

## SAC Policy Gradient

原文中提到了两种解法：likelihood-ratio gradient estimator和reparameterization，接下来我们简单推导一下这两个解法。

我们优化的函数是：

$$J(\phi) = E_{a \sim \pi_\phi} [\log \pi_\phi(a | s) - Q_\theta(s, a)]$$

### Gradient estimator

$$\begin{aligned} \nabla_\phi J(\phi) &= \nabla_\phi \left[ \sum_a \pi_\phi(a | s) (\log \pi_\phi(a | s) - Q_\theta(s, a)) \right] \\ &= \sum_a [\nabla_\phi \pi_\phi(a | s) (\log \pi_\phi(a | s) - Q_\theta(s, a)) + \pi_\phi(a | s) \nabla_\phi \log \pi_\phi(a | s)] \\ &= E_{a \sim \pi_\phi} [\nabla_\phi \log \pi_\phi(a | s) (\log \pi_\phi(a | s) - Q_\theta(s, a) + 1)] \end{aligned}$$

**注意：**这里是把a当作常数处理，然后直接求出目标函数的导数，由于是一个期望形式，我们可以直接用采样的方式来近似估计这个梯度值。

### Reparameterization

令  $\epsilon \sim N(0, 1)$ ,  $\mu = \pi_{\phi_\mu}(s)$ ,  $\delta = \pi_{\phi_\delta}(s)$ , 可以定义：

$$\begin{aligned} a_\phi^1 &= \delta * \epsilon + \mu \\ a_\phi^2 &= \tanh(a_\phi^1) \\ a_\phi^3 &= scale * a_\phi^2 + bias \\ J(\phi) &= E_{\epsilon \sim N} [\log \pi_\phi(a_\phi^3 | s) - Q_\theta(s, a_\phi^3)] \end{aligned}$$

**注意：**如果不把a当作常数处理，而是把他转化为一个关于 $\phi$ 的函数，这个函数一般是通过一个正态分布输入噪声构造的。这里我们做的事情实际上是让 $J(\phi)$ 可以通过采样来近似，这就需要让分布不依赖于 $\phi$ ，所以我们用了一个 $\epsilon \sim N(0, 1)$ 。