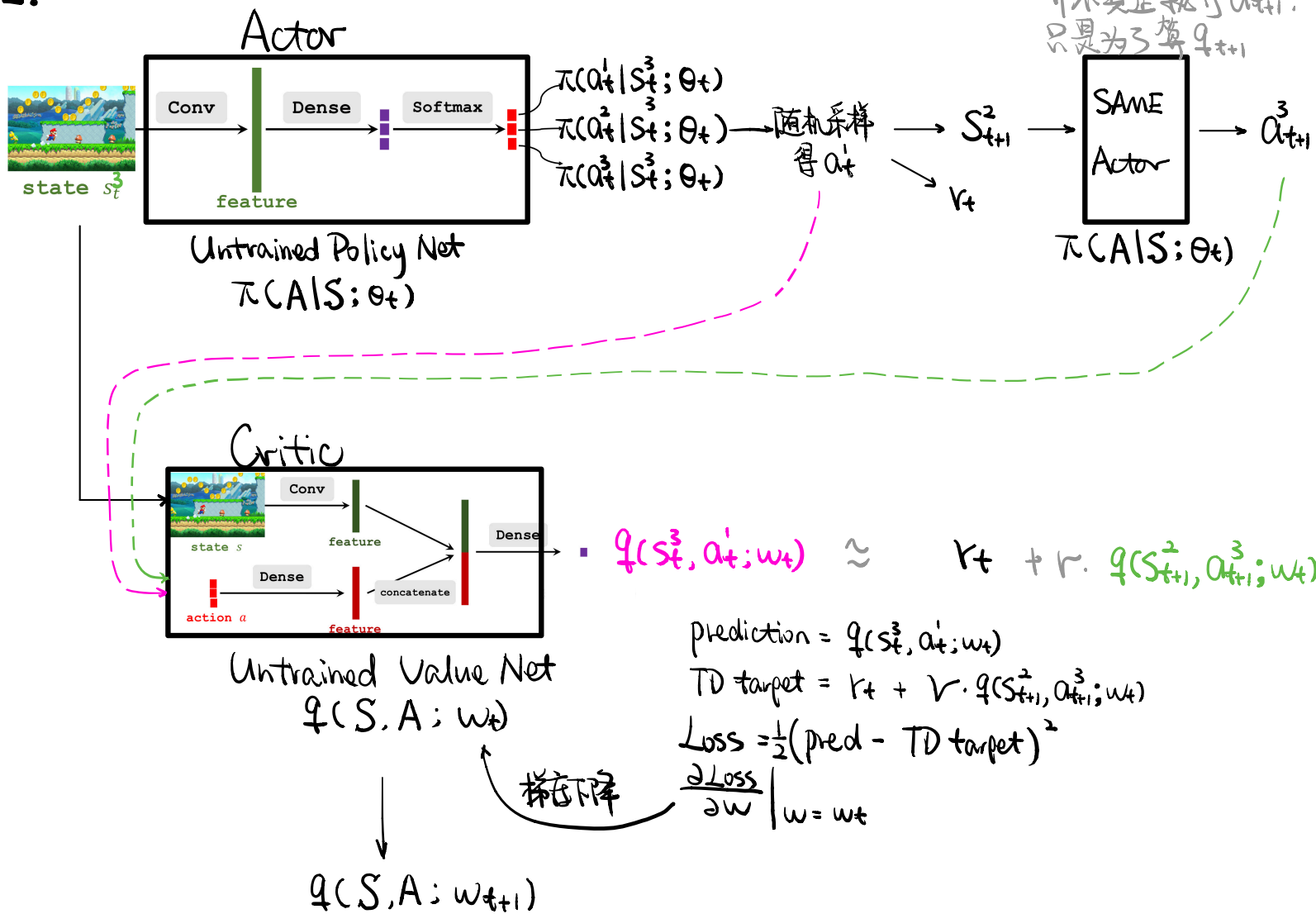
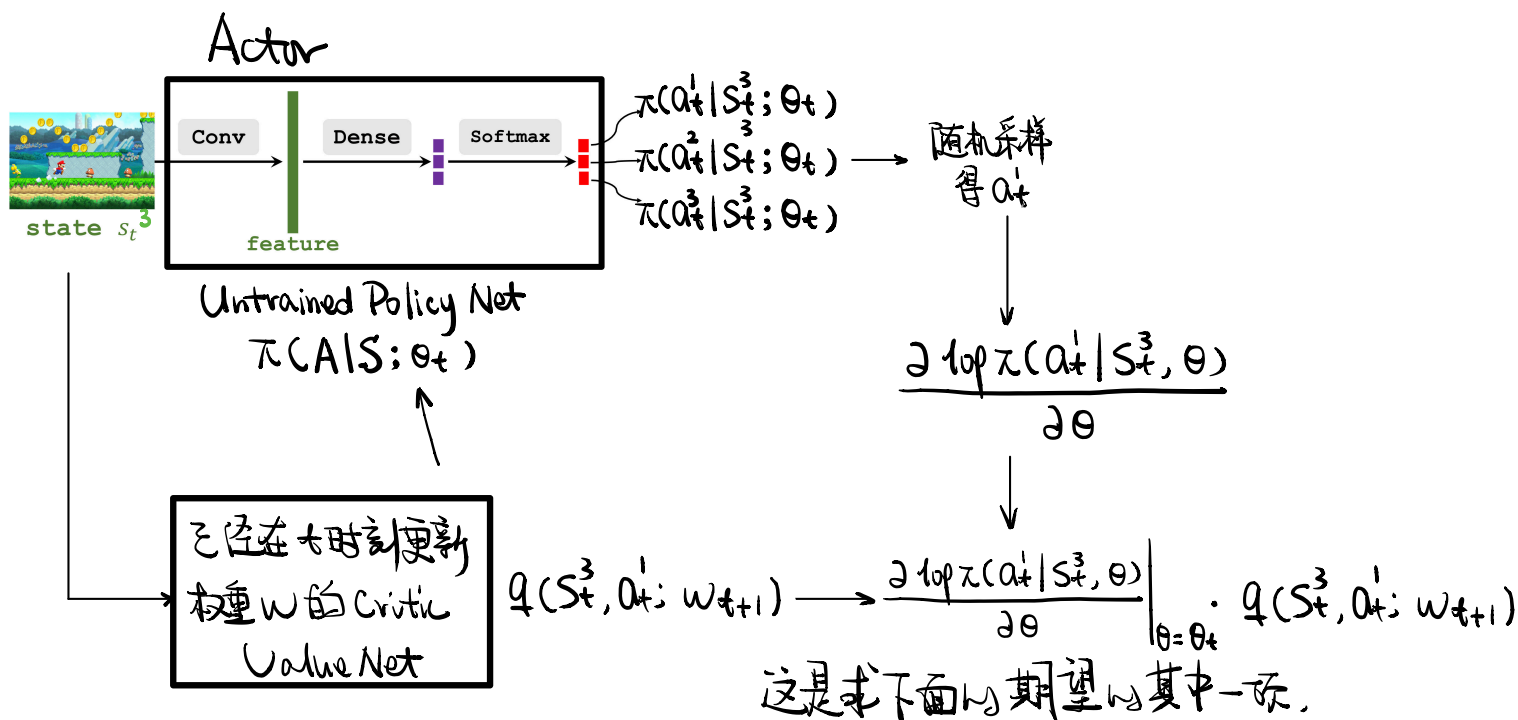


# 1. Update $w$ (in value network) using temporal difference (TD)



# 2. Update $\theta$ (in policy network) using policy gradient.



Form 2:  $\frac{\partial V(s; \theta)}{\partial \theta} = \mathbb{E}_{A \sim \pi(\cdot | s; \theta)} \left[ \frac{\partial \log \pi(A | s, \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right]$

梯度上升

我们用这一项来表示对期望的蒙特卡罗估计

$$\begin{aligned} \frac{\partial V(S_t^3; \theta)}{\partial \theta} &\approx \frac{\partial \log \pi(a_t^1 | S_t^3, \theta)}{\partial \theta} \Big|_{\theta=\theta_t} \cdot Q(S_t^3, a_t^1; w_{t+1}) \\ &= d_{\theta, t} \cdot q_t = g(a_t^1, \theta_t) \end{aligned}$$

Policy Net 训练时的“监督”完全来自 Critic, 即  $q(S_t^3, a_t^1; w_{t+1})$ .  
Critic 的好坏将会影响到 Policy Net 的训练效果;  
Critic 训练时的“监督”完全来自环境中的 Reward.

## Summary of Algorithm

1. Observe state  $s_t$  and randomly sample  $a_t \sim \pi(\cdot | s_t; \theta_t)$ .
2. Perform  $a_t$ ; then environment gives new state  $s_{t+1}$  and reward  $r_t$ .
3. Randomly sample  $\tilde{a}_{t+1} \sim \pi(\cdot | s_{t+1}; \theta_t)$ . (Do not perform  $\tilde{a}_{t+1}$ !)
4. Evaluate value network:  $q_t = q(s_t, a_t; w_t)$  and  $q_{t+1} = q(s_{t+1}, \tilde{a}_{t+1}; w_t)$ .
5. Compute TD error:  $\delta_t = q_t - (r_t + \gamma \cdot q_{t+1})$ .
6. Differentiate value network:  $d_{w, t} = \frac{\partial q(s_t, a_t; w)}{\partial w} \Big|_{w=w_t}$ .
7. Update value network:  $w_{t+1} = w_t - \alpha \cdot \delta_t \cdot d_{w, t}$ .
8. Differentiate policy network:  $d_{\theta, t} = \frac{\partial \log \pi(a_t | s_t, \theta)}{\partial \theta} \Big|_{\theta=\theta_t}$ .
9. Update policy network:  $\theta_{t+1} = \theta_t + \beta \cdot q_t \cdot d_{\theta, t}$ .