

策略  $\pi(a|s, \theta)$  参数  $\theta$  的随机梯度下降减小。

• 对于某一个  $s$ ,

$$\begin{aligned}\nabla J(\theta) &= \sum_a \nabla_{\theta} \pi(a|s, \theta) g_{\pi}(s, a) \\ &= \sum_a \pi(a|s, \theta) [\nabla_{\theta} \log \pi(a|s, \theta) g_{\pi}(s, a)] \\ &= E_a [\nabla_{\theta} \log \pi(a|s, \theta) g_{\pi}(s, a)] \quad ①\end{aligned}$$

依据证明可知,  $g_{\pi}(s, a) - b(s)$  与  $\nabla J(\theta)$  不变,

$$\Leftrightarrow E_a [\nabla_{\theta} \log \pi(a|s, \theta) (g_{\pi}(s, a) - b(s))] \quad ②$$

则: ① 和 ② 两个随机梯度, 具有相同的期望。

$\nabla_{\theta} \log \pi(a|s, \theta)$  具有方向性,  $g_{\pi}(s, a) - b(s)$  将降低了方差。

④ ⑤

• 减 baseline  $b(s)$ , 使得 随机变化值 更稳定。

$$\nabla J(\theta) = E_a [\text{随机变化值}]$$

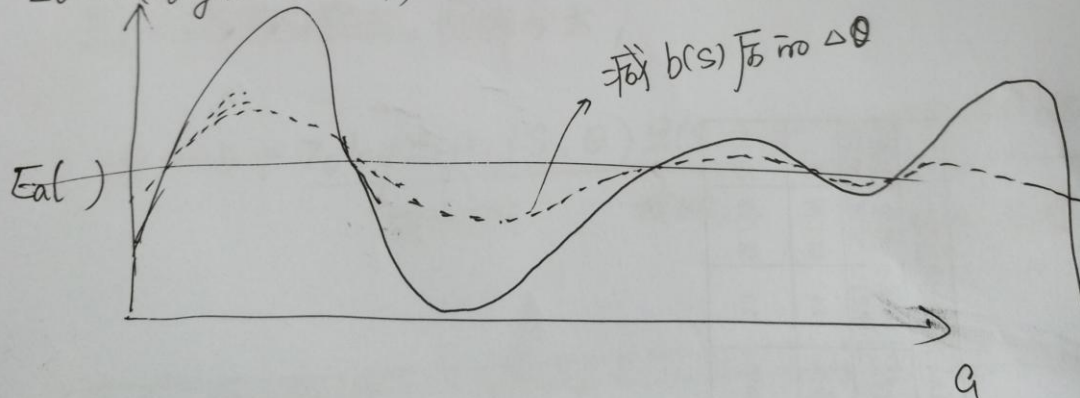
两方面进行讨论

① 同一  $s$  情况

② 不同  $s$  情况

(1) 初同S情况

$$\Delta\theta = \alpha (\nabla_{\theta} \log \pi(a|S, \theta) \mathcal{Q}(S, a))$$

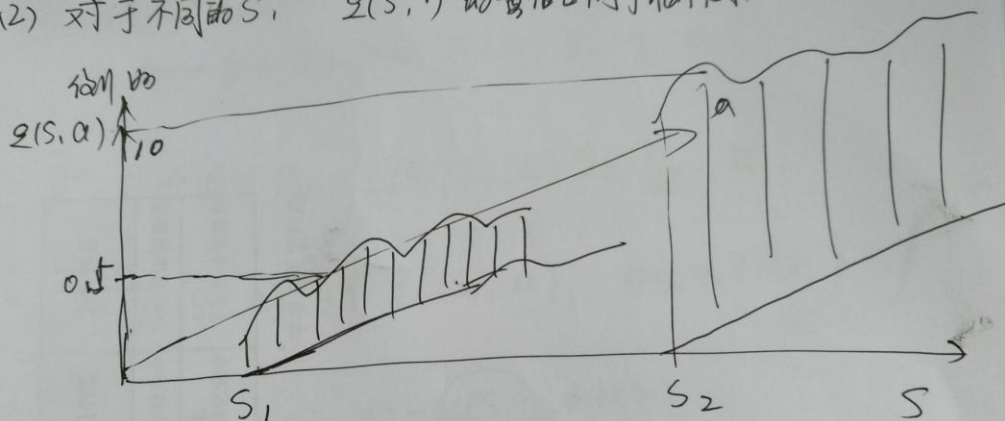


从效果上看, 可理解为: 如果  $\mathcal{Q}(S, a_1) > \mathcal{Q}(S, a_2)$

那么,  $\pi(a|S)$  的分布肯定倾向于选  $a_1$  大, 那么  $(S, a_1)$  被  
采样的 ~~概率会更大~~, 频率会大,

$$\theta \leftarrow \theta + \underbrace{\nabla_{\theta} \log \pi(a_1|S, \theta)}_{\text{朝该方向}} \underbrace{\mathcal{Q}(S, a_1)}_{\text{增加很大}} - \underbrace{\nabla_{\theta} \log \pi(a_1|S, \theta) b(s)}_{\text{抑制作用}}$$

12) 对于不同的  $S$ ,  $g(S, \cdot)$  的取值区间可能不同.



所以建议  $b(S_1) \neq b(S_2)$

⑤

~~导致情况如下~~: 不同的  $S$ , 对  $\theta$  值的变化影响大.

$$\theta \leftarrow \theta + \nabla_{\theta} \log \pi(a_i / \underline{S}_1, \theta) g_{\pi}(\underline{S}_1, a_i)$$

$$\text{和 } \theta \leftarrow \theta + \nabla_{\theta} \log \pi(a_i / \underline{S}_2, \theta) g_{\pi}(\underline{S}_2, a_i)$$

由于  $g_{\pi}(S_2, a) \gg g_{\pi}(S_1, a)$ ,  $\theta$  会以不同的量级变。  
 $\frac{g(S_2)}{g(S_1)}$

$$\checkmark \text{ 而非 } g_{\pi}(S_1, a) - \underline{V_{\pi}(S_1)} \rightarrow b(S_1)$$

$$g_{\pi}(S_2, a) - \underline{V_{\pi}(S_2)} \rightarrow b(S_2)$$

则  $g(S_1, \cdot)$  和  $g(S_2, \cdot)$  都在收敛.

⑥