$$\hat{\eta}_\pi(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_\pi(s) \sum_a \tilde{\pi}(a|s) A_\pi(s,a). \tag{1}$$

we can derive Eq.1 and Eq.2 by replacing $\tilde{\pi}$ with parameterized $\pi_\theta$ due to $\sum_s \rho_\pi(s) \sum_a \pi_\theta(a|s) A_\pi(s,a) = 0$.

Let $p(\theta) = \sigma(\tau(r_t(\theta) - 1))$, $r_t(\theta) = \pi_\theta(a_t|s_t)/\pi_{\theta_{old}}(a_t|s_t)$, then we have

$$\nabla_\theta p(\theta) = \sigma(\tau(r_t(\theta) - 1))(1 - \sigma(\tau(r_t(\theta) - 1)))\tau \tag{2}$$
$$= p(\theta)(1 - p(\theta))\tau \tag{3}$$

Then use $\nabla_\theta p(\theta)$ to abbreviate $\nabla_\theta L^{sc}$, we derive

$$\nabla_\theta L^{sc} := \nabla_\theta E_{a \sim \pi_{\theta_{old}}}[p(\theta)\frac{4}{\tau}\hat{A}] \tag{4}$$

$$= \nabla_\theta \int_{\pi_{\theta_{old}}} p(\theta)\frac{4}{\tau}\hat{A}a \tag{5}$$

$$= \int_{\pi_{\theta_{old}}} \nabla_\theta p(\theta)\frac{4}{\tau}\hat{A}a \tag{6}$$

$$= \int_{\pi_{\theta_{old}}} 4p(\theta)(1 - p(\theta)\nabla_\theta r_t(\theta)\hat{A}a \tag{7}$$

$$= E_{\pi_{\theta_{old}}}[4p(\theta)(1 - p(\theta)\nabla_\theta r_t(\theta)\hat{A}] \tag{8}$$