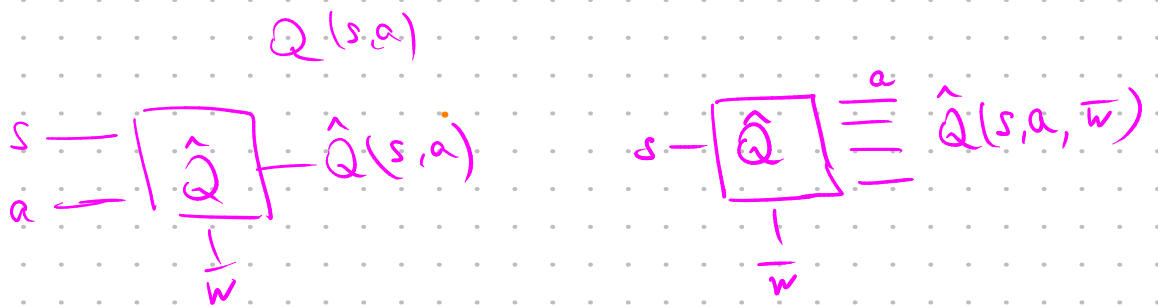
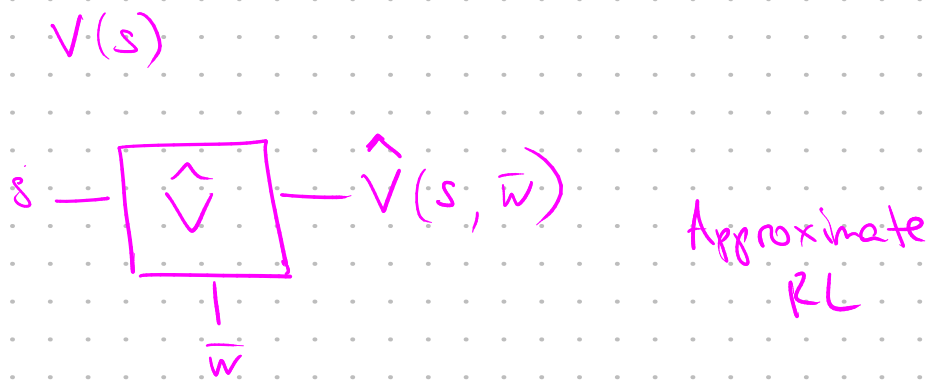
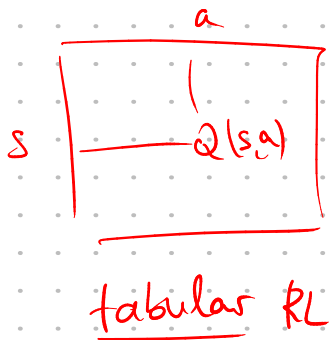


$$Q(s, a) = Q(s, a) + \alpha (r + \max_{a'} Q(s', a') - Q(s, a))$$



$$VE(\bar{w}) = \sum_s \underbrace{\mu(s)}_{\substack{\text{Pr}[s|b] \\ \uparrow \text{behaviour}}} (V_\pi(s) - \hat{V}_\pi(s, \bar{w}))^2 \xrightarrow{\bar{w}} \min$$

SGD:  $s_1, s_2, \dots, s_t, \dots$   $V_\pi(s_t)$  Gradient App. RL

$$\bar{w}_{t+1} = \bar{w}_t + \alpha (V_\pi(s_t) - \hat{V}_\pi(s_t, \bar{w}_t)) \cdot \nabla_{\bar{w}} \hat{V}(s_t, \bar{w}_t)$$

$$\bar{w}_{t+1} = \bar{w}_t + \alpha (V'_t - \hat{V}_\pi(s_t, \bar{w}_t)) \cdot \nabla_{\bar{w}} \hat{V}(s_t, \bar{w}_t)$$

1) Gradient MC estimation  $v(s_t) \approx V'_t = G_t$

-  $s_0, a_0, R_1, s_1, \dots, R_T, s_T$

-  $t \leq T-1, \dots \rightarrow 0$

$$\bar{w} := \bar{w} + \alpha (G_t - \hat{V}(s_t, \bar{w})) \nabla_{\bar{w}} \hat{V}(s_t, \bar{w})$$

2) Semi-gradient TD(0) estimation  $v(s_t) \approx V'_t = R_{t+1} + \gamma V(s_{t+1})$

-  $(s, a, s', r)$ :

$$\bar{w} := \bar{w} + \alpha (r + \gamma \hat{V}(s', \bar{w}) - \hat{V}(s, \bar{w})) \nabla_{\bar{w}} \hat{V}(s, \bar{w})$$

### 3) Semi-gradient TD control

$$\bar{w} := \bar{w} + \alpha (\hat{Q}' - \hat{Q}(s, a, \bar{w})) \nabla_{\bar{w}} \hat{Q}(s, a, \bar{w})$$

Sarsa:  $\bar{w} := \bar{w} + \alpha (r + \gamma \hat{Q}(s', a', \bar{w}) - \hat{Q}(s, a, \bar{w})) \nabla_{\bar{w}} \hat{Q}(s, a, \bar{w})$

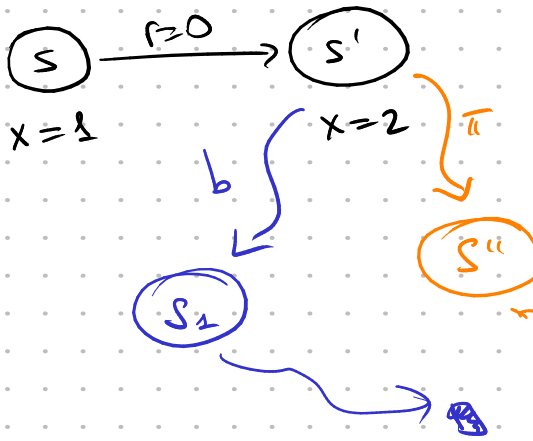
n-step  $R_{t+n} + \gamma R_{t+n+1} + \dots + \gamma^{k-1} R_{t+k} + \gamma^k \hat{Q}(s_{t+k}, A_{t+k}, \bar{w})$

off policy  $p_t = \frac{\pi(A_t | S_t)}{b(A_t | S_t)}$  importance weights

TD(0):  $\bar{w} := \bar{w} + \alpha p_t (r + \gamma \hat{V}(s', \bar{w}) - \hat{V}(s, \bar{w})) \nabla_{\bar{w}} \hat{V}(s, \bar{w})$

Expected Sarsa  $\bar{w} := \bar{w} + \alpha (r + \gamma \sum_{a'} \pi(a' | s') \hat{Q}(s', a', \bar{w}) - \hat{Q}(s, a, \bar{w})) \nabla_{\bar{w}} \hat{Q}(s, a, \bar{w})$

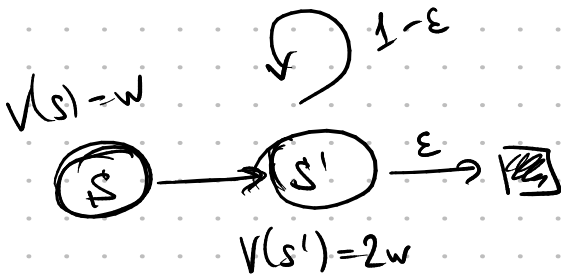
$V(s) = w$        $V(s') = 2w$



TD =  $\gamma \cdot 2w - w = (2\gamma - 1)w$

$w := w + \alpha p_t (2\gamma - 1)w \cdot \nabla_{\bar{w}} \hat{V}(s, \bar{w}) = (1 + \alpha(2\gamma - 1))w$

$V(s) = 0$



$w_{k+1} = \operatorname{argmin}_{\bar{w}} \sum_s (\hat{V}(s, \bar{w}) - \mathbb{E}[r + \gamma \hat{V}(s', \bar{w})])^2$

$= \operatorname{argmin}_{\bar{w}} \left[ (w - \gamma \cdot 2w_k)^2 + (2w - (1-\epsilon) \cdot \gamma \cdot 2w_k)^2 \right]$

$= \operatorname{argmin}_{\bar{w}} (w^2 + 4w_k^2 - 2w(2w_k\gamma + 4w_k\gamma(1-\epsilon))) + \dots$

$\gamma > \frac{5}{6-4\epsilon}$

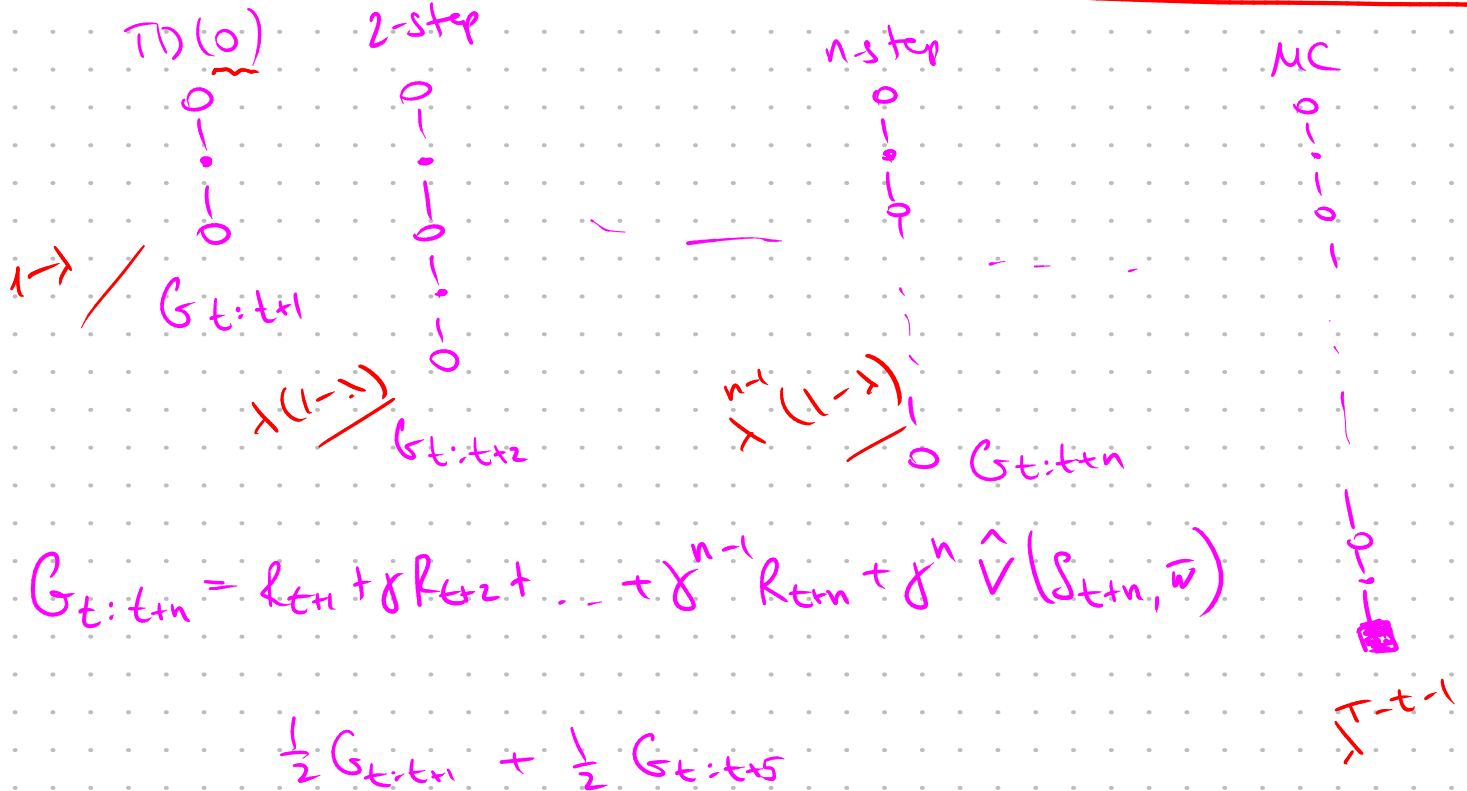
$w_{k+1} = \frac{6-4\epsilon}{5} w_k \gamma$

# The Deadly Triad

function approximation

bootstrapping

Off-policy training



$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{V}(S_{t+n}, \bar{w})$$

$$\frac{1}{2} G_{t:t+n} + \frac{1}{2} G_{t:t+5}$$

Eligibility traces

TD( $\lambda$ )

$$\begin{aligned} G_t^\lambda &= (1-\lambda) G_{t:t+1} + (1-\lambda) \cdot \lambda \cdot G_{t:t+2} + \\ &\quad + (1-\lambda) \lambda^2 G_{t:t+3} + \dots = \\ &= (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t:t+n} \end{aligned}$$