

Model-free Reinforcement Learning Summary

2020.11.23 v1

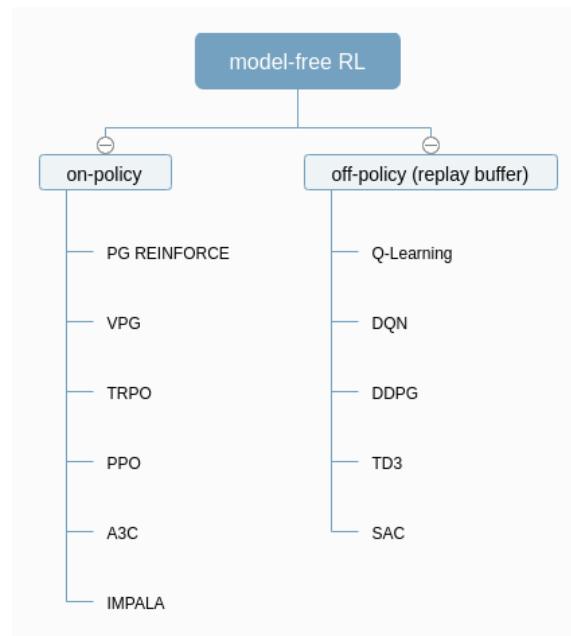
2020.11.27 v2 (more details on Q-learning and TRPO)

李高

Contents

1 Algorithms.....	2
1.1 On-policy RL.....	5
1.1.1 Policy Gradient (PG) 1992.....	5
1.1.2 Vanilla Policy Gradient (VPG) 1999.....	7
1.1.3 Trust Region Policy Optimization (TRPO) 2017.....	10
1.1.4 Proximal Policy Optimization (PPO) 2017.....	17
1.1.5 Synchronous Advantage Actor-Critic (A2C), Asynchronous Advantage Actor-Critic (A3C) 2016.....	19
1.1.6 Importance Weighted Actor-Learner Architecture (IMPALA) 2018.....	20
1.2 Off-policy RL.....	22
1.2.1 Q-learning 1992.....	22
1.2.2 Deep Q-network (DQN) 2015.....	24
1.2.3 Deep Deterministic policy Gradient (DDPG) 2016.....	26
1.2.4 Twin Delayed DDPG (TD3) 2018.....	27
1.2.5 Soft Actor-Critic (SAC) 2018.....	29
2 References.....	33

1 Algorithms



1) 基本概念

1. Observation:

A state s is a complete description of the state of the world.

An observation o is a partial description of a state, which may omit information.

2. Action:

Discrete action spaces, where only a finite number of moves are available to the agent.

Continuous action spaces, actions are real-valued vectors.

3. Policies:

deterministic policy

$$a_t = \mu(s_t)$$

stochastic policy

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

Categorical Policies (for discrete action): you have one final linear layer that gives you logits for each action, followed by a softmax to convert the logits into probabilities.

Diagonal Gaussian Policies (for continuous action): always has a neural network that maps from observations to mean actions, $\mu_\theta(s)$.

4. Trajectories

trajectory τ 由一组 actions 及 states 组成:

$$\tau = (s_0, a_0, s_1, a_1, \dots)$$

其中 s_0 通过 start-state distribution 采样获得: $s_0 = \rho_0(\bullet)$

5. Reward and Return

Reward function R: $r_t = R(s_t, a_t, s_{t+1})$

Return: $R(\tau) = \sum_{t=0}^T r_t$ 或 $R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t$

2) 基本理论

The goal in RL is to select a policy which maximizes expected return when the agent acts according to it.

Let's suppose that both the environment transitions and the policy are stochastic. In this case, the probability of a T -step trajectory is:

$$P(\tau|\pi) = \rho_0(s_0) \prod_{t=0}^{T-1} P(s_{t+1}|s_t, a_t) \pi(a_t|s_t).$$

The expected return (for whichever measure), denoted by $J(\pi)$, is then:

$$J(\pi) = \int_{\tau} P(\tau|\pi) R(\tau) = \mathbb{E}_{\tau \sim \pi} [R(\tau)].$$

The central optimization problem in RL can then be expressed by

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

with π^* being the **optimal policy**.

Value Functions: value means the expected return if you start in that state or state-action pair, and then act according to a particular policy forever after.

There are four main functions of note here.

1. The **On-Policy Value Function**, $V^\pi(s)$, which gives the expected return if you start in state s and always act according to policy π :

$$V^\pi(s) = \underset{\tau \sim \pi}{\mathbb{E}} [R(\tau) | s_0 = s]$$

2. The **On-Policy Action-Value Function**, $Q^\pi(s, a)$, which gives the expected return if you start in state s , take an arbitrary action a (which may not have come from the policy), and then forever after act according to policy π :

$$Q^\pi(s, a) = \underset{\tau \sim \pi}{\mathbb{E}} [R(\tau) | s_0 = s, a_0 = a]$$

3. The **Optimal Value Function**, $V^*(s)$, which gives the expected return if you start in state s and always act according to the *optimal* policy in the environment:

$$V^*(s) = \max_{\pi} \underset{\tau \sim \pi}{\mathbb{E}} [R(\tau) | s_0 = s]$$

4. The **Optimal Action-Value Function**, $Q^*(s, a)$, which gives the expected return if you start in state s , take an arbitrary action a , and then forever after act according to the *optimal* policy in the environment:

$$Q^*(s, a) = \max_{\pi} \underset{\tau \sim \pi}{\mathbb{E}} [R(\tau) | s_0 = s, a_0 = a]$$

two key connections between the value function and the action-value function:

$$V^\pi(s) = \underset{a \sim \pi}{\mathbb{E}} [Q^\pi(s, a)],$$

$$V^*(s) = \max_a Q^*(s, a).$$

The Optimal Q-Function and the Optimal Action:

The optimal policy in s will select whichever action maximizes the expected return from starting in s . As a result, if we have Q^* , we can directly obtain the optimal action, $a^*(s)$, via

$$a^*(s) = \arg \max_a Q^*(s, a).$$

Bellman Equations:

The Bellman equations for the on-policy value functions are

$$\begin{aligned} V^\pi(s) &= \underset{\substack{a \sim \pi \\ s' \sim P}}{\mathbb{E}} [r(s, a) + \gamma V^\pi(s')], \\ Q^\pi(s, a) &= \underset{s' \sim P}{\mathbb{E}} \left[r(s, a) + \gamma \underset{a' \sim \pi}{\mathbb{E}} [Q^\pi(s', a')] \right], \end{aligned}$$

where $s' \sim P$ is shorthand for $s' \sim P(\cdot | s, a)$, indicating that the next state s' is sampled from the environment's transition rules; $a \sim \pi$ is shorthand for $a \sim \pi(\cdot | s)$; and $a' \sim \pi$ is shorthand for $a' \sim \pi(\cdot | s')$.

The Bellman equations for the optimal value functions are

$$\begin{aligned} V^*(s) &= \max_a \underset{s' \sim P}{\mathbb{E}} [r(s, a) + \gamma V^*(s')], \\ Q^*(s, a) &= \underset{s' \sim P}{\mathbb{E}} \left[r(s, a) + \gamma \max_{a'} Q^*(s', a') \right]. \end{aligned}$$

The crucial difference between the Bellman equations for the on-policy value functions and the optimal value functions, is the absence or presence of the max over actions. Its inclusion reflects the fact that whenever the agent gets to choose its action, in order to act optimally, it has to pick whichever action leads to the highest value.

on-policy 比较稳定， off-policy 不那么稳定的原因：

For more information about how and why Q-learning methods can fail, see 1) this classic paper by Tsitsiklis and van Roy, 2) the (much more recent) review by Szepesvari (in section 4.3.2), and 3) chapter 11 of Sutton and Barto, especially section 11.3 (on “the deadly triad” of function approximation, bootstrapping, and off-policy data, together causing instability in value-learning algorithms).

1.1 On-policy RL

1.1.1 Policy Gradient (PG) 1992

1) 理论介绍

直接根据 expected return 对 policy 进行优化： $J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$

由于目标是是要让 $J(\pi_\theta)$ 最大化，所以这种优化 policy 的方法被成为 policy ascent：

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

接下来就是怎么计算梯度 $\nabla_\theta J(\pi_\theta)$ 的问题，需要用的公式如下：

1. Probability of a Trajectory. The probability of a trajectory $\tau = (s_0, a_0, \dots, s_{T+1})$ given that actions come from π_θ is

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

2. 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 is $1/x$. When rearranged and combined with chain rule, we get:

$$\nabla_\theta P(\tau|\theta) = P(\tau|\theta) \nabla_\theta \log P(\tau|\theta).$$

3. Log-Probability of a Trajectory. The log-prob of a trajectory is just

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

4. Gradients of Environment Functions. The environment has no dependence on θ , so gradients of $\rho_0(s_0)$, $P(s_{t+1}|s_t, a_t)$, and $R(\tau)$ are zero.

5. Grad-Log-Prob of a Trajectory. The gradient of the log-prob of a trajectory is thus

$$\begin{aligned} \nabla_\theta \log P(\tau|\theta) &= \nabla_\theta \log \rho_0(s_0) + \sum_{t=0}^T \left(\nabla_\theta \log P(s_{t+1}|s_t, a_t) + \nabla_\theta \log \pi_\theta(a_t|s_t) \right) \\ &= \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t|s_t). \end{aligned}$$

所以推导如下：

Derivation for Basic Policy Gradient

$$\begin{aligned}
\nabla_{\theta} J(\pi_{\theta}) &= \nabla_{\theta} \underset{\tau \sim \pi_{\theta}}{\text{E}} [R(\tau)] \\
&= \nabla_{\theta} \int_{\tau} P(\tau | \theta) R(\tau) && \text{Expand expectation} \\
&= \int_{\tau} \nabla_{\theta} P(\tau | \theta) R(\tau) && \text{Bring gradient under integral} \\
&= \int_{\tau} P(\tau | \theta) \nabla_{\theta} \log P(\tau | \theta) R(\tau) && \text{Log-derivative trick} \\
&= \underset{\tau \sim \pi_{\theta}}{\text{E}} [\nabla_{\theta} \log P(\tau | \theta) R(\tau)] && \text{Return to expectation form} \\
\therefore \nabla_{\theta} J(\pi_{\theta}) &= \underset{\tau \sim \pi_{\theta}}{\text{E}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau) \right] && \text{Expression for grad-log-prob}
\end{aligned}$$

在 policy gradient 算法中可以对 Return 进行修改，但不影响 Expected Return mean 的原因
(如具体算法中的 **reward-to-go Policy Gradient, Baselines Policy Gradients**) :

EGLP Lemma. Suppose that P_{θ} is a parameterized probability distribution over a random variable, x . Then:

$$\underset{x \sim P_{\theta}}{\text{E}} [\nabla_{\theta} \log P_{\theta}(x)] = 0.$$

Proof

Recall that all probability distributions are *normalized*:

$$\int_x P_{\theta}(x) = 1.$$

Take the gradient of both sides of the normalization condition:

$$\nabla_{\theta} \int_x P_{\theta}(x) = \nabla_{\theta} 1 = 0.$$

Use the log derivative trick to get:

$$\begin{aligned}
0 &= \nabla_{\theta} \int_x P_{\theta}(x) \\
&= \int_x \nabla_{\theta} P_{\theta}(x) \\
&= \int_x P_{\theta}(x) \nabla_{\theta} \log P_{\theta}(x) \\
\therefore 0 &= \underset{x \sim P_{\theta}}{\text{E}} [\nabla_{\theta} \log P_{\theta}(x)].
\end{aligned}$$

2) 具体算法

REINFORCE policy based

所以最基本的 Policy Gradient 算法 REINFORCE 就是可以根据下面这个公式获得：

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

其中 $|\mathcal{D}|$ 表示 trajectories 的数量。

具体的 REINFORCE 算法如下：

Initialize the policy parameter θ at random.

Generate one trajectory on policy $\pi_\theta: s_1, a_1, r_2, s_2, a_2, \dots, s_T$.

For $t=1, 2, \dots, T$:

Update policy parameters: $\theta \leftarrow \theta + \alpha R_t \nabla_\theta \log \pi_\theta(a_t | s_t)$

reward-to-go Policy Gradient (从统计上来看这样是不会改变 expectation 的) :

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

Baselines Policy Gradients

根据 EGLP lemma 可以有:

$$\underset{a_t \sim \pi_\theta}{\mathbb{E}} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)] = 0.$$

This allows us to add or subtract any number of terms like this from our expression for the policy gradient, without changing it in expectation:

$$\nabla_\theta J(\pi_\theta) = \underset{\tau \sim \pi_\theta}{\mathbb{E}} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \left(\sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}) - b(s_t) \right) \right].$$

Any function b used in this way is called a **baseline**.

baseline 是不能跟 a_t 有关, 根据 EGLP, 这样处理也是不会改变 expectation 的

3) 实现上的折中及不足

算法的 expectation 是没有问题的, 但每次都要去通过累加的方式计算 Return 还是有点麻烦的, 并且这样单次累加的方式, 算法的 variance 是会比较大的。

1.1.2 Vanilla Policy Gradient (VPG) 1999

1) 算法特点

VPG is an on-policy algorithm.

VPG can be used for environments with either discrete or continuous action spaces.

2) 产生的原因

主要是想通过 baseline 的方式及估计 Q 值去减小 REINFORCE 算法中 variance

理论介绍

What we have seen so far is that the policy gradient has the general form

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \Phi_t \right],$$

where Φ_t could be any of

$$\Phi_t = R(\tau),$$

or

$$\Phi_t = \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}),$$

or

$$\Phi_t = \sum_{t'=t}^T R(s_{t'}, a_{t'}, s_{t'+1}) - b(s_t).$$

All of these choices lead to the same expected value for the policy gradient, despite having different variances. It turns out that there are two more valid choices of weights Φ_t which are important to know.

1. On-Policy Action-Value Function.

$$\text{The choice } \Phi_t = Q^{\pi_{\theta}}(s_t, a_t)$$

2. The Advantage Function.

Recall that the advantage of an action, defined by $A^{\pi}(s_t, a_t) = Q^{\pi}(s_t, a_t) - V^{\pi}(s_t)$, describes how much better or worse it is than other actions on average (relative to the current policy). This choice,

$$\Phi_t = A^{\pi_{\theta}}(s_t, a_t)$$

另外，关于 Q 或 V 的获得方式如下：

In practice, $V^{\pi}(s_t)$ cannot be computed exactly, so it has to be approximated. This is usually done with a neural network, $V_{\phi}(s_t)$, which is updated concurrently with the policy (so that the value network always approximates the value function of the most recent policy).

The simplest method for learning V_{ϕ} , used in most implementations of policy optimization algorithms (including VPG, TRPO, PPO, and A2C), is to minimize a mean-squared-error objective:

$$\phi_k = \arg \min_{\phi} \mathbb{E}_{s_t, \hat{R}_t \sim \pi_k} \left[(V_{\phi}(s_t) - \hat{R}_t)^2 \right],$$

where π_k is the policy at epoch k . This is done with one or more steps of gradient descent, starting from the previous value parameters ϕ_{k-1} .

3) 具体算法

Algorithm 1 Vanilla Policy Gradient Algorithm

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 4: Compute rewards-to-go \hat{R}_t .
- 5: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 6: Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta_k} \hat{A}_t.$$

- 7: Compute policy update, either using standard gradient ascent,

$$\theta_{k+1} = \theta_k + \alpha_k \hat{g}_k,$$

or via another gradient ascent algorithm like Adam.

- 8: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k| T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 9: **end for**
-

4) 实现上的折中及不足

在训练的后期可能会存在 exploration 不足的问题：

VPG trains a stochastic policy in an on-policy way. This means that it explores by sampling actions according to the latest version of its stochastic policy. The amount of randomness in action selection depends on both initial conditions and the training procedure. Over the course of training, the policy typically becomes progressively less random, as the update rule encourages it to exploit rewards that it has already found. This may cause the policy to get trapped in local optima.

并且对于 policy 更新的步长是有比较高的要求的，不能设置太大，不然效果会比较差，设置太小则会很难收敛：

This method uses the first-order derivative and approximates the surface to be flat. If the surface has high curvature, we can make horrible moves.

Large policy change destroys training.

Cannot map changes between policy and parameter space easily.

Improper learning rate causes vanishing or exploding gradient, and Poor sample efficiency.

1.1.3 Trust Region Policy Optimization (TRPO) 2017

1) 算法特点

TRPO is an on-policy algorithm.

TRPO can be used for environments with either discrete or continuous action spaces.

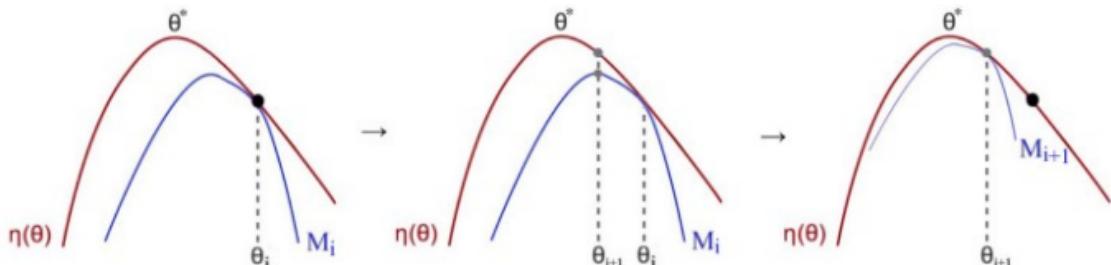
2) 产生的原因

TRPO 主要解决的就是 PG 难收敛的问题，可以保证算法在训练过程中一定能往好的方向优化，但也并不是说保证收敛到全局最优。

相比于 PG，TRPO 是可以多次利用 trajectories，只要能满足相关约束就行。

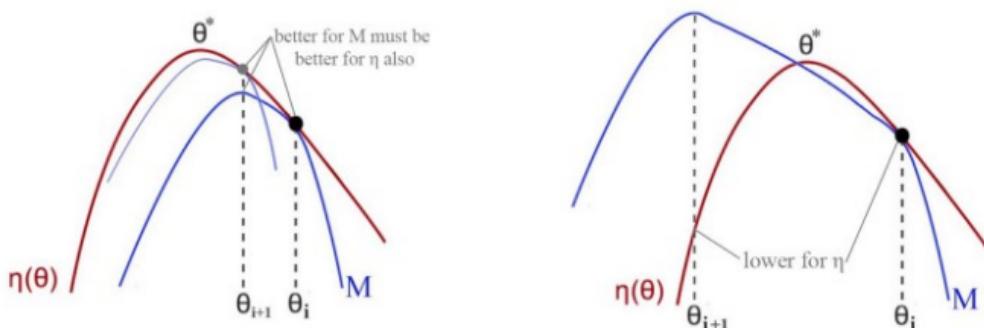
3) 理论介绍

1. Minorize-Maximization MM algorithm



如上图，假如 object function 是 $\eta(\theta)$ ，想要优化 $\eta(\theta)$ 中的 θ ，可以先找到 $\eta(\theta)$ 的一个 lower bound function M_i ，可以去优化 M_i 中的 θ 来保证 $\eta(\theta)$ 中的 θ 也是有被优化到的，迭代这个过程，最终是是可以获得 $\eta(\theta)$ 上的局部或全局最有 θ 的。

但是非 lower bound function 是没有这样的保证的，如下图所示：



2. 然后就是构造 TRPO 的 object function:

$$\underset{\pi'}{\text{maximize}} \eta(\pi') = \underset{\pi'}{\text{maximize}} \eta(\pi') - \eta(\pi)$$

其中：

$$\eta(\pi) = \mathbb{E}_{s_0, a_0, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right], \text{ where}$$

$$s_0 \sim \rho_0(s_0), a_t \sim \pi(a_t|s_t), s_{t+1} \sim P(s_{t+1}|s_t, a_t).$$

上面公式的右边是 TRPO 要优化的 object function，其中 π' 是要更新的 policy (是一个变量)， π 是要旧的 policy (是一个常量，用于获得 trajectories)， $\eta(\pi)$ 是一般 RL 的目标函数。由于 $\eta(\pi)$ 是常量，所以这个等式是成立的。

3. $\eta(\pi)$ 与 $\eta(\pi')$ 的关系有 (推导可见 TRPO appendix A) :

$$\eta(\tilde{\pi}) = \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \tilde{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) \right] \quad (1)$$

可以写成：

$$\begin{aligned} \eta(\tilde{\pi}) &= \eta(\pi) + \sum_{t=0}^{\infty} \sum_s P(s_t = s | \tilde{\pi}) \sum_a \tilde{\pi}(a|s) \gamma^t A_{\pi}(s, a) \\ &= \eta(\pi) + \sum_s \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \tilde{\pi}) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a) \\ &= \eta(\pi) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a). \end{aligned} \quad (2)$$

其中：

$$\rho_{\pi}(s) = P(s_0 = s) + \gamma P(s_1 = s) + \gamma^2 P(s_2 = s) + \dots,$$

如果采用类似 Q-learning 的 deterministic policy 的方式，在每个 state 选的 action 对于的 advantage 值都是非负的，那个根据公式 (2) 更新的 policy 是可以保证持续被优化的。但是如果用 network 及 stochastic policy，advantage 值都是非负的假设是很难保证的，就无法用公式 (2) 对 policy 进行持续优化。所以文中采用 local approximation to η 的方式：

$$L_{\pi}(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a). \quad (3)$$

公式 (2) 与公式 (3) 的区别是在于 state 获取的分布是不同的，公式 (2) 使用的是 $\rho_{\pi}(s)$ ，公式 (3) 使用的是 $\rho_{\pi}(s)$ 。

4. 当 π 与 $\tilde{\pi}$ 相等时，式 (2) 与 (3) 及其一阶导数是相等的：

$$\begin{aligned} L_{\pi_{\theta_0}}(\pi_{\theta_0}) &= \eta(\pi_{\theta_0}), \\ \nabla_{\theta} L_{\pi_{\theta_0}}(\pi_{\theta}) \Big|_{\theta=\theta_0} &= \nabla_{\theta} \eta(\pi_{\theta}) \Big|_{\theta=\theta_0}. \end{aligned} \quad (4)$$

所以可以认为，只要 π 与 $\tilde{\pi}$ 接近，优化公式 (3) 就相当于是在优化公式 (2)。

5. 所以接下来要解释的是 π 跟 $\tilde{\pi}$ 在怎样的数学关系，才是接近，并且还能满足公式 (4)。早期 Kakade & Langford (2002) 的做法如下：

At each iteration update, let π_{old} denote the current policy, and let $\pi' = \arg \max_{\pi'} L_{\pi_{\text{old}}}(\pi')$. The new policy π_{new} was defined to be the following mixture:

$$\pi_{\text{new}}(a|s) = (1 - \alpha)\pi_{\text{old}}(a|s) + \alpha\pi'(a|s). \quad (5)$$

Kakade and Langford derived the following lower bound:

$$\begin{aligned} \eta(\pi_{\text{new}}) &\geq L_{\pi_{\text{old}}}(\pi_{\text{new}}) - \frac{2\epsilon\gamma}{(1-\gamma)^2}\alpha^2 \\ \text{where } \epsilon &= \max_s |\mathbb{E}_{a \sim \pi'(a|s)} [A_{\pi}(s, a)]|. \end{aligned} \quad (6)$$

公式 (5) 的 policy 叫 mixture policies，在实际使用中是不太方便的。TRPO 中想要用的是公式 (6)，为了让算法更容易实现，他们将 α 用 π 与 $\tilde{\pi}$ 的距离 (total variation divergence) 进行替代。并证明这种替代是可行的 (推导可见 TRPO appendix A)，结论如下：

Theorem 1. Let $\alpha = D_{\text{TV}}^{\max}(\pi_{\text{old}}, \pi_{\text{new}})$. Then the following bound holds:

$$\begin{aligned} \eta(\pi_{\text{new}}) &\geq L_{\pi_{\text{old}}}(\pi_{\text{new}}) - \frac{4\epsilon\gamma}{(1-\gamma)^2}\alpha^2 \\ \text{where } \epsilon &= \max_{s,a} |A_{\pi}(s, a)| \end{aligned} \quad (8)$$

由于 KL divergence 与 total variation divergence 有如下关系：

$$D_{\text{TV}}(p \parallel q)^2 \leq D_{\text{KL}}(p \parallel q).$$

所以可以用 KL divergence 替换公式 8 中的 total variation divergence：

$$\eta(\tilde{\pi}) \geq L_{\pi}(\tilde{\pi}) - CD_{\text{KL}}^{\max}(\pi, \tilde{\pi}),$$

where $C = \frac{4\epsilon\gamma}{(1-\gamma)^2}$. (9)

6. 到这一步，就又可以去说明基于公式 (9) 的 policy 优化是可以持续变好的，其实这也就涉及到了 minorization-maximization (MM) algorithm，如下：

It follows from Equation (9) that Algorithm 1 is guaranteed to generate a monotonically improving sequence of policies $\eta(\pi_0) \leq \eta(\pi_1) \leq \eta(\pi_2) \leq \dots$. To see this, let $M_i(\pi) = L_{\pi_i}(\pi) - CD_{\text{KL}}^{\max}(\pi_i, \pi)$. Then

$$\begin{aligned}\eta(\pi_{i+1}) &\geq M_i(\pi_{i+1}) \text{ by Equation (9)} \\ \eta(\pi_i) &= M_i(\pi_i), \text{ therefore,} \\ \eta(\pi_{i+1}) - \eta(\pi_i) &\geq M_i(\pi_{i+1}) - M_i(\pi_i).\end{aligned}\quad (10)$$

Thus, by maximizing M_i at each iteration, we guarantee that the true objective η is non-decreasing. This algorithm

所以对于 policy 的优化，可以写成如下形式：

$$\underset{\theta}{\text{maximize}} [L_{\theta_{\text{old}}}(\theta) - CD_{\text{KL}}^{\max}(\theta_{\text{old}}, \theta)].$$

7、但是为了实现方便，通常是用 constraint 的方法，这也就是 a trust region constraint 了：

$$\begin{aligned}\underset{\theta}{\text{maximize}} \quad &L_{\theta_{\text{old}}}(\theta) \\ \text{subject to} \quad &D_{\text{KL}}^{\max}(\theta_{\text{old}}, \theta) \leq \delta.\end{aligned}\quad (11)$$

但公式 (11) 还是不太方便，因为要求 max，就要穷举所有 states，所以 TRPO 用了 average KL divergence 进行近似，可以避免穷举，所以最终的优化公式为：

$$\begin{aligned}\underset{\theta}{\text{maximize}} \quad &L_{\theta_{\text{old}}}(\theta) \\ \text{subject to} \quad &\overline{D}_{\text{KL}}^{\rho_{\theta_{\text{old}}}}(\theta_{\text{old}}, \theta) \leq \delta.\end{aligned}\quad (12)$$

所以就有了如下形式的 object function:

$$\max_{\pi'} \mathcal{L}_{\pi'}(\pi') - C \sqrt{\mathbb{E}_{s \sim d^{\pi_k}} [D_{KL}(\pi' || \pi)[s]]}$$

or

$$\begin{aligned} & \max_{\pi'} \mathcal{L}_{\pi'}(\pi') \\ & \text{s.t. } \mathbb{E}_{s \sim d^{\pi}} [D_{KL}(\pi' || \pi)[s]] \leq \delta \end{aligned}$$

8. 接下来就是 $\mathcal{L}_{\pi}(\pi')$ 的具体计算过程:

通过泰勒公式将目标函数转化为二次方程，然后就可以直接解了。

To solve it, we can use Taylor's series to expand both terms above up to the second-order. But the second-order of \mathcal{L} is much smaller than the KL-divergence term and will be ignored.

$$\begin{aligned} \mathcal{L}_{\theta_k}(\theta) &\approx \mathcal{L}_{\theta_k}^0 + g^T (\theta - \theta_k) + \dots \\ \bar{D}_{KL}(\theta || \theta_k) &\approx \bar{D}_{KL}^0(\theta_k || \theta_k) + \nabla_{\theta} \bar{D}_{KL}(\theta || \theta_k) |_{\theta_k} (\theta - \theta_k) + \frac{1}{2} (\theta - \theta_k)^T H (\theta - \theta_k) \end{aligned}$$

Taking out all the zero terms, it becomes:

$$\begin{aligned} \mathcal{L}_{\theta_k}(\theta) &\approx g^T (\theta - \theta_k) & g &\doteq \nabla_{\theta} \mathcal{L}_{\theta_k}(\theta) |_{\theta_k} \\ \bar{D}_{KL}(\theta || \theta_k) &\approx \frac{1}{2} (\theta - \theta_k)^T H (\theta - \theta_k) & H &\doteq \nabla_{\theta}^2 \bar{D}_{KL}(\theta || \theta_k) |_{\theta_k} \end{aligned}$$

where g is the policy gradient (the same we learn in Policy Gradient) and H measure the sensitivity (curvature) of the policy relative to the model parameter θ . Our objective turns to:

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

This is a quadratic equation and can be solved analytically:

$$\theta_{k+1} = \theta_k + \underbrace{\sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g}_{\text{natural policy gradient}}$$

9. 上面的优化公式就构成了 Natural Policy Gradient 算法，但是要计算出 Hessian 矩阵及其逆，是没办法用在神经网络上的，所以下面要将 Hessian 矩阵去除：

Instead of finding the inverse of FIM, we want to compute the following combined term directly

$$x_k \approx \hat{H}_k^{-1} \hat{g}_k$$

where x can be solved as:

$$\hat{H}_k x_k \approx \hat{g}$$

We convert this equation into an optimization problem for a quadratic equation:

$$\text{Solving } Ax = b$$

is equivalent to

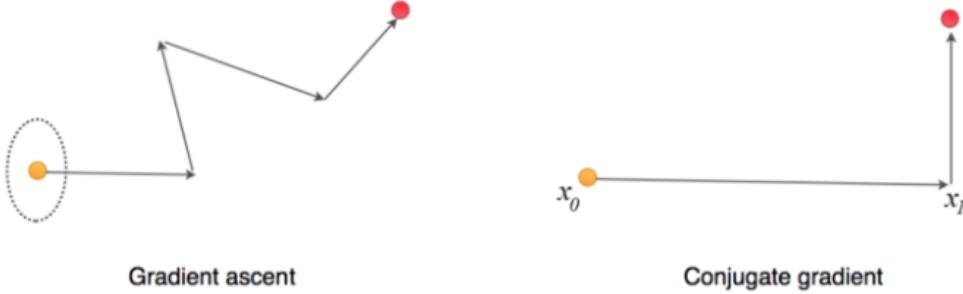
$$\underset{x}{\text{minimize}} \quad f(x) = \frac{1}{2} x^T A x - b^T x$$

$$\text{since } f'(x) = Ax - b = 0$$

In short, we can transform our problem as optimizing the quadratic equation below:

$$\underset{x \in \mathbb{R}^n}{\min} \quad \frac{1}{2} x^T H x - g^T x$$

这样的形式就可以不求逆了，由于上式是一个二次方程，所以用 conjugate gradient 来进行优化会比 gradient ascent 要快（至于二次方程为什么不会受系数相关性的影响我不是很清楚），这样我们就有 Truncated Natural Policy Gradient (TNPG) 算法了。



10. 但是上面的做法中 trust region 是比较小的，所以收敛会比较慢，TRPO 会放宽这个条件：

做法是添加一个步长 α 系数，使优化的速度可以更快一点。但如果这个步长随便给的话，算法就会退回到一般的 PG 算法，所以 TRPO 上的做法是加入一定限制的。这种做法叫 backtracking line search，具体如下：

Algorithm 2 Line Search for TRPO

```

Compute proposed policy step  $\Delta_k = \sqrt{\frac{2\delta}{\hat{g}_k^T \hat{H}_k^{-1} \hat{g}_k}} \hat{H}_k^{-1} \hat{g}_k$ 
for  $j = 0, 1, 2, \dots, L$  do
    Compute proposed update  $\theta = \theta_k + \alpha^j \Delta_k$ 
    if  $\mathcal{L}_{\theta_k}(\theta) \geq 0$  and  $\bar{D}_{KL}(\theta || \theta_k) \leq \delta$  then
        accept the update and set  $\theta_{k+1} = \theta_k + \alpha^j \Delta_k$ 
        break
    end if
end for

```

4) 具体算法

Algorithm 2 Trust Region Policy Optimization

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: Hyperparameters: KL-divergence limit δ , backtracking coefficient α , maximum number of backtracking steps K
- 3: **for** $k = 0, 1, 2, \dots$ **do**
- 4: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 5: Compute rewards-to-go \hat{R}_t .
- 6: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 7: Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)|_{\theta_k} \hat{A}_t.$$

- 8: Use the conjugate gradient algorithm to compute

$$\hat{x}_k \approx \hat{H}_k^{-1} \hat{g}_k,$$

where \hat{H}_k is the Hessian of the sample average KL-divergence.

- 9: Update the policy by backtracking line search with

$$\theta_{k+1} = \theta_k + \alpha^j \sqrt{\frac{2\delta}{\hat{x}_k^T \hat{H}_k \hat{x}_k}} \hat{x}_k,$$

where $j \in \{0, 1, 2, \dots, K\}$ is the smallest value which improves the sample loss and satisfies the sample KL-divergence constraint.

- 10: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k| T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 11: **end for**
-

5) 实现上的折中及不足

同样到训练的后期会存在 exploration 不足的问题，一般会在 policy 的 object function 中加入 entropy。

并且实现起来会相对麻烦一点，并且要计算 Hessian 等原因，在 deep neural network 上是不是很适合的。

6) 参考

<https://jonathan-hui.medium.com/rl-trust-region-policy-optimization-trpo-explained-a6ee04eeee9>

<https://jonathan-hui.medium.com/rl-trust-region-policy-optimization-trpo-part-2-f51e3b2e373a>

1.1.4 Proximal Policy Optimization (PPO) 2017

1) 算法特点

PPO is an on-policy algorithm.

PPO can be used for environments with either discrete or continuous action spaces.

There are two primary variants of PPO: PPO-Penalty and PPO-Clip.

PPO-Penalty approximately solves a KL-constrained update like TRPO, but penalizes the KL-divergence in the objective function instead of making it a hard constraint, and automatically adjusts the penalty coefficient over the course of training so that it's scaled appropriately.

PPO-Clip doesn't have a KL-divergence term in the objective and doesn't have a constraint at all. Instead relies on specialized clipping in the objective function to remove incentives for the new policy to get far from the old policy.

2) 产生的原因

算是 TRPO 实现的一种简化版本，但效果还是比较好的。

Because of this scalability issue, TRPO is not practical for large deep networks.

PPO&ACKTR are introduced to address these problems.

TRPO tries to solve this problem with a complex second-order method, PPO is a family of first-order methods that use a few other tricks to keep new policies close to old.

3) 理论介绍

优化的主要目标还是跟 TRPO 一样的，但处理 KL 的方式有点不一样：

PPO-clip updates policies via

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

typically taking multiple steps of (usually minibatch) SGD to maximize the objective. Here L is given by

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

也有更简化的写法：

PPO-clip updates policies via

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

typically taking multiple steps of (usually minibatch) SGD to maximize the objective. Here L is given by

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

直观上的表示如下：

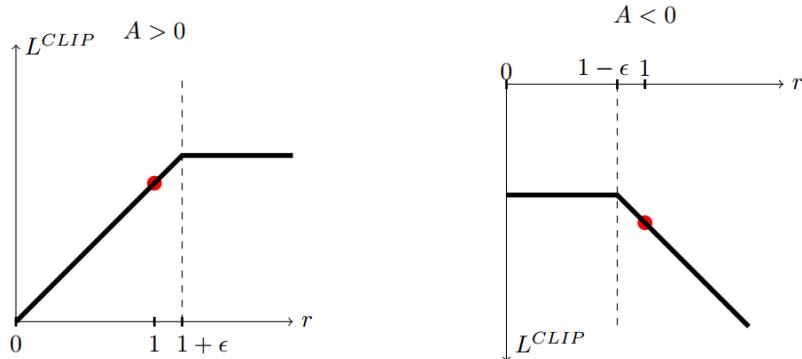


Figure 1: Plots showing one term (i.e., a single timestep) of the surrogate function L^{CLIP} as a function of the probability ratio r , for positive advantages (left) and negative advantages (right). The red circle on each plot shows the starting point for the optimization, i.e., $r = 1$. Note that L^{CLIP} sums many of these terms.

4) 具体算法

Algorithm 3 PPO-Clip

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 4: Compute rewards-to-go \hat{R}_t .
- 5: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 6: Update the policy by maximizing the PPO-Clip objective:

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

typically via stochastic gradient ascent with Adam.

- 7: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 8: **end for**
-

5) 实现上的折中及不足

同样到训练的后期会存在 exploration 不足的问题

1.1.5 Synchronous Advantage Actor-Critic (A2C), Asynchronous Advantage Actor-Critic (A3C) 2016

1) 算法特点

an on-policy algorithm.

can be used for environments with either discrete or continuous action spaces.

2) 产生的原因

主要是为了进行并行运算，弥补 PG 算法 sampling 不足的问题。算法如果使用了神经网络，通过并行也可以缓解数据之间的相关性，防止过拟合。

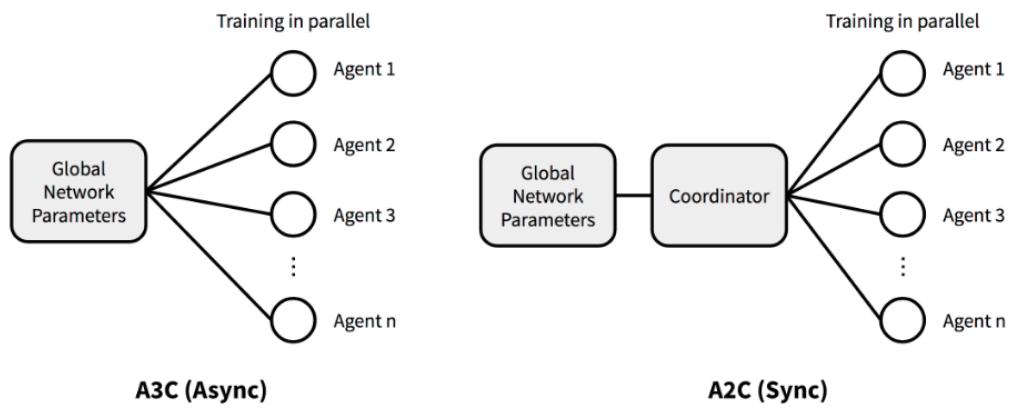
3) 理论介绍

并行部分属于偏重实验的工作

优化部分主要还是基于原有 PG 的理论

4) 具体算法

A3C 采用异步的并行方式，A2C 是同步的并行方式：



Algorithm S3 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  $T > T_{max}$ 

```

1.1.6 Importance Weighted Actor-Learner Architecture (IMPALA) 2018

1) 算法特点

算是一种 off-policy 的算法，但这边还是放在 on-policy 这个部分来介绍，是因为其基本理论还是 PG 那一套。

算法可用于 continue 及 discrete action 的任务上。

2) 产生的原因

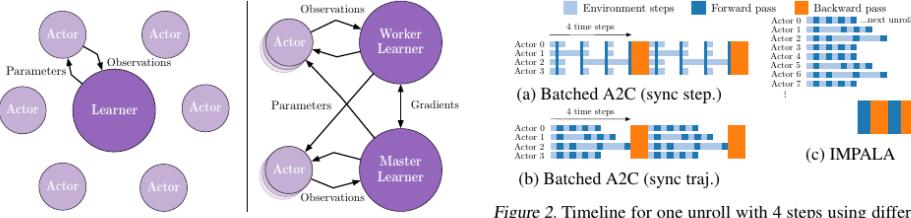


Figure 1. Left: Single Learner. Each actor generates trajectories and sends them via a queue to the learner. Before starting the next trajectory, actor retrieves the latest policy parameters from learner. Right: Multiple Synchronous Learners. Policy parameters are distributed across multiple learners that work synchronously.

Figure 2. Timeline for one unroll with 4 steps using different architectures. Strategies shown in (a) and (b) can lead to low GPU utilisation due to rendering time variance within a batch. In (a), the actors are synchronised after every step. In (b) after every n steps. IMPALA (c) decouples acting from learning.

IMPLALA 是一个大规模强化学习训练的框架，具有较高的性能（high throughput）、较好的扩展性（scalability）和较高的效率（data-efficiency）。

在大规模异步的并行计算框架下，采样（behavior policy 的 V_μ ）和策略（target policy V_π ）更新会有一些错位（不再是完全的 on-policy），在这种情况下，文章通过 V-trace 技术来使用这种 off-policy 样本进行训练。

A3C 是传递的 gradient，IMPALA 传递的是 trajectories。

3) 理论介绍

V- TRACE 就是为了解决 V 值在这种高并行的架构下估计不准的问题，具体如下：

根本的做法还是基于 Importance Sampling

Consider a trajectory $(x_t, a_t, r_t)_{t=s}^{t=s+n}$ generated by the actor following some policy μ . We define the n -steps V-trace target for $V(x_s)$, our value approximation at state x_s , as:

$$v_s \stackrel{\text{def}}{=} V(x_s) + \sum_{t=s}^{s+n-1} \gamma^{t-s} \left(\prod_{i=s}^{t-1} c_i \right) \delta_t V, \quad (1)$$

where $\delta_t V \stackrel{\text{def}}{=} \rho_t (r_t + \gamma V(x_{t+1}) - V(x_t))$ is a temporal difference for V , and $\rho_t \stackrel{\text{def}}{=} \min(\bar{\rho}, \frac{\pi(a_t|x_t)}{\mu(a_t|x_t)})$ and $c_i \stackrel{\text{def}}{=} \min(\bar{c}, \frac{\pi(a_i|x_i)}{\mu(a_i|x_i)})$ are truncated importance sampling (IS) weights (we make use of the notation $\prod_{i=s}^{t-1} c_i = 1$ for $s = t$). In addition we assume that the truncation levels are such that $\bar{\rho} \geq \bar{c}$.

如果 $\rho_i=1 \quad c_i=1$ ，那么上述公式是会回退到 on-policy n-step Bellman 方程。

另外就是系数 $\bar{\rho} \quad \bar{c}$ 的作用的说明：

$\bar{\rho}$ 影响的是 Value 能收敛到的位置，如过 $\bar{\rho}=\infty$ ，那么 value function 将会收敛到 target policy 的 V_π ；而如果 $\bar{\rho}$ 设成接近于 0 的值，那么 value function 将会收敛到 behavior policy 的 V_μ ；如过介于 0 与 ∞ ，那么 value function 会收敛到介于 target policy 与 behavior policy 的 V 上。

c 主要影响的是收敛的速度，当然也是不能太大或太小的。（具体证明可以参考原文）

另外就是可以将上述公式写成一个递归的解法：

Remark 1. *V-trace targets can be computed recursively:*

$$v_s = V(x_s) + \delta_s V + \gamma c_s (v_{s+1} - V(x_{s+1})).$$

4) 具体算法

Consider a parametric representation V_θ of the value function and the current policy π_ω . Trajectories have been generated by actors following some behaviour policy μ . The V-trace targets v_s are defined by (1). At training time s , the value parameters θ are updated by gradient descent on the $l2$ loss to the target v_s , i.e., in the direction of

$$(v_s - V_\theta(x_s)) \nabla_\theta V_\theta(x_s),$$

and the policy parameters ω in the direction of the policy gradient:

$$\rho_s \nabla_\omega \log \pi_\omega(a_s|x_s) (r_s + \gamma v_{s+1} - V_\theta(x_s)).$$

In order to prevent premature convergence we may add an entropy bonus, like in A3C, along the direction

$$-\nabla_\omega \sum_a \pi_\omega(a|x_s) \log \pi_\omega(a|x_s).$$

The overall update is obtained by summing these three gradients rescaled by appropriate coefficients, which are hyperparameters of the algorithm.

1.2 Off-policy RL

1.2.1 Q-learning 1992

1) 算法特点

off-policy，说以会有更高效 sampling

主要是用于处理 discrete action 的问题

Policy 是隐式的，action 的选取直接依据于 Q 值的大小

2) 产生的原因

Q-learning 源于之前提到的 optimal value function 的 Bellman Equation。

3) 理论介绍

optimal value function 的 Bellman Equation:

The Bellman equations for the optimal value functions are

$$V^*(s) = \max_a \mathbb{E}_{s' \sim P} [r(s, a) + \gamma V^*(s')],$$
$$Q^*(s, a) = \mathbb{E}_{s' \sim P} [r(s, a) + \gamma \max_{a'} Q^*(s', a')].$$

从 Dynamic Programming 到上面的 Bellman Equation 还是有一些推导过程的 (Policy iteration → Policy iteration with dynamic programming → value iteration → fitted value iteration → fitted Q-iteration → Q-Learning) , 可以参考 UCB CS285 [Lecture 7: Value Function Methods](#)。

根据 Q 的定义，可以写成： $Q(s, a) = r(s, a) + \gamma E(V(s'))$

近似有： $E(V(s')) \approx \max_a Q(s', a')$

所以有 Q 的 Bellman Equation： $Q(s, a) = \mathbb{E}_{s' \sim P} [r(s, a) + \gamma \max_a Q(s', a')]$

optimal Bellman equation 的收敛证明：

Q-Learning

$$\pi'(s) = \arg \max_a Q^\pi(s, a)$$
$$V^{\pi'}(s) \geq V^\pi(s), \text{ for all state } s$$
$$V^\pi(s) = Q^\pi(s, \pi(s))$$
$$\leq \max_a Q^\pi(s, a) = Q^\pi(s, \pi'(s))$$
$$V^\pi(s) \leq Q^\pi(s, \pi'(s))$$
$$= E[r_{t+1} + V^\pi(s_{t+1}) | s_t = s, a_t = \pi'(s_t)]$$
$$\leq E[r_{t+1} + Q^\pi(s_{t+1}, \pi'(s_{t+1})) | s_t = s, a_t = \pi'(s_t)]$$
$$= E[r_{t+1} + r_{t+2} + V^\pi(s_{t+2}) | \dots]$$
$$\leq E[r_{t+1} + r_{t+2} + Q^\pi(s_{t+2}, \pi'(s_{t+2})) | \dots] \dots \leq V^{\pi'}(s)$$

Created with EverCam
http://www.camdemmy.com

4) 具体算法

下面给出 Q-learning 跟 Sarsa 的对比：

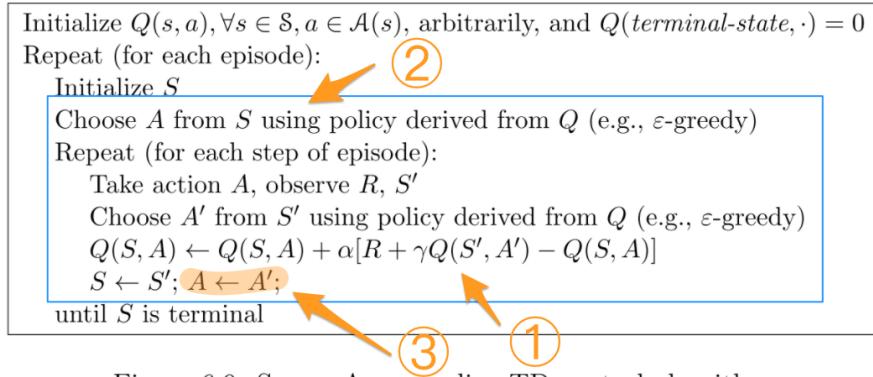


Figure 6.9: Sarsa: An on-policy TD control algorithm.

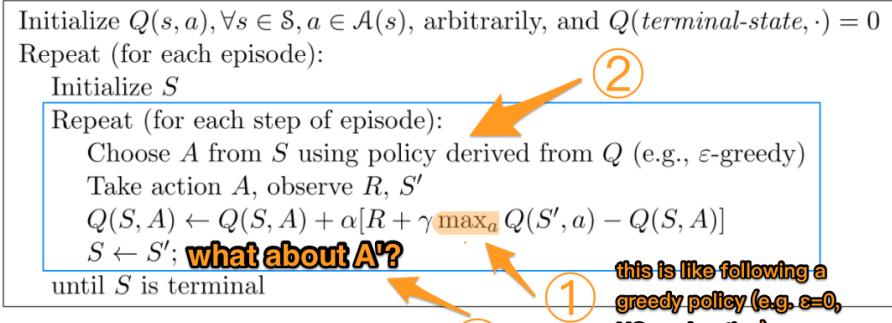


Figure 6.12: Q-learning: An off-policy TD control algorithm.

5) 实现上的折中及不足

处理 continuous action 的任务会比较麻烦

Q 容易被高估

适用于 tabular 的问题，所以 observation space 是不能太大的

1.2.2 Deep Q-network (DQN) 2015

1) 算法特点

off-policy

适用于 discrete action 的问题

2) 产生的原因

为了能直接处理图像输入的 RL 问题

3) 理论介绍

想处理图像作为 observation 的问题，就需要借助 CNN。如果 batch 里面的数据关联性太强，就会造成过拟合，所以就需要采用 replay buffer，并且 DQN 的 replay buffer 通常会是比较大的。

另外就是 target value 需要使用一个 separated Network 来估计，主要是原因如下：

如果不使用一个 separated network 来估计 target value，那么 target value 是跟 value network 相关的，但 Q-learning 的形式又是采用 regression 的方式进行优化的（regression 中的 label 数据需要独立于网络参数），同时 target value 上并没有去计算 gradient，这样算法的收敛会很难保证。

4) 具体算法

Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory D to capacity N

Initialize action-value function Q with random weights θ

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 ε select a random action a_t

 otherwise select $a_t = \text{argmax}_a Q(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, r_t, x_{t+1}$ 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 θ

 Every C steps reset $\hat{Q} = Q$

End For

End For

5) 实现上的折中及不足

因为使用了神经网络，loss 是无法达到 0 的，所以相对于 tabular 的方式，policy 也是无法达到最优的，所以虽然 Bellman Equation 跟 L2 (supervised learning 的形式) 都是可以保证收敛的，但这三者结合起来，却是无法保证收敛的（具体证明可以参考 UCB CS285 [Lecture 7: Value Function Methods](#) 及 [Lecture 8: Deep RL with Q-functions](#)）

1.2.3 Deep Deterministic policy Gradient (DDPG) 2016

1) 算法特点

off policy

只能用于处理 continue action 的问题

2) 产生的原因

算是 Q-learning 用于处理 continue action 问题的一种形式

3) 理论介绍

Critic 部分就 mean-squared Bellman error (MSBE) function:

$$L(\phi, \mathcal{D}) = \underset{(s, a, r, s', d) \sim \mathcal{D}}{\text{E}} \left[\left(Q_\phi(s, a) - \left(r + \gamma(1 - d) \max_{a'} Q_\phi(s', a') \right) \right)^2 \right]$$

Policy 部分:

Policy learning in DDPG is fairly simple. We want to learn a deterministic policy $\mu_\theta(s)$ which gives the action that maximizes $Q_\phi(s, a)$. Because the action space is continuous, and we assume the Q-function is differentiable with respect to action, we can just perform gradient ascent (with respect to policy parameters only) to solve

$$\max_{\theta} \underset{s \sim \mathcal{D}}{\text{E}} [Q_\phi(s, \mu_\theta(s))].$$

Note that the Q-function parameters are treated as constants here.

4) 具体算法

Algorithm 4 Deep Deterministic Policy Gradient

```
1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi$ , empty replay buffer  $\mathcal{D}$ 
2: Set target parameters equal to main parameters  $\theta_{\text{targ}} \leftarrow \theta$ ,  $\phi_{\text{targ}} \leftarrow \phi$ 
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 however many updates do
11:      Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ 
12:      Compute targets

$$y(r, s', d) = r + \gamma(1 - d)Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))$$

13:      Update Q-function by one step of gradient descent using

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

14:      Update policy by one step of gradient ascent using

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

15:      Update target networks with

$$\begin{aligned} \phi_{\text{targ}} &\leftarrow \rho \phi_{\text{targ}} + (1 - \rho) \phi \\ \theta_{\text{targ}} &\leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta \end{aligned}$$

16:    end for
17:  end if
18: until convergence
```

5) 实现上的折中及不足

具体存在的问题就是跟 Q-learning 一样的。

Frequently brittle with respect to hyperparameters and other kinds of tuning.

1.2.4 Twin Delayed DDPG (TD3) 2018

1) 算法特点

off policy

只能用在 continue action 的任务中

2) 产生的原因

主要是为了解决 DDPG 中存在的一些问题，如 Q 值的估计不准，并由此使 policy 很容易学得不好。

3) 理论介绍

主要的解决方法如下：

Trick One: 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.

Trick Two: “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.

Trick Three: 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.

Together, these three tricks result in substantially improved performance over baseline DDPG.

First: **target policy smoothing.** Actions used to form the Q-learning target are based on the target policy, $\mu_{\theta_{\text{targ}}}$, but with clipped noise added on each dimension of the action. After adding the clipped noise, the target action is then clipped to lie in the valid action range (all valid actions, a , satisfy $a_{\text{Low}} \leq a \leq a_{\text{High}}$). The target actions are thus:

$$a'(s') = \text{clip}(\mu_{\theta_{\text{targ}}}(s') + \text{clip}(\epsilon, -c, c), a_{\text{Low}}, a_{\text{High}}), \quad \epsilon \sim \mathcal{N}(0, \sigma)$$

Target policy smoothing essentially serves as a regularizer for the algorithm. It addresses a particular failure mode that can happen in DDPG: if the Q-function approximator develops an incorrect sharp peak for some actions, the policy will quickly exploit that peak and then have brittle or incorrect behavior. This can be averted by smoothing out the Q-function over similar actions, which target policy smoothing is designed to do.

Next: **clipped double-Q learning.** Both Q-functions use a single target, calculated using whichever of the two Q-functions gives a smaller target value:

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

and then both are learned by regressing to this target:

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

$$L(\phi_2, \mathcal{D}) = \underset{(s, a, r, s', d) \sim \mathcal{D}}{\mathbb{E}} \left[\left(Q_{\phi_2}(s, a) - y(r, s', d) \right)^2 \right].$$

Using the smaller Q-value for the target, and regressing towards that, helps fend off overestimation in the Q-function.

Lastly: the policy is learned just by maximizing Q_{ϕ_1} :

$$\max_{\theta} \underset{s \sim \mathcal{D}}{\mathbb{E}} [Q_{\phi_1}(s, \mu_{\theta}(s))],$$

4) 具体算法

Algorithm 5 Twin Delayed DDPG

1: Input: initial policy parameters θ , Q-function parameters ϕ_1, ϕ_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

$$a'(s') = \text{clip}(\mu_{\theta_{\text{targ}}}(s') + \text{clip}(\epsilon, -c, c), a_{\text{Low}}, a_{\text{High}}), \quad \epsilon \sim \mathcal{N}(0, \sigma)$$

 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

$$\phi_{\text{targ},i} \leftarrow \rho \phi_{\text{targ},i} + (1 - \rho) \phi_i \quad \text{for } i = 1, 2$$

$$\theta_{\text{targ}} \leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta$$

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

5) 实现上的折中及不足

具体存在的问题就是跟 Q-learning 一样的。

另外由于是 deterministic policy，算法的开始的时候通常是探索不足的，最好在 action 上加入一定的 noise，但后期为了可以获得更好的 trajectory，这个 noise 可以逐渐去减小。

1.2.5 Soft Actor-Critic (SAC) 2018

1) 算法特点

off policy

可以用于 discrete 跟 continue action 的任务中

2) 产生的原因

采用 stochastic policies，所以 policy 可以更灵活一点，可以用来处理 discrete 跟 continue action 的问题

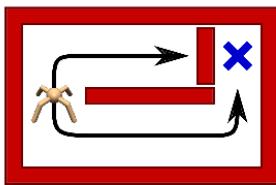
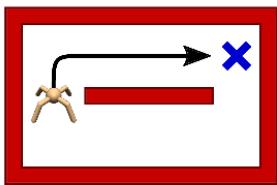
最主要还是 entropy regularization，这样可以更好平衡 exploration 跟 exploitation

具有 on policy 算法的稳定性及 off policy 算法的高效性

3) 理论介绍

Energy-based policy

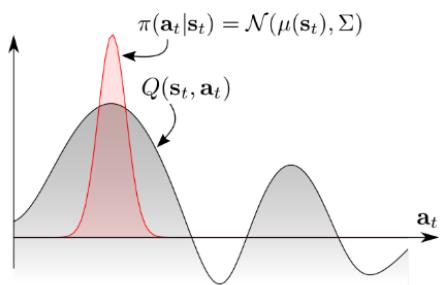
使用 energy-based policy 的原因：



RL 常规的目标函数定义如下：

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{\pi} \left[\sum_{t=0}^T r_t \right]$$

按照这种方式获得的 policy (如 Q-learning)，是一种 unimodal policy，这种 policy 在解决左图的 task 是没什么问题的，但比较难处理右图的 task

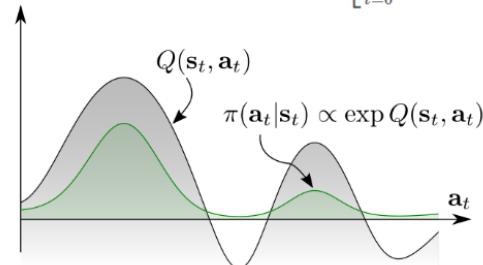


为了解决右边这种 task，需要一种 multi-modal 的 policy：

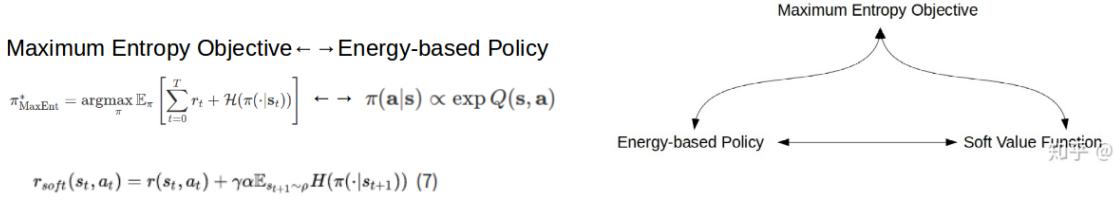
$$\pi(a|s) \propto \exp Q(s, a)$$

在这种 multimodal policy 的定义下是可以推出一种 entropy-based 的 policy 的：

$$\pi_{\text{MaxEnt}}^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{\pi} \left[\sum_{t=0}^T r_t + \mathcal{H}(\pi(\cdot|s_t)) \right]$$



从 Maximum Entropy Object，我们可以推出 Q、V 的形式及更 policy 的关系：



我们将 (7) 带入到公式 (5) :

$$\begin{aligned} Q_{soft}(s_t, a_t) &= r(s_t, a_t) + \gamma \alpha \mathbb{E}_{s_{t+1} \sim \rho} H(\pi(\cdot|s_{t+1})) + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}} [Q_{soft}(s_{t+1}, a_{t+1})] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho, a_{t+1} \sim \pi} [Q_{soft}(s_{t+1}, a_{t+1})] + \gamma \alpha \mathbb{E}_{s_{t+1} \sim \rho} H(\pi(\cdot|s_{t+1})) \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho, a_{t+1} \sim \pi} [Q_{soft}(s_{t+1}, a_{t+1})] + \gamma \mathbb{E}_{s_{t+1} \sim \rho} \mathbb{E}_{a_{t+1} \sim \pi} [-\alpha \log \pi(a_{t+1}|s_{t+1})] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} [\mathbb{E}_{a_{t+1} \sim \pi} [Q_{soft}(s_{t+1}, a_{t+1})] - \alpha \log \pi(a_{t+1}|s_{t+1})] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}} [Q_{soft}(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1}|s_{t+1})] \end{aligned}$$

与此同时，我们知道：

前面我们已经推导得到了公式 (10)，那么根据公式 (10)，我们可以直接推导得到policy的形式：

$$Q(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} [V(s_{t+1})] \quad (9)$$

因此，我们有：

$$\begin{aligned} \pi(s_t, a_t) &= \exp\left(\frac{1}{\alpha}(Q_{soft}(s_t, a_t) - V_{soft}(s_t))\right) \\ &= \frac{\exp\left(\frac{1}{\alpha}Q_{soft}(s_t, a_t)\right)}{\exp\left(\frac{1}{\alpha}V_{soft}(s_t)\right)} \end{aligned} \quad (14)$$

$$V_{soft}(s_t) = \mathbb{E}_{a_t \sim \pi} [Q_{soft}(s_t, a_t) - \alpha \log \pi(a_t|s_t)] \quad (10)$$

接下来就可以设计具体的算法了，主要是因为要计算 entropy，所以使用的 stochastic policy。另外，在 policy 优化的时候，由于要直接对 action 进行 back propagation，所以不能采用 sampling 的 action，SAC 采用了 reparameterization trick 的方法。（这边给出的具体算法跟 Soft Actor Critic 文章上的有点差别，主要是这边没有单独去求 V 值）

Q-function 的训练，具体做的时候，SAC 也跟 TD3 一样，用了两个 Q network 来估计 Q 值：

Before we give the final form of the Q-loss, let's take a moment to discuss how the contribution from entropy regularization comes in. We'll start by taking our recursive Bellman equation for the entropy-regularized Q^π from earlier, and rewriting it a little bit by using the definition of entropy:

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

The RHS is an expectation over next states (which come from the replay buffer) and next actions (which come from the current policy, and **not** the replay buffer). Since it's an expectation, we can approximate it with samples:

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

所以 Q Function 的最终的优化目标如下：

Putting it all together, the loss functions for the Q-networks in SAC are:

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

where the target is given by

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

policy 的训练：

Learning the Policy. The policy should, in each state, act to maximize the expected future return plus expected future entropy. That is, it should maximize $V^\pi(s)$, which we expand out into

$$\begin{aligned} V^\pi(s) &= \underset{a \sim \pi}{\mathbb{E}} [Q^\pi(s, a)] + \alpha H(\pi(\cdot|s)) \\ &= \mathbb{E}_\pi [Q^\pi(s, a) - \alpha \log \pi(a|s)]. \end{aligned}$$

reparameterization trick, SAC 中 action 的产生方式跟 PG 算法有点不一样，一是对 variance σ 添加了另外的 noise，另外使用 tanh 对 action 进行了归一化（为什么要这么做，文章没有给出解释。对 σ 添加 noise，主要是为了模拟 stochastic policy 中的 sampling；tanh 可以将 action 限制在一定范围，可以让训练比较稳定，并且是可导的）：

The way we optimize the policy makes use of the **reparameterization trick**, in which a sample from $\pi_\theta(\cdot|s)$ is drawn by computing a deterministic function of state, policy parameters, and independent noise. To illustrate: following the authors of the SAC paper, we use a squashed Gaussian policy, which means that samples are obtained according to

$$\tilde{a}_\theta(s, \xi) = \tanh(\mu_\theta(s) + \sigma_\theta(s) \odot \xi), \quad \xi \sim \mathcal{N}(0, I).$$

所以 policy 最终的优化目标如下（下面的 policy 的公式是可以根据 KL divergence 推出来的）：

The reparameterization trick allows us to rewrite the expectation over actions (which contains a pain point: the distribution depends on the policy parameters) into an expectation over noise (which removes the pain point: the distribution now has no dependence on parameters):

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

To get the policy loss, the final step is that we need to substitute Q^{π_θ} with one of our function approximators. Unlike in TD3, which uses Q_{ϕ_1} (just the first Q approximator), SAC uses $\min_{j=1,2} Q_{\phi_j}$ (the minimum of the two Q approximators). The policy is thus optimized according to

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

which is almost the same as the DDPG and TD3 policy optimization, except for the min-double-Q trick, the stochasticity, and the entropy term.

4) 具体算法

Algorithm 6 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
```

5) 实现上的折中及不足

算是目前比较好的 model free RL 算法

6) 参考

Haarnoja, Tuomas, et al. "Reinforcement learning with deep energy-based policies." arXiv preprint arXiv:1702.08165 (2017).

Haarnoja, Tuomas, et al. "Reinforcement learning with deep energy-based policies." arXiv preprint arXiv:1702.08165 (2017).

<https://bair.berkeley.edu/blog/2017/10/06/soft-q-learning/>

<https://zhuanlan.zhihu.com/p/70360272>

https://blog.csdn.net/geter_CS/article/details/85068407

2 References

<https://lilianweng.github.io/lil-log/2018/04/08/policy-gradient-algorithms.html>

<https://sites.google.com/view/deep-rl-bootcamp/lectures>

<http://rail.eecs.berkeley.edu/deeprlcourse/>

<https://www.davidsilver.uk/teaching/>