

- $\omega_i' = \omega_i$ for $i = 1, 2, \dots, d$

- For $t = 0, 1, \dots, T-1$:

$$G_{t:T} = \sum_{k=t}^T (r^k \cdot \gamma^{k-t}) \quad // \text{ long term reward}$$

For $i = 1, 2, \dots, d$:

$$\omega_i' = \omega_i + \alpha \cdot (G_{t:T}) \cdot \text{derivative}_i(s^t, a^t, \omega = (\omega_1, \dots, \omega_d))$$

- $\text{derivative}_i(s, a, \omega = (\omega_1, \dots, \omega_d))$:

// derivative of $\ln \pi_\omega(s, a)$ wrt ω_i

$$\omega \cdot \text{dot} \cdot \text{phi} = 0.$$

for $i = 1, 2, \dots, d$:

$$\omega \cdot \text{dot} \cdot \text{phi} += \omega_i \cdot \phi_i(s)$$

$$\text{grad} = \frac{+1}{(1 + e^{-\omega \cdot \text{dot} \cdot \text{phi}})^2} \cdot e^{-\omega \cdot \text{dot} \cdot \text{phi}} \cdot \phi_i(s)$$

// equivalently $\text{grad} = \pi_\omega(s, 0) \cdot \pi_\omega(s, 1) \cdot \phi_i(s)$

if $a == 0$:

return $\text{grad} / \pi_\omega(s, 0)$

if $a == 1$:

return $-\text{grad} / \pi_\omega(s, 1)$