are likely to be exploring different parts of the environment. Moreover, one can explicitly use different exploration policies in each actor-learner to maximize this diversity. By running different exploration policies in different threads, the overall changes being made to the parameters by multiple actor-learners applying online updates in parallel are likely to be less correlated in time than a single agent applying online updates. Hence, we do not use a replay buffer and rely on parallel actors employing different exploration policies to perform the stabilizing role undertaken by experience replay in the DQN training algorithm.

**KG:** On-policy needs data from the current policy, what happen if the actor learners follows different exploration policies?

In addition to stabilizing learning, using multiple parallel actor learners has multiple practical benefits. First, we obtain a reduction in training time that is roughly linear in the number of parallel actor-learners. Second, since we no longer rely on experience replay, we are able to use on-policy reinforcement learning methods such as actor-critic to train neural network in a stable way.

**KG:** Asynchronous advantage actor-critic **A3C**

The algorithm maintains a policy  $\pi(a_t|s_t, \theta)$  and an estimate of the value function  $V(s_t, \theta_v)$ . This algorithm operates in the forward view and uses the same mix of n-step returns to update both the

policy and the value-function. The policy and the value function are updated after every  $t_{max}$  actions or when a terminal state is reached. The update performed by the algorithm can be seen as

$$\nabla_{\theta'} \log \pi(a_t | s_t, \theta') A(s_t, a_t, \theta, \theta'_v) \quad (66)$$

where  $A(s_t, a_t, \theta, \theta'_v)$  is an estimate of the advantage function given by

$$\sum_{i=0}^{k-1} \gamma^i r_{t+1} + \gamma^k V(s_{t+k}, \theta_v) - V(s_t, \theta_v) \quad (67)$$

where  $k$  can vary from state to state and is upper-bounded by  $t_{max}$ .

As with the value-based methods we rely on parallel actor-learners and accumulated updates to improving training stability. Note that while the parameters  $\theta$  of the policy and  $\theta_v$  of the value function are shown as being separate for generality, we always share some of the parameters in practice. We typically use a convolutional neural network that has one softmax output for the policy  $\pi(a_t | s, \theta)$  and one linear output for the value function  $V(s_t, \theta_v)$  with all non-output layers shared.

We also found that adding the entropy of the policy to the objective function improved exploration by discouraging premature convergence to sub-optimal deterministic policies. The gradient of the full objective function including the entropy regularization term with respect to the policy parameters take the form

$$\nabla_{\theta'} \log \pi(s_t | s_t, \theta') (R_t - V(s_t, \theta_v) + \beta \nabla_{\theta'} H(\pi(s_t; \theta'))) \quad (68)$$

where  $H(\cdot)$  is the entropy. The hyper-parameter  $\beta$  controls the strength of the entropy regularization term.

#### KG: Other Tricks

- We use a shared and slowly changing target network in computing the loss.
- We also accumulate gradients over multiple time-steps before they are applied, which is similar to using mini-batches. This reduces the chance of multiple actor learners overwriting each other's updates. Accumulating updates over several steps also provides some ability to trade-off computational efficiency for data efficiency.
- Finally, we found that giving each thread a different exploration policy helps improve the robustness. While there are many possible ways of making the exploration policies differ we experiment with using  $\epsilon$ -greedy exploration with  $\epsilon$  periodically sampled from some distribution by each thread.
- We use the same mix n-step returns to update the policy and the value-function.
- We typically use a convolutional neural network that has one soft-max output for the policy and one linear output for the value function, with all non-linear layers shared.
- We also found that adding entropy of the policy to the objective function improved exploration by discouraging premature convergence to sub-optimal deterministic policies.

## 7.8 A2C

**A2C** is a synchronous, deterministic variant of Asynchronous Advantage Actor Critic (A3C) which OpenAI found gives equal performance. This algorithm is shown in the [\[Link\]](#).

## 7.9 ACER

**ACER** is the off-policy counterpart of **A3C**.

## 7.10 TRPO

TRPO updates policies by taking the largest possible step to improve performance, while satisfying a special constraint on how close the new and old policies are allowed to be. The constraint is expressed in terms of KL-Divergence, a measure of distance between probability distributions.

Note that in normal policy gradient, small differences in parameter space can have very large differences in performance - so a single bad step can collapse the policy performance. This makes

---

**Algorithm S2** Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.

---

// Assume global shared parameter vectors  $\theta$  and  $\theta_v$  and global shared counter  $T = 0$

// Assume thread-specific parameter vectors  $\theta'$  and  $\theta'_v$

Initialize thread step counter  $t \leftarrow 1$

**repeat**

Reset gradients:  $d\theta \leftarrow 0$  and  $d\theta_v \leftarrow 0$ .

Synchronize thread-specific parameters  $\theta' = \theta$  and  $\theta'_v = \theta_v$

$t_{start} = t$

Get state  $s_t$

**repeat**

Perform  $a_t$  according to policy  $\pi(a_t|s_t; \theta')$

Receive reward  $r_t$  and new state  $s_{t+1}$

$t \leftarrow t + 1$

$T \leftarrow T + 1$

**until** terminal  $s_t$  **or**  $t - t_{start} == t_{max}$

$R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t // \text{ Bootstrap from last state} \end{cases}$

**for**  $i \in \{t-1, \dots, t_{start}\}$  **do**

$R \leftarrow r_i + \gamma R$

Accumulate gradients wrt  $\theta'$ :  $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i|s_i; \theta')(R - V(s_i; \theta'_v))$

Accumulate gradients wrt  $\theta'_v$ :  $d\theta_v \leftarrow d\theta_v + \partial (R - V(s_i; \theta'_v))^2 / \partial \theta'_v$

**end for**

Perform asynchronous update of  $\theta$  using  $d\theta$  and of  $\theta_v$  using  $d\theta_v$ .

**until**  $T > T_{max}$

---

Figure 2: A3C Algorithm

it dangerous to use large step sizes with vanilla policy gradients, thus hurting its sample efficiency. TRPO nicely avoids this kind of collapse, and tends to quickly and monotonically improve performance.

The TRPO update is

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}(\theta_k, \theta) \text{ s.t., } D_{KL}(\theta || \theta_k) \leq \delta \quad (69)$$

**KG:** The intuition is that find the maximize step to optimize the performance while follows some constraints.

where  $\mathcal{L}(\theta_k, \theta)$  is surrogate advantage, a measure of how policy  $\pi_{\theta}$  performs relative to the old policy  $\pi_{\theta_k}$ , using data from the old policy

$$\mathcal{L}(\theta_k, \theta) = \mathbb{E}_{(s,a) \sim \pi_{\theta_k}} \left[ \frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a) \right] \quad (70)$$

$$D_{KL}(\theta || \theta_k) = \mathbb{E}_{s \sim \pi_{\theta_k}} \left[ D_{KL}(\pi_{\theta}(\cdot|s) || \pi_{\theta_k}(\cdot|s)) \right] \quad (71)$$

**TROP** makes some approximation

$$\mathcal{L}(\theta_k, \theta) \approx g^T(\theta - \theta_k) \quad (72)$$

$$D_{KL}(\theta || \theta_k) \approx \frac{1}{2}(\theta - \theta_k)^T H(\theta - \theta_k) \quad (73)$$

The optimization problem is

$$\theta_{k+1} = \arg \max_{\theta} g^T(\theta - \theta_k) \text{ s.t. } \frac{1}{2}(\theta - \theta_k)^T H(\theta - \theta_k) \leq \delta \quad (74)$$

The update rule is

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g \quad (75)$$

Due to the approximation errors introduced by the Taylor expansion, this may not satisfy the KL constraints.

$$\theta_{k+1} = \theta_k + \alpha^j \sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g \quad (76)$$

where  $\alpha = (0, 1)$  and  $j$  nonnegative integer. This is backtracking line search.



## 7.11 PPO

**PPO**: Proximal Policy Optimiaztion is an on-policy policy gradient algorithm like **Actor-Critic**, **A3C** and **A2C**. **PPO** is an improvement over **A3C**, which changes the expression used to estimate the policy gradient. Instead of the gradient of logarithm probability of the action taken, the PPO method uses a different objective: the ratio between the new and old policy scales by the advantage.

To understand **PPO**, we start from policy gradient method, and then move to **TRPO**. Next, we discuss the improvement of **PPO** over **TRPO**.

### Policy gradient methods:

Policy gradient methods compute the the gradient of policy and uses stochastic gradient ascent to update the policy. The most commonly used gradient estimator is

$$\hat{g} = \hat{\mathbb{E}}_t[\nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \hat{A}_t] \quad (77)$$

where  $\pi_{\theta}$  is a stochastic policy and  $\hat{A}_t$  is an estimator of the advantage function at timestep  $t$ . Here, the expectation indicates the empirical average over a finite batch of samples, in an algorithm that alternates between sampling and optimization.

The estimator  $\hat{g}$  is obtained by differentiating the objective

$$L^{PG}(\theta) = \hat{\mathbb{E}}_t[\log \pi_{\theta}(a_t|s_t) \hat{A}_t] \quad (78)$$

While it is appealing to perform multiple steps of optimization on this loss using the same trajectory, doing so is not well-justified, and empirically it often leads to destructively large policy update.

In **TRPO** an objective is maximized subject to a constraint on the size of policy update.

$$\max_{\theta} \hat{\mathbb{E}}_t \left[ \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t \right] \quad (79)$$

$$\text{subject to } \hat{\mathbb{E}}_t[KL[\pi_{\theta_{old}}(\cdot|s_t), \pi_{\theta}(\cdot|s_t)]] \leq \delta$$

Here  $\theta_{old}$  is the vector of policy parameters before the update.

The main objective proposed by **PPO** is

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[ \min \left( \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t, \text{clip} \left( \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right) \right] \quad (80)$$

where  $\epsilon$  is a hyperparamter. The motivation of this objective is as follows:

- The first term is the same as **TRPO**.
- The second term modifies the surrogate objective by clipping the probability ratio, which removes the incentive for moving  $r_t = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$  outside of the interval  $[1 - \epsilon, 1 + \epsilon]$ .
- Finally, we take the minimum of the clipped and unclipped objective, so the final objective is a lower bound on the unclipped objective.

### KG: What does the objective function try to achieve?

With this scheme, we only ignore the change in the ratio when it would make the objective improve and we include it makes the objective worse. Note that the ratio is clipped at  $1 - \epsilon$  or  $1 + \epsilon$  depending on whether the advantage is positive or negative.

### KG: How to compute the advantage?

Regarding to computing the advantage, most techniques for computing variance-reduced advantage-function estimation make use a learned stat-value function  $V(s)$ . We use a truncated version of generalised advantage estimation, which reduce to n-step unrolling advantage when  $\lambda = 1$

$$\hat{A}_t = \delta_t + \gamma \lambda \delta_{t+1} + \dots + (\gamma \lambda)^{T-t+1} \delta_{T-1} \quad (81)$$

where  $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ .

In contrast, **A3C** uses  $n$ -step unrolling advantage.

$$A_t = r_t + \gamma r_{t+1} + \dots + \gamma^{T-t+1} r_{T-1} + \gamma^{T-1} V(s_T) - V(s_t) \quad (82)$$



**KG:** The overall objective function if actor and critic share the same neural network architecture

If using a neural network architecture that shares parameters between the policy and value function, we must use a loss function that combines the policy surrogate and a value function error term. This objective can further be augmented by adding an entropy bonus to ensure sufficient exploration.

$$L(\theta) = \hat{\mathbb{E}}_t \left[ L_t^{clip} - c_1 L_t^{VF} + c_2 S[\pi_\theta](s_t) \right] \quad (83)$$

$S$  denotes an entropy bonus, and  $L_t^{VF}$  is a squared-error loss  $(\hat{V}(s_t) - V_t^{target})^2$ .

**KG:** How to updated via SGD?

**PPO** algorithm that uses fixed-length trajectory segments is shown below. Each iteration, each of  $N$  actors collect  $T$  time-steps of data. Then we construct the surrogate loss on these  $NT$  time-steps of data, and optimize it with mini-batch SGD for  $K$  epoches. **KG:** The length of trajectory is  $NT$ , which is then divided in to  $K$  segments, each segment is used to calculate the loss and gradient to update the model (mini-batch SGD).

---

**Algorithm 1** PPO, Actor-Critic Style

---

```

for iteration=1, 2, ... do
  for actor=1, 2, ...,  $N$  do
    Run policy  $\pi_{\theta_{old}}$  in environment for  $T$  timesteps
    Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
  end for
  Optimize surrogate  $L$  wrt  $\theta$ , with  $K$  epochs and minibatch size  $M \leq NT$ 
   $\theta_{old} \leftarrow \theta$ 
end for

```

---

Figure 3: PPO Algorithm

**PPO** is motivated to solve the same problem as **TROP**. Namely, how can we take the biggest possible policy improvement step without stepping so far that the policy collapse. **TROP** is a second-order method, while **PPO** is a first-order method.

The update rule of **PPO** is

$$\theta_{k+1} = \arg \max_{\theta} \mathbb{E}_{(s,a) \sim \pi_{\theta_k}} [L(s, a, \theta_k, \theta)] \quad (84)$$

The objective function is defined as

$$L(s, a, \theta_k, \theta) = \min \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a), \text{clip} \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A^{\pi_{\theta_k}}(s, a) \right) \quad (85)$$

Another way to express the objective function is

$$L(s, a, \theta_k, \theta) = \min \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a), g(\epsilon, A^{\pi_{\theta_k}}(s, a)) \right) \quad (86)$$

where

$$g(\epsilon, A) = \begin{cases} (1 + \epsilon)A & A \geq 0 \\ (1 - \epsilon)A & A < 0 \end{cases} \quad (87)$$

When advantage is positive

$$L(s, a, \theta_k, \theta) = \min \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, (1 + \epsilon) \right) A^{\pi_{\theta_k}}(s, a) \quad (88)$$

Because the advantage is positive, the objective will increase if the action becomes more likely. i.e.,  $\pi_\theta(a|s)$  increases. The min term puts a limit to how much the objective function can increase.

When the advantage is negative

$$L(s, a, \theta_k, \theta) = \max \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, (1 - \epsilon) \right) A^{\pi_{\theta_k}}(s, a) \quad (89)$$

Because the advantage is negative, the objective will increase if the action becomes less likely. i.e.,  $\pi_\theta(a|s)$  decreases. The max term puts a limit to how much the objective increase.

Steps of **PPO**:

- Collect a set of trajectories  $D_k = \{\tau_i\}$  by running the policy  $\pi_{\theta_k}$ .
- Compute rewards to go  $\hat{R}^t$ .
- Compute Value function  $V_{\phi_k}$ .
- Compute advantage  $\hat{A}_t$ .
- Update policy parameter

$$\theta_{k+1} = \arg \max_{\theta} \frac{1}{|D_k|T} \sum_{\tau \sim D_k} \sum_{t=0}^T \min \left( \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}}(a_t|s_t) A^{\pi_{\theta_k}}(s_t, a_t), g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t)) \right) \quad (90)$$

- Update value function

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|D_k|T} \sum_{\tau \sim D_k} \sum_{t=0}^T \left( V_{\phi}(s_t) - \hat{R}_t \right)^2 \quad (91)$$

## 7.12 DPG

**DPG**: Deterministic policy gradients which is a variation of **Actor-Critic**, but being off-policy. It means we can have a huge replay buffer which improve the sample efficiency greatly.

This paper proposed an off-policy deterministic policy gradient algorithm for RL problems with continuous action space. The algorithm is an **Actor-Critic** typed algorithm which has an actor provides the action (output) to take under an state (input) and a critic while estimate the action value via Q-learning. The exploration is ensured by off-policy learning following behaviour policies (more exploration).

The key difference between **DPG** and **Actor-Critic** that the actor of **Actor-Critic** estimate the stochastic policy which returns the distribution of actions while actor of **DPG** estimate the deterministic policy which returns the action to take.

**KG: Why important sampling can be avoided?**

We note that stochastic off-policy actor-critic algorithms typically use important sampling for both actor and critic. However, because the deterministic policy gradient removes the integral over actions, we can avoid importance sampling in the actor, and by using Q-learning, we can avoid importance sampling in the critic.

**KG: How to update actor via SGD?**

Let  $a = u(s, \theta)$  denote the actor which returns the action  $a$  to take under state  $s$  and  $\hat{Q}(s, a)$  be the critic which estimates the action-value. We can substitute the actor into the critic as  $\hat{Q}(s, a) = \hat{Q}(s, u(s, \theta))$ . Note the objective function of policy gradient algorithm is exact the  $Q$ -value, therefore, we can update the parameter of actor  $u_{\theta}(s)$  to maximize  $Q$ -value. The gradient is

$$\nabla_{\theta} \mathbb{E}_{\pi} [Q(s, \mu(s, \theta))] = \nabla_{\theta} Q(s, \mu(s, \theta)) \nabla_{\theta} \mu(s, \theta) \quad (92)$$

The actor update rule is

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} Q(s, \mu(s, \theta)) \nabla_{\theta} \mu(s, \theta) \quad (93)$$

**KG: The difference with A2C**

Note that despite both **A2C** and **DPG** belonging to policy-gradient family, the way that critic is used is different. In **A2C**, we used the critic as a baseline for the reward, so the critic is a optimal piece (without it, we get **REINFORCE**) and is used to improve the stability.

In **DPG**, as our policy is deterministic, we can now calculate the gradient from  $Q$ , obtained from the critic up to the actor weights, so the whole system is differentiable and could be optimized end-to-end with SGD. To update the critic network, we can use the Bellman equation to find the approximation of  $Q(s, a)$  and minimize the MSE objective.

In summary, the critic is update in the same way as in **A2C** and the actor (policy) is updated in a way to maximize the critic's output. The essence is that the method is off-policy, so we can have a replay buffer for training.

KG: How to ensure exploration?

However, note that unlike in A2C the policy is stochastic, the exploration can be encourage by cross-entropy loss, the policy of DPG is deterministic, we need a mechanism to take care exploration. This can be done by adding noise to the action returned by the actor (policy) or following more exploratory behaviour policies

### 7.13 DDPG

DDPG is an extension of DPG, which combines DPG with neural network, replay buffer and delayed updated target network and batch normalization. It is an off-policy, policy gradient deterministic policy algorithm.

KG: The objective function

Q-learning, a commonly used off-policy algorithm, uses the greedy policy  $\mu(s) = \arg \max_a Q(s, a)$ . We consider function approximation parameterized  $\theta^Q$ , which we optimize by minimizing the loss

$$L(\theta^Q) = \mathbb{E}_{s_t \sim \rho^\beta, a_t \sim \beta, r_t \sim \mathbb{E}}[(Q(s_t, a_t | \theta^Q) - y_t)^2] \quad (94)$$

where

$$y_t = r(s_t, a_t) + \gamma Q(s_{t+1}, \mu(s_{t+1}) | \theta^Q) \quad (95)$$

while  $y_t$  is also dependent on  $\theta^Q$ , this is typically ignored.

The use of large, non-linear function approximation for learning value or action-value functions has often been avoided in the past since theoretical performance guarantees are impossible, a practically learning tend so be unstable. Recently DQN adapted the Q-learning algorithm in order to make effective use of large neural networks as function approximators. In order to scale Q-learning they introduce two major changes: the use of a replay buffer, and a separate target network for calculating  $y_t$ . We employ these in the context of DDPG and explain their implementation in the next section.

DPG algorithm maintains a parameterized actor function  $\mu(s | \theta^\mu)$  which specifies the current policy by deterministically mapping states to a specific action. The critic  $Q(s, a)$  is learned using the Bellman equation as in Q-learning. The actor is updated by following the applying the chain rule to the expected return from the start distribution  $J$  with respect to the actor parameters:

$$\begin{aligned} \nabla_{\theta^\mu} J &\approx \mathbb{E}_{s_t \sim \rho^\beta} [\nabla_{\theta^\mu} Q(s, a | \theta^Q) | s = s_t, a = \mu(s_t | \theta^\mu)] \\ &= \mathbb{E}_{s_t \sim \rho^\beta} [\nabla_a Q(s, a | \theta^Q) | s = s_t, a = \mu(s_t) \nabla_{\theta^\mu} \mu(s | \theta^\mu) | s = s_t] \end{aligned} \quad (96)$$

KG: The update rule

Update the actor policy using the sampled policy gradient

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a | \theta^Q) \nabla_{\theta^\mu} \mu(s | \theta^\mu) \quad (97)$$

Update the target network

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'} \quad (98)$$

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'} \quad (99)$$

KG: Replay Buffer

As in DQN, we used an replay buffer to address the issue of high correlated samples. Transitions are sampled from the environment according to the exploration policy and stored in the replay buffer. At each time step the actor and critic are updated by sampling a minibatch uniformly from the buffer.

KG: Why target net?

We also used the target net in DQN. We create a copy of the actor and critic networks respectively, that are used for calculating the target values. The weights of these target networks are then updated by having them slowly track the learned networks:

$$\theta' \leftarrow \tau \theta + (1 - \tau) \theta' \quad (100)$$

where  $\tau < 1$ . This means that the target values are constrained to change slowly, greatly learning the action-value function closer to the case of supervised learning, a problem for which robust solutions exist.

### KG: Why batch normalization?

When learning from low dimensional feature vector observations, the different components of the observation may have different physical units and the range may vary across environment. We address this issue by adapting batch normalization. Batch Normalization normalizes each dimension across the samples in a mini-batch to have unit mean and variance. In addition, it maintains a running average of the mean and the variance to use for normalization during testing. In deep network, it is used to minimize co-variance shift during training by ensuring that each layer receives whitened input. In the low-dimensional case, we used batch normalization on the state input and all layers of actor network and all layers of critic network prior to the action input. With batch normalization, we were able to learn effectively across many different tasks with differing types of units, without needing to manually ensure the units were within a set range.

### KG: Exploration

A major challenge of learning in continuous actions space is exploration. We construct an exploration policy  $\mu'$  by adding noise sampled from a noise process  $\mathcal{N}$  to the actor policy

$$\mu'(s_t) = \mu(s_t) + \mathcal{N} \quad (101)$$

The noise process is Ornstein-Uhlenbeck process.

---

#### Algorithm 1 DDPG algorithm

---

Randomly initialize critic network  $Q(s, a|\theta^Q)$  and actor  $\mu(s|\theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$ .  
Initialize target network  $Q'$  and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q, \theta^{\mu'} \leftarrow \theta^\mu$   
Initialize replay buffer  $R$   
**for** episode = 1, M **do**  
  Initialize a random process  $\mathcal{N}$  for action exploration  
  Receive initial observation state  $s_1$   
  **for** t = 1, T **do**  
    Select action  $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$  according to the current policy and exploration noise  
    Execute action  $a_t$  and observe reward  $r_t$  and observe new state  $s_{t+1}$   
    Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $R$   
    Sample a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$   
    Set  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'}))$   
    Update critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$   
    Update the actor policy using the sampled policy gradient:  

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$
  
    Update the target networks:  

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$
  

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$$
  
  **end for**  
**end for**

---

Figure 4: DDPG Algorithm [7]

When the action space is continuous, the function  $Q^*(s, a)$  is presumed to be differentiable with respect to the action. This allows us to set up an efficient, gradient-based learning rule to policy  $\mu(s)$ .

$$Q^*(s, a) = \mathbb{E}_{s' \sim P} \left[ r(s, a) + \gamma \max_{a'} Q^*(s', a') \right] \quad (102)$$

The Q-value function is learned by MSE

$$L(\phi, D) = \mathbb{E}_{(s, a, r, s') \sim D} \left[ \left( Q_\phi(s, a) - (r + \gamma \max_{a'} Q_\phi(s', a')) \right)^2 \right] \quad (103)$$

Suppose the policy is parameterized as  $\mu_\theta(s)$ .

$$L(\phi, D) = \mathbb{E}_{(s, a, r, s') \sim D} \left[ \left( Q_\phi(s, a) - (r + \gamma \max_{a'} Q_{\phi_{\text{target}}}(s', \mu_{\theta_{\text{target}}}(s'))) \right)^2 \right] \quad (104)$$

where  $\mu_{\theta_{\text{target}}}$  is the target policy.

We want to learn a deterministic policy  $\pi_\theta(s)$  which gives the action that maximizes  $Q_\phi(s, a)$ . To learn such policy, the objective function is simple

$$\max_{\theta} \mathbb{E}_{s \sim D}[Q_\phi(s, \mu_\theta(s))] \quad (105)$$

Steps of **DDPG**:

- Randomly sample a batch of transition  $B = \{(s, a, r, s')\}$  fro buffer.
- Compute value target

$$y(r, s') = r + \gamma Q_{\phi_{target}}(s', \mu_{\theta_{target}}(s')) \quad (106)$$

- Update value function

$$\nabla_{\phi} \frac{1}{|B|} \sum_{(s, a, r, s')} (Q_\phi(s, a) - y(r, s'))^2 \quad (107)$$

- Update policy

$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_\phi(s, \mu_\theta(s)) \quad (108)$$

- update target network

$$\phi_{target} = \rho \phi_{target} + (1 - \rho) \phi \quad (109)$$

$$\theta_{target} = \rho \theta_{target} + (1 - \rho) \theta \quad (110)$$

## 7.14 TD3

**TD3**: Twin Delayed DDPG is an improvement over **DDPG** which focus on the overestimation of action-value in policy-based algorithms. The issue is only exploited in value-based algorithms like **Double DQN** before.

**KG: The theoretical insight**

In the value-based RL methods such as deep Q-learning, function approximation errors are known to lead to overestimated value estimates and sub-optimal policies. We show that this problem persists in an actor-critic setting and propose novel mechanisms to minimize its effects on both actor and critic.

**KG: The motivation of TD3**

While **DDPG** can achieve great performance sometimes, it is frequently brittle with respect to hyper-parameters and other kinds of tuning. A common failure mode for **DDPG** is that the learned Q-function begins to dramatically overestimate Q-values, which then leads to the policy breaking, because it exploits the errors in the Q-function.

**KG: The ket components of TD3**

Twin Delayed DDPG (**TD3**) is an algorithm that addresses this issue by introducing three critical tricks:

- **Clipped Double-Q Learning**. **TD3** learns two Q-functions instead of one (hence “twin”), and uses the smaller of the two Q-values to form the targets in the Bellman error loss functions.
- **“Delayed” Policy Updates**. **TD3** updates the policy (and target networks) less frequently than the Q-function. The paper recommends one policy update for every two Q-function updates.
- **Target Policy Smoothing**. **TD3** adds noise to the target action, to make it harder for the policy to exploit Q-function errors by smoothing out Q along changes in action.

**KG: Update Rule of TD3**

**TD3** aims to solve the overestimate Q-values by proposing three tricks.

**Target Policy Smoothing**: Add clipped noise to the action given by the target policy

$$a'(s') = clip(\mu_{\theta_{target}}(s') + clip(\epsilon, -c, c), a_{low}, a_{high}), \quad \epsilon \sim \mathcal{N}(0, \sigma) \quad (111)$$

---

**Algorithm 1** Twin Delayed DDPG

---

```
1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi_1, \phi_2$ , empty replay buffer  $\mathcal{D}$ 
2: Set target parameters equal to main parameters  $\theta_{\text{targ}} \leftarrow \theta$ ,  $\phi_{\text{targ},1} \leftarrow \phi_1$ ,  $\phi_{\text{targ},2} \leftarrow \phi_2$ 
3: repeat
4:   Observe state  $s$  and select action  $a = \text{clip}(\mu_\theta(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$ , where  $\epsilon \sim \mathcal{N}$ 
5:   Execute  $a$  in the environment
6:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal
7:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$ 
8:   If  $s'$  is terminal, reset environment state.
9:   if it's time to update then
10:    for  $j$  in range(however many updates) do
11:      Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ 
12:      Compute target actions
          
$$a'(s') = \text{clip}(\mu_{\phi_{\text{targ}}}(s') + \text{clip}(\epsilon, -c, c), a_{\text{Low}}, a_{\text{High}}), \quad \epsilon \sim \mathcal{N}(0, \sigma)$$

13:      Compute targets
          
$$y(r, s', d) = r + \gamma(1 - d) \min_{i=1,2} Q_{\phi_{\text{targ},i}}(s', a'(s'))$$

14:      Update Q-functions by one step of gradient descent using
          
$$\nabla_{\phi_i} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi_i}(s, a) - y(r, s', d))^2 \quad \text{for } i = 1, 2$$

15:      if  $j \bmod \text{policy\_delay} = 0$  then
16:        Update policy by one step of gradient ascent using
          
$$\nabla_\theta \frac{1}{|B|} \sum_{s \in B} Q_{\phi_1}(s, \mu_\theta(s))$$

17:        Update target networks with
          
$$\begin{aligned} \phi_{\text{targ},i} &\leftarrow \rho \phi_{\text{targ},i} + (1 - \rho) \phi_i \\ \theta_{\text{targ}} &\leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta \end{aligned} \quad \text{for } i = 1, 2$$

18:      end if
19:    end for
20:  end if
21: until convergence
```

---

Figure 5: TD3 Algorithm

**Clipped double-Q Learning:** the value target is

$$y(r, s') = r + \gamma \min_{i=1,2} Q_{\phi_i, \text{target}}(s', a'(s')) \quad (112)$$

$$L(\phi_1, D) = \mathbb{E}_{(s,a,r,s') \sim D} \left[ (Q_{\phi_1}(s, a) - y(r, s'))^2 \right] \quad (113)$$

$$L(\phi_2, D) = \mathbb{E}_{(s,a,r,s') \sim D} \left[ (Q_{\phi_2}(s, a) - y(r, s'))^2 \right] \quad (114)$$

The policy is learned

$$\max_{\theta} \mathbb{E}_{s \in D} [Q_{\phi_1}(s, \mu_\theta(s))] \quad (115)$$

The policy is updated less frequently than the Q-functions. This helps damp the volatility that normally arises in DDPG because of how a policy update change the target.

### 7.15 Soft Q-Learning

We show that there is a precise equivalence between Q-learning and policy gradient methods in the setting of entropy-regularized reinforcement learning, that **soft Q-learning** is exactly equivalent to a policy gradient method. We also point out a connection between Q-learning methods and natural policy gradient methods.

**Entropy-Regularized RL** Let us define the entropy-augmented return

$$R(s_0) = \sum_{t=0}^{\infty} \gamma^t (r_t - \tau D_{KL}[\pi_t(s_t) || \bar{\pi}(s_t)]) \quad (116)$$

where  $r_t$  is the reward,  $\gamma \in [0, 1]$  is the discount factor,  $\tau$  is a scalar coefficient and  $KL_t = D_{KL}[\pi_t(s_t) || \bar{\pi}(s_t)]$  is the KL divergence between the current policy  $\pi_t$  and a reference policy  $\bar{\pi}$ . **KG: What is the motivation of entropy-regularized RL setting?. KG: Why the reference policy is needed?.**

Maximum entropy reinforcement learning alters the RL objective, though the original objective can be recovered using a temperature parameter. More importantly, the maximum entropy formulation provides a substantial improvement in exploration and robustness. Maximum entropy policies are robust in the face of model and estimation errors, and improve exploration by acquiring diverse behaviors.

The state-value function

$$V_{\pi}(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t (r_t - \tau D_{KL}[\pi_t(s_t) || \bar{\pi}(s_t)]) | s_0 = s \right] \quad (117)$$

The Q-function as

$$Q_{\pi}(s, a) = \mathbb{E} \left[ r_0 + \sum_{t=0}^{\infty} \gamma^t (r_t - \tau D_{KL}[\pi_t(s_t) || \bar{\pi}(s_t)]) | s_0 = s, a_0 = a \right] \quad (118)$$

The relationship between  $Q_{\pi}$  and  $V_{\pi}$

$$V_{\pi}(s) = \mathbb{E}_{a \sim \pi} [Q_{\pi}(s, a)] - \tau KL(s) \quad (119)$$

In the entropy-regularized setting, we also need to define the greedy policy, as the optimal policy is stochastic.

The Boltzmann policy

$$\pi_Q^B(\cdot | s) = \arg \max_{\pi} \{ \mathbb{E}_{a \sim \pi} [Q_{\pi}(s, a)] - \tau KL(s) \} = \frac{\bar{\pi}(a|s) \exp(Q(s, a)/\tau)}{\mathbb{E}_{a' \sim \pi} [\exp(Q(s, a')/\tau)]} \quad (120)$$

### Bellman equation

We can generalize the  $\mathcal{T}_{\pi}$  operators in standard reinforcement learning correspond to computing the expected return with a non-step look-ahead: they take the expectation over one step of dynamics, and then fall back on the value function at the next time-step. We can easily generalize these operators to the entropy-regularized setting. We define

$$\mathcal{T}_{\pi} V(s) = \mathbb{E}_{a \sim \pi, (r, s') \sim P(r, s' | s, a)} [r - \tau KL(s) + \gamma V(s')] \quad (121)$$

$$\mathcal{T}_{\pi} Q(s, a) = \mathbb{E}_{(r, s') \sim P(r, s' | s, a)} [r + \gamma (\mathbb{E}_{a' \sim \pi} [Q(s', a')] - \tau KL(s'))] \quad (122)$$

The 'Bellman error' is defined as

$$\delta_t = (r_t - \tau KL_t) + \gamma V(s_{t+1}) - V(s_t) \quad (123)$$

then we have

$$\mathcal{T}_{\pi}^n V(s) = \mathbb{E} \left[ \sum_{t=0}^{n-1} \gamma^t \delta_t + \gamma^n V(s_n) | s_0 = s \right] \quad (124)$$

### Policy Gradient Algorithm

## 7.16 SAC

**KG:** SAC is an off-policy algorithm to learn stochastic policy, while DDPG/TD3 is off-policy algorithm to learn deterministic policy. In this paper, we propose soft-actor-critic, an off-policy actor-critic deep RL algorithm based on the maximum entropy RL framework. In this framework, the actors aims to maximize expected reward while also maximizing entropy. That is, to. succeed at the task while acting as randomly as possible. Prior deep RL methods based on this framework have been formulated as Q-learning methods. By combining off-policy updates with a stable stochastic actor-critic formulation, our method achieves state-of-the-art performance on a range of continuous control benchmark tasks.

We explore how to design an efficient and stable model-free deep RL algorithm for continuous state and action spaces. To that end, we draw on the maximum reward framework, which augments the standard maximum reward RL objectives with an entropy maximization term.

In this paper, we demonstrate that we can devise an off-policy maximum entropy actor-critic algorithm, which we call **SAC**, which provides for both sample-efficient learning and stability. The algorithm consists of three key ingredients: an actor-critic architecture with separate policy and value function networks, an off-policy formulation that enables reuse of previously collected data for efficiency, and entropy maximization to enable stability and exploration.

### KG: Maximum Entropy Reinforcement Learning

Standard RL maximizes the expected sum of rewards

$$\sum_t \mathbb{E}_{(s_t, a_t) \sim \rho_\pi} [r(s_t, a_t)] \quad (125)$$

We consider a more general maximum entropy objective which favors stochastic policies by augmenting the objective with the expected entropy of the policy over  $\rho_\pi(s_t)$ .

$$J(\pi) = \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim \rho_\pi} [r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot|s_t)))] \quad (126)$$

The objective has a number of conceptual and practical advantages:

- First, the policy is incentivized to explore more widely, while giving up on clearly unpromising avenues.
- Second, the policy can capture multiple modes of near-optimal behavior. In problem settings where multiple actions seems equally attractive, the policy will commit equal probability mass to those actions.
- Lastly, prior work has observed improved exploration with this objective, and in our experiments we observe that it considerably improves learning speed.

### KG: Soft policy iteration

In the policy evaluation step of soft policy iteration, we compute the value of a policy  $\pi$  according to the maximum entropy policy objective. The soft Q-value can be computed iteratively.

$$\mathcal{T}^\pi Q(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} [V(s_{t+1})] \quad (127)$$

where

$$V(s_t) = \mathbb{E}_{a_t \sim \pi} [Q(s_t, a_t) - \log \pi(a_t|s_t)] \quad (128)$$

is the soft state value function.

In the policy improvement step, we update the policy towards the exponential of the new Q-function. We update the policy according to

$$\pi_{new} = \arg \min_{\pi' \in \Pi} D_{KL} \left( \pi'(\cdot|s_t) \parallel \frac{\exp(Q^{\pi_{old}}(s_t, \cdot))}{Z^{\pi_{old}}(s_t)} \right) \quad (129)$$

Since in practice we prefer policies that are tractable, we will additionally restrict the policy to some set of policies  $\Pi$ , which can correspond to a parameterized family of distributions such as Gaussians.

The full soft policy iteration algorithm alternates between the soft policy evaluation and the soft policy improvement step, and it will provably converge to the optimal maximum entropy policy among the policies in  $\Pi$ .



### KG: Soft Actor-Critic SAC

We use function approximators for both the Q-function and policy. We consider a parameterized state value function  $V_\psi(s_t)$ , soft Q-function  $Q_\theta(s_t, a_t)$  and a policy  $\pi_\phi(a_t|s_t)$ .

The soft value function  $V_\psi(s_t)$  is trained to minimize the squared residual error

$$J_V(\psi) = \mathbb{E}_{s_t \sim \mathcal{D}} \left[ \frac{1}{2} (V_\psi(s_t) - \mathbb{E}_{a_t \sim \pi_\phi} [Q_\theta(s_t, a_t) - \log \pi_\phi(a_t|s_t)])^2 \right] \quad (130)$$

The gradient is

$$\nabla_\psi J_V(\psi) = \nabla_\psi V_\psi(s_t) (V_\psi(s_t) - Q_\theta(s_t, a_t) + \log \pi_\phi(a_t|s_t)) \quad (131)$$

The soft Q-function is trained to minimize the soft Bellman residual.

$$J_Q(\theta) = \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} \left[ \frac{1}{2} (Q_\theta(s_t, a_t) - \hat{Q}(s_t, a_t))^2 \right] \quad (132)$$

with

$$\hat{Q}(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} [V_{\bar{\psi}}(s_{t+1})] \quad (133)$$

The gradient is

$$\nabla_\theta J_Q(\theta) = \nabla_\theta Q_\theta(s_t, a_t) (Q_\theta(s_t, a_t) - r(s_t, a_t) - \gamma V_{\bar{\psi}}(s_{t+1})) \quad (134)$$

The policy is learned by minimizing the expected KL-divergence

$$J_\pi(\phi) = \mathbb{E}_{s_t \sim \mathcal{D}} [D_{KL}(\pi_\phi(\cdot|s_t) || \frac{\exp(Q_\theta(s_t, \cdot))}{Z_\theta(s_t)})] \quad (135)$$

We use a reparameterization trick resulting in a lower variance estimator. To that end, we reparameterize the policy using a neural network transformation

$$a_t = f_\phi(\epsilon_t, s_t) \quad (136)$$

where  $\epsilon_t$  is an input noise vector, sampled from some fixed distribution, such as spherical Gaussian. The objective is

$$J_\pi(\phi) = \mathbb{E}_{s_t \sim \mathcal{D}, \epsilon_t \sim \mathcal{N}} [\log \pi_\phi(f_\phi(\epsilon_t, s_t)|s_t) - Q_\theta(s_t, f_\phi(\epsilon_t, s_t))] \quad (137)$$

where  $\pi_\phi$  is defined implicitly in terms of  $f_\phi$ . The gradient is

$$\nabla_\phi J_\pi(\phi) = \nabla_\phi \log \pi_\phi(a_t|s_t) + (\nabla_{a_t} \log \pi_\phi(a_t|s_t) - \nabla_{a_t} Q_\theta(s_t, a_t)) \nabla_\phi f_\phi(\epsilon_t, s_t) \quad (138)$$

Where  $a_t$  is evaluated at  $f_\phi(\epsilon_t, s_t)$ .

This algorithm also makes use of two Q-functions to mitigate positive bias (overestimate) in the policy improvement step that is known to degrade the performance of value based methods.

### Entropy-Regularized Reinforcement Learning

The entropy of a random variable is

$$H(P) = \mathbb{E}_{x \sim P} [-\log P(x)] \quad (139)$$

The state value is defined as

$$V^\pi(s) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R(s_t, a_t, s_{t+1}) + H(\pi(\cdot|s_t))) | s_0 = s \right] \quad (140)$$

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} R(s_t, a_t, s_{t+1}) + \alpha \sum_{t=1}^{\infty} \gamma^t H(\pi(\cdot|s_t)) | s_0 = s, a_0 = a \right] \quad (141)$$

The relationship between  $V^\pi$  and  $Q^\pi$  is defined as

$$V^\pi(s) = \mathbb{E}_{a \sim \pi} [Q^\pi(s, a)] + \alpha H(\pi(\cdot|s)) \quad (142)$$

The Bellman equation is

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E}_{s' \sim P, a' \sim \pi} [R(s, a, s') + \gamma (Q^\pi(s', a') + \alpha H(\pi(\cdot|s')))] \\ &= \mathbb{E}_{s' \sim P} [R(s, a, s') + \gamma V^\pi(s')] \\ &= \mathbb{E}_{s' \sim P, a' \sim \pi} [R(s, a, s') + \gamma (Q^\pi(s', a') - \alpha \log(\pi(a'|s')))] \end{aligned} \quad (143)$$

This can be estimated from sampled

$$Q^\pi(s, a) \approx r + \gamma(Q^\pi(s', \bar{a}') - \alpha \log \pi(\bar{a}'|s')), \bar{a}' \sim \pi(\cdot|s') \quad (144)$$

The next state is from the replay buffer and next actions from the current policy.

**SAC** maintains a policy  $\pi_\theta$  and two Q-functions  $Q_{\phi_1}$  and  $Q_{\phi_2}$ .

The loss function for Q-networks is

$$L(\phi_i, D) = \mathbb{E}_{(s,a,r,s') \sim D} \left[ \left( Q_{\phi_i}(s, a) - y(r, s') \right)^2 \right] \quad (145)$$

where the target is computed using clipped double-Q trick.

$$y(r, s') = r + \gamma \left( \min_{j=1,2} Q_{\phi_{target}(s', \bar{a}') - \alpha \log \pi_\theta(\bar{a}'|s')} \right), \bar{a}' \sim \pi_\theta(\cdot|s') \quad (146)$$

To learn the policy, it should maximize  $V^\pi(s)$ .

$$V^\pi(s) = \mathbb{E}_{a \sim \pi} [Q^\pi(s, a) - \alpha \log \pi(a|s)] \quad (147)$$

The action is sampled by squashed Gaussian policy

$$\bar{a}_\theta(s, \xi) = \tanh(\mu_\theta(s) + \sigma_\theta(s) \otimes \xi), \xi \sim \mathcal{N}(0, I) \quad (148)$$

$$\mathbb{E}_{a \sim \pi_\theta} [Q^{\pi_\theta}(s, a) - \alpha \log \pi_\theta(a|s)] = \mathbb{E}_{\xi \sim \mathcal{N}} [Q^{\pi_\theta}(s, \bar{a}_\theta(s, \xi) - \alpha \log \pi_\theta(\bar{a}_\theta(s, \xi)|s)] \quad (149)$$

The policy is optimized according to

$$\max_{\theta} \mathbb{E}_{s \sim D, \xi \sim \mathcal{N}} \left[ \min_{j=1,2} Q_{\phi_j}(s, \bar{a}_\theta(s, \xi)) - \alpha \log \pi_\theta(\bar{a}_\theta(s, \xi)|s) \right] \quad (150)$$

## 7.17 State-of-The-Art Algorithms

### On-Policy Algorithms:

*Vanilla Policy Gradient is the most basic, entry-level algorithm in the deep RL space because it completely predates the advent of deep RL altogether. The core elements of VPG go all the way back to the late 80s / early 90s. It started a trail of research which ultimately led to stronger algorithms such as TRPO and then PPO soon after.*

*A key feature of this line of work is that all of these algorithms are on-policy: that is, they don't use old data, which makes them weaker on sample efficiency. But this is for a good reason: these algorithms directly optimize the objective you care about—policy performance—and it works out mathematically that you need on-policy data to calculate the updates. So, this family of algorithms trades off sample efficiency in favor of stability—but you can see the progression of techniques (from VPG to TRPO to PPO) working to make up the deficit on sample efficiency.*

### Off-Policy Algorithms:

*DDPG is a similarly foundational algorithm to VPG, although much younger—the theory of deterministic policy gradients, which led to DDPG, wasn't published until 2014. DDPG is closely connected to Q-learning algorithms, and it concurrently learns a Q-function and a policy which are updated to improve each other.*

*Algorithms like DDPG and Q-Learning are off-policy, so they are able to reuse old data very efficiently. They gain this benefit by exploiting Bellman's equations for optimality, which a Q-function can be trained to satisfy using any environment interaction data (as long as there's enough experience from the high-reward areas in the environment).*

*But problematically, there are no guarantees that doing a good job of satisfying Bellman's equations leads to having great policy performance. Empirically one can get great performance—and when it happens, the sample efficiency is wonderful—but the absence of guarantees makes algorithms in this class potentially brittle and unstable. TD3 and SAC are descendants of DDPG which make use of a variety of insights to mitigate these issues.*

---

**Algorithm 1** Soft Actor-Critic

---

1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi_1, \phi_2$ , empty replay buffer  $\mathcal{D}$   
2: Set target parameters equal to main parameters  $\phi_{\text{targ},1} \leftarrow \phi_1, \phi_{\text{targ},2} \leftarrow \phi_2$   
3: **repeat**  
4:   Observe state  $s$  and select action  $a \sim \pi_\theta(\cdot|s)$   
5:   Execute  $a$  in the environment  
6:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal  
7:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$   
8:   If  $s'$  is terminal, reset environment state.  
9:   **if** it's time to update **then**  
10:     **for**  $j$  in range(however many updates) **do**  
11:       Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$   
12:       Compute targets for the Q functions:  
$$y(r, s', d) = r + \gamma(1 - d) \left( \min_{i=1,2} Q_{\phi_{\text{targ},i}}(s', \tilde{a}') - \alpha \log \pi_\theta(\tilde{a}'|s') \right), \quad \tilde{a}' \sim \pi_\theta(\cdot|s')$$
  
13:       Update Q-functions by one step of gradient descent using  
$$\nabla_{\phi_i} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi_i}(s, a) - y(r, s', d))^2 \quad \text{for } i = 1, 2$$
  
14:       Update policy by one step of gradient ascent using  
$$\nabla_\theta \frac{1}{|B|} \sum_{s \in B} \left( \min_{i=1,2} Q_{\phi_i}(s, \tilde{a}_\theta(s)) - \alpha \log \pi_\theta(\tilde{a}_\theta(s)|s) \right),$$
  
      where  $\tilde{a}_\theta(s)$  is a sample from  $\pi_\theta(\cdot|s)$  which is differentiable wrt  $\theta$  via the reparametrization trick.  
15:       Update target networks with  
$$\phi_{\text{targ},i} \leftarrow \rho \phi_{\text{targ},i} + (1 - \rho) \phi_i \quad \text{for } i = 1, 2$$
  
16:     **end for**  
17:   **end if**  
18: **until** convergence

---

Figure 6: SAC Algorithm

## 8 Papers

- [Sutton]Ref: Reinforcement learning: An introduction [17]  
[DQN]Ref: Human-level control through deep reinforcement learning[12]  
[Double DQN]Ref: Deep reinforcement learning with double q-learning [18]  
[Noisy Network]Ref: Noisy networks for exploration[5]  
[Prioritized Replay Buffer]Ref: Prioritized experience replay [13]  
[Dueling DQN]Ref: Dueling network architectures for deep reinforcement learning[19]  
[Categorical DQN]Ref: A distributional perspective on reinforcement learning [4]  
[RainBow]Ref: Rainbow: Combining improvements in deep reinforcement learning [8]  
[DPG]Ref: Deterministic policy gradient algorithms[16]  
[DDPG]Ref: Continuous control with deep reinforcement learning[10]  
[D4PG]Ref: Distributed distributional deterministic policy gradients[3]  
[TD3]Ref: Addressing function approximation error in actor-critic methods[6]  
[A3C]Ref: Asynchronous methods for deep reinforcement learning [11]  
[A2C][Code | Link]  
[PPO]Ref: Proximal policy optimization algorithms[15]

[Soft Q-learning] Ref: Equivalence between policy gradients and soft q-learning [14]

[SAC]Ref: Soft Actor-Critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor [7]

[NGU] Ref: Never Give Up: Learning Directed Exploration Strategies[2]

[R2D2] Ref: Recurrent experience replay in distributed reinforcement learning[9]

[Agent57] Ref: Agent57: Outperforming the atari human benchmark[1]

## References

- [1] Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskyi, Daniel Guo, and Charles Blundell. Agent57: Outperforming the atari human benchmark. *arXiv preprint arXiv:2003.13350*, 2020.
- [2] Adrià Puigdomènech Badia, Pablo Sprechmann, Alex Vitvitskyi, Daniel Guo, Bilal Piot, Steven Kapturowski, Olivier Tieleman, Martín Arjovsky, Alexander Pritzel, Andrew Bolt, et al. Never give up: Learning directed exploration strategies. *arXiv preprint arXiv:2002.06038*, 2020.
- [3] Gabriel Barth-Maron, Matthew W Hoffman, David Budden, Will Dabney, Dan Horgan, Dhruva Tb, Alistair Muldal, Nicolas Heess, and Timothy Lillicrap. Distributed distributional deterministic policy gradients. *arXiv preprint arXiv:1804.08617*, 2018.
- [4] Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. *arXiv preprint arXiv:1707.06887*, 2017.
- [5] Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves, Vlad Mnih, Remi Munos, Demis Hassabis, Olivier Pietquin, et al. Noisy networks for exploration. *arXiv preprint arXiv:1706.10295*, 2017.
- [6] Scott Fujimoto, Herke Van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. *arXiv preprint arXiv:1802.09477*, 2018.
- [7] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*, 2018.
- [8] Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. *arXiv preprint arXiv:1710.02298*, 2017.
- [9] Steven Kapturowski, Georg Ostrovski, John Quan, Remi Munos, and Will Dabney. Recurrent experience replay in distributed reinforcement learning. In *International conference on learning representations*, 2018.
- [10] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [11] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillcrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937, 2016.
- [12] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [13] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.
- [14] John Schulman, Xi Chen, and Pieter Abbeel. Equivalence between policy gradients and soft q-learning. *arXiv preprint arXiv:1704.06440*, 2017.
- [15] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [16] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. 2014.
- [17] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [18] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. *arXiv preprint arXiv:1509.06461*, 2015.
- [19] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In *International conference on machine learning*, pages 1995–2003, 2016.