

深度确定策略梯度 (Deep Deterministic Policy Gradient, DDPG)

- 策略网络  $\pi(s; \theta)$ : 价值网络  $q(s, a; \omega)$ ; target network  $q(s, a; \omega^-)$ ; 其中  $\theta^+ = \theta$ ,  $\omega^- = \omega$
- 根据当前状态  $s_t$  执行动作  $a_t = \pi(s_t; \theta)$ , 得到奖励  $r_t$  和下一状态  $s_{t+1}$
- 计算确定策略梯度  $\frac{\partial q(s, a; \omega)}{\partial \theta} = \frac{\partial q(s, a; \omega)}{\partial \theta} = \frac{\partial q(s, a; \omega)}{\partial \theta} = \frac{\partial q(s, a; \omega)}{\partial \theta}$
- 使用策略梯度上升更新策略网络  $\theta_{t+1} = \theta_t + \beta \cdot g$
- 使用价值网络计算  $q_t = q(s_t, a_t; \omega)$
- 使用 target network 计算  $q_{t+1} = q(s_{t+1}, a_{t+1}; \omega^-)$ , 其中  $a_{t+1} = \pi(s_{t+1}; \theta_t)$
- 计算 TD target  $y_t = r_t + \gamma \cdot q_{t+1}$
- 计算 TD error  $\delta_t = q_t - y_t$
- 使用梯度下降更新价值网络  $\omega_{t+1} = \omega_t - \alpha \cdot \delta_t \cdot \frac{\partial q(s, a; \omega)}{\partial \omega} |_{s=s_t, a=a_t}$

每  $N_{\text{update}}$  个 step 后更新 target network 的参数:

- $\theta^- \leftarrow \tau \cdot \theta + (1 - \tau) \cdot \theta^-$
- $\omega^- \leftarrow \tau \cdot \omega + (1 - \tau) \cdot \omega^-$

确定策略梯度  $\frac{\partial q(s, \pi(s; \theta); \omega)}{\partial \theta} = \frac{\partial q(s, a; \omega)}{\partial \theta}$

计算策略梯度  $\frac{\partial q(s, \pi(s; \theta); \omega)}{\partial \theta} = \frac{\partial q(s, a; \omega)}{\partial \theta}$

使用高斯函数作为策略函数

每个动作分量从正态分布中采样得到

$\theta_{t+1} \leftarrow \theta_t + \beta \cdot \frac{\partial q(s, \pi(s; \theta); \omega)}{\partial \theta}$

假设动作空间维度为  $d$

$\mathbf{a} = [a_1, a_2, \dots, a_d]$

$\mathbf{a}_i \sim \mathcal{N}(\mu_i, \sigma_i^2), i = 1, 2, \dots, d$

动作  $\mathbf{a}$  相互独立  $\rightarrow \pi(\mathbf{a}|\mathbf{s}) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_i}} \cdot \exp\left(-\frac{(a_i - \mu_i)^2}{2\sigma_i^2}\right)$

为了求策略梯度 对策略函数取对数  $\rightarrow \ln \pi(\mathbf{a}|\mathbf{s}) = \sum_{i=1}^d \left[ -\ln \sigma_i - \frac{(a_i - \mu_i)^2}{2\sigma_i^2} \right] + \text{const}$

设  $\rho_i = \ln \sigma_i^2, i = 1, 2, \dots, d \rightarrow \ln \pi(\mathbf{a}|\mathbf{s}) = \sum_{i=1}^d \left[ -\frac{\rho_i}{2} - \frac{(a_i - \mu_i)^2}{2 \exp(\rho_i)} \right] + \text{const}$

定义辅助函数  $f(s, \mathbf{a}) = \sum_{i=1}^d \left[ -\frac{\rho_i}{2} - \frac{(a_i - \mu_i)^2}{2 \exp(\rho_i)} \right]$

$\ln \pi(\mathbf{a}|\mathbf{s}) = f(s, \mathbf{a}) + \text{const}$

使用 REINFORCE with baseline 或 A2C 的方法估计完优势函数  $A_\pi(s, \mathbf{a}) \approx A$  后 目标函数  $J$  的值与  $\tau$  有关  $J = \min[\tau \cdot A, \text{clip}(\tau, 1 - \epsilon, 1 + \epsilon) \cdot A]$

等价形式

$J = \min[r \cdot A, h(\epsilon, A)], \text{ where } h(\epsilon, A) = \begin{cases} (1 + \epsilon) \cdot A, & A \geq 0; \\ (1 - \epsilon) \cdot A, & A < 0. \end{cases}$

PPO 最大化目标函数

$J(\theta) = \mathbb{E}_{S, A} \{ \min[r(\theta) \cdot A_\pi(S, A), \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \cdot A_\pi(S, A)] \}$

令  $r(\theta) = \frac{\pi(A|S; \theta)}{\pi(A|S; \theta_{\text{old}})}$ ,  $\theta$  离  $\theta_{\text{old}}$  越远,  $r(\theta)$  越接近 1,  $r(\theta_{\text{old}}) = 1$

将动作价值函数  $Q_\pi(S, A)$  替换为优势函数  $A_\pi(S, A)$

置信域策略优化 (Trust Region Policy Optimization, TRPO)

- 使用策略  $A \sim \pi(\cdot|s; \theta_{\text{old}})$  完成一局游戏  $s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T, a_T, r_T$
- 计算每一步  $i = 1, 2, \dots, T$  的折扣回报  $u_i = \sum_{k=i}^T \gamma^k \cdot r_k$
- 使用均值估计期望近似目标函数  $J(\theta) \approx L(\theta) = \frac{1}{T} \sum_{i=1}^T \min[r(\theta) \cdot \hat{A}_i, \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \cdot \hat{A}_i]$
- 求解带约束条件的优化问题  $\theta_{\text{new}} \leftarrow \text{argmax}_{\theta} L(\theta; \theta_{\text{old}}), \text{ s.t. } \theta \in \mathcal{N}(\theta_{\text{old}})$ , 其中  $\mathcal{N}(\theta_{\text{old}})$  可以选为:
  - $\|\theta - \theta_{\text{old}}\|_2 < \Delta$
  - $\sum_{i=1}^T D_{KL}[\pi(\cdot|s_i; \theta) || \pi(\cdot|s_i; \theta_{\text{old}})] < \Delta$
- 每个 Episode 重复迭代上述 4 步

TRPO 限制了策略网络参数  $\theta$  的更新范围, 使得策略网络不会因为更新幅度过大导致策略网络发散

近端策略优化 (Proximal Policy Optimization, PPO)

目标函数中取最小值的第一项是 TRPO 中的目标函数

目标函数中取最小值的第二项是 TRPO 目标函数的 clipped version ( $\epsilon$  通常取 0.2)

$\text{clip}(x, a, b) = \begin{cases} x, & \text{if } x < a; \\ a, & \text{if } a \leq x \leq b; \\ b, & \text{if } x > b. \end{cases}$

训练流程:

- 使用神经网络  $v(s; \omega)$  近似状态价值函数  $V_\pi(s)$
- 使用策略  $A \sim \pi(\cdot|s; \theta_{\text{old}})$  完成一局游戏  $s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T, a_T, r_T$
- 使用 REINFORCE with baseline 或 A2C 的方法估计每一步的优势函数  $\hat{A}_i$ , 其中  $i = 1, 2, \dots, T$
- 使用均值估计期望近似目标函数  $J(\theta) \approx L(\theta) = \frac{1}{T} \sum_{i=1}^T \min[r(\theta) \cdot \hat{A}_i, \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \cdot \hat{A}_i]$
- 求解无约束优化问题  $\theta_{\text{new}} \leftarrow \text{argmax}_{\theta} L(\theta)$
- 每个 Episode 重复迭代 2-5 步

PPO 无需计算 KL 散度  $D_{KL}$ , 因此与 TRPO 相比减少了计算量; 同时 PPO 目标函数中的 clip 函数同样限制了策略网络参数的更新幅度

随机高斯策略网络

使用神经网络  $\mu(s; \theta^\mu)$  输出给定状态下各个动作分量的均值

使用神经网络  $\sigma(s; \theta^\sigma)$  输出给定状态下各个动作分量的方差

记每个网络的参数为  $\theta = (\theta^\mu, \theta^\sigma)$ , 则  $\ln \pi(\mathbf{a}|\mathbf{s}; \theta) = f(s, \mathbf{a}; \theta) + \text{const}$

等式两边分别对  $\theta$  求偏导, 可得策略梯度中对数策略函数对网络参数的导数等于辅助函数对网络函数的导数:

$\frac{\partial \ln \pi(\mathbf{a}|\mathbf{s}; \theta)}{\partial \theta} = \frac{\partial f(s, \mathbf{a}; \theta)}{\partial \theta}$

随机策略梯度  $g(s, \theta) = \frac{\partial \ln \pi(\mathbf{a}|\mathbf{s}; \theta)}{\partial \theta}$ ,  $Q_\pi(s, \mathbf{a}) = \frac{\partial V_\pi(s)}{\partial \theta} \cdot Q_\pi(s, \mathbf{a})$

- 使用 REINFORCE 更新参数  $\theta$
- 使用 Actor-Critic 更新参数  $\theta$

带基线的策略梯度  $g(s, \theta) = \frac{\partial \ln \pi(\mathbf{a}|\mathbf{s}; \theta)}{\partial \theta}$ ,  $Q_\pi(s, \mathbf{a}) = \frac{\partial V_\pi(s)}{\partial \theta} \cdot [Q_\pi(s, \mathbf{a}) - V_\pi(s)]$

- 使用 REINFORCE with baseline 更新参数  $\theta$
- 使用 A2C 更新参数  $\theta$

解决 DDPG 的高估问题

Clipped Double-Q Learning

- 使用两个价值网络  $q(s, a; \omega_1)$  和  $q(s, a; \omega_2)$
- 对应两个 target network  $q(s, a; \omega_1^-)$  和  $q(s, a; \omega_2^-)$
- 在计算 TD target 时,  $q_{t+1}$  的值选择两个 target network 输出中较小的值  $q_{t+1} = \min_{i=1,2} q(s_{t+1}, a_{t+1}; \omega_i^-)$

"Delayed" Policy Updates

- 策略网络的更新频率比价值网络低
- 每更新两次价值网络之后, 更新一次策略网络
- 每更新两次价值网络的 target network 之后, 更新一次策略网络的 target network

Target Policy Smoothing

- 在计算  $q_{t+1}$  时所用到的 target action  $a_{t+1}^-$  上添加高斯噪声
- 设  $\tilde{a}_{t+1}^- = \pi(s_{t+1}; \theta_t^-) + \text{clip}(\epsilon \cdot \omega, 0)$ , 其中  $\epsilon \sim \mathcal{N}(0, \sigma^2)$
- 通过给  $\tilde{a}_{t+1}^-$  添加噪声, 可以避免由价值网络估计误差导致的过度自信的问题
- 如果价值网络在更新过程中对某个动作的评分过高, 那么策略网络会误以为该动作最好, 这时给策略网络输出的动作添加噪声可以避免该动作, 从而避免该问题。

孪生延迟深度确定性策略梯度 (Twin Delayed DDPG, TD3)

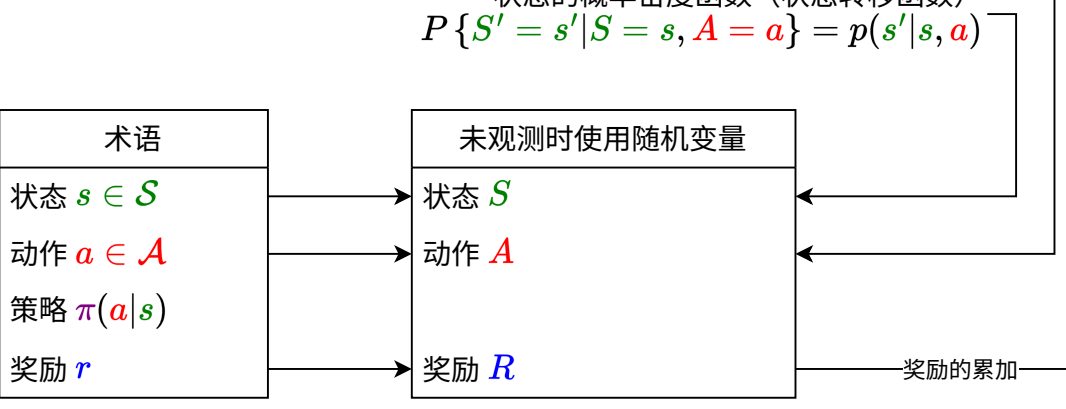
初始化策略网络  $\pi(s; \theta)$ ; 两个价值网络  $q(s, a; \omega_1)$ ,  $q(s, a; \omega_2)$

三个网络对应的 target network  $\pi(s; \theta^-)$ ,  $q(s, a; \omega_1^-)$ ,  $q(s, a; \omega_2^-)$ ; 其中  $\theta^- = \theta$ ,  $\omega_1^- = \omega_1$ ,  $\omega_2^- = \omega_2$

设置策略网络, 所有 target network 的更新频率为  $N_{\text{update}}$

- 根据当前状态  $s_t$  执行动作  $a_t = \pi(s_t; \theta)$ , 得到奖励  $r_t$  和下一个状态  $s_{t+1}$
  - 选择 target action  $\tilde{a}_{t+1}^- = \pi(s_{t+1}; \theta_t^-) + \text{clip}(\epsilon \cdot \omega, 0)$ , 其中  $\epsilon \sim \mathcal{N}(0, \sigma^2)$
  - 使用 target network 计算  $q_{t+1} = \min_{i=1,2} q(s_{t+1}, \tilde{a}_{t+1}^-; \omega_i^-)$
  - 计算 TD target  $y_t = r_t + \gamma \cdot q_{t+1}$
  - 分别计算两个价值网络的 TD error  $\delta_{i,t} = q(s_t, a_t; \omega_i) - y_t$ , 其中  $i = 1, 2$
  - 使用梯度下降更新两个价值网络  $\omega_{i,t+1} = \omega_{i,t} - \alpha \cdot \delta_{i,t} \cdot \frac{\partial q(s, a; \omega_i)}{\partial \omega_i} |_{s=s_t, a=a_t}$ , 其中  $i = 1, 2$
  - 如果当前 step 计数  $N_{\text{update}} = 0$ 
    - 使用第一个价值网络计算确定策略梯度  $g = \frac{\partial q(s, a; \omega_1)}{\partial \theta} = \frac{\partial q(s, a; \omega_1)}{\partial \theta} = \frac{\partial q(s, a; \omega_1)}{\partial \theta} = \frac{\partial q(s, a; \omega_1)}{\partial \theta}$
    - 使用策略梯度上升更新策略网络  $\theta_{t+1} = \theta_t + \beta \cdot g$
    - 更新所有 target network 参数
      - $\theta^- \leftarrow \tau \cdot \theta + (1 - \tau) \cdot \theta^-$
      - $\omega_i^- \leftarrow \tau \cdot \omega_i + (1 - \tau) \cdot \omega_i^-$ , 其中  $i = 1, 2$
- 实现算法时应当先判断、更新策略网络, 再更新价值网络, 因为更新策略网络需要用到价值网络更新前的参数  $\omega_{i,t}$

深度强化学习



奖励  $R_t$  与当前状态  $S_t$  和做出的动作  $A_t$  有关

在  $t$  时刻观测到状态  $s_t$  并做出动作  $a_t$

对未来的状态和动作  $S_{t+1}, S_{t+2}, \dots, A_{t+1}, A_{t+2}, \dots$  求期望

动作价值函数  $Q_\pi(s_t, a_t) = \mathbb{E}_{S_{t+1}, S_{t+2}, \dots, A_{t+1}, A_{t+2}, \dots} [U_t | S_t = s_t, A_t = a_t]$

对动作  $A$  求期望

状态价值函数  $V_\pi(s_t) = \mathbb{E}_A [Q_\pi(s_t, A)]$

表示当前状态  $s_t$  的评分

使用神经网络  $Q(s, a; \theta)$  近似策略函数  $Q_\pi(s, a)$

$V_\pi(s_t) = \mathbb{E}_A [Q_\pi(s_t, A)] = \sum_a \pi(a|s_t) \cdot Q_\pi(s_t, a)$

指导  $\pi(a|s; \theta)$  的训练

对于任意一个与动作  $A$  无关的 baseline 函数  $b$  具有下面的性质:

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \mathbb{E}_A \left[ \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \left[ b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} \right] = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta} = b \cdot \frac{\partial \ln \pi(A|s; \theta)}{\partial \theta}$

$\mathbb{E}_A \$