

① TD-learning

$(s_t, a_t, r_t, s_{t+1})$

оценка действия Q

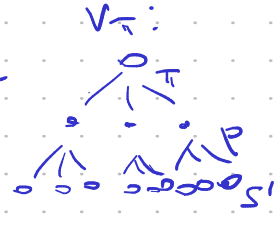
$$Q_{\pi}(S_t, A_t) := Q_{\pi}(S_t, A_t) + \alpha (R_{t+1} + \gamma Q_{\pi}(S_{t+1}, A_{t+1}) - Q_{\pi}(S_t, A_t))$$



backbone diagrams

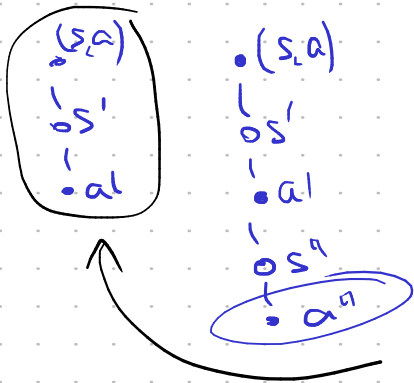


Bellman:

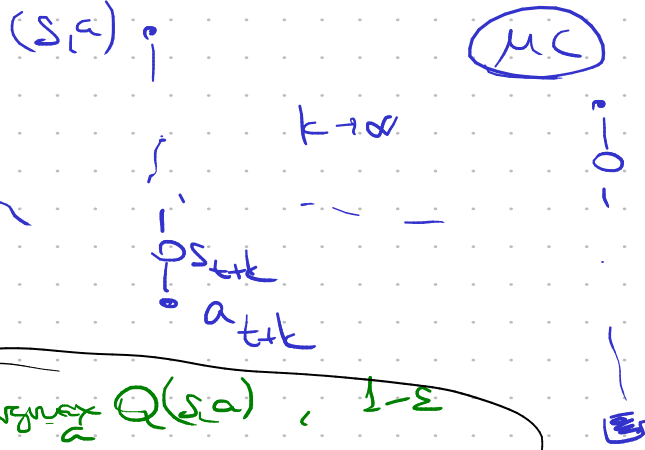


n-step TD:

$$Q_{\pi}(S_t, A_t) := \dots + \alpha (R_{t+1} + \gamma R_{t+2} + \gamma^2 Q_{\pi}(S_{t+2}, A_{t+2}) - Q_{\pi}(S_t, A_t))$$



SARSA



On-policy TD control:

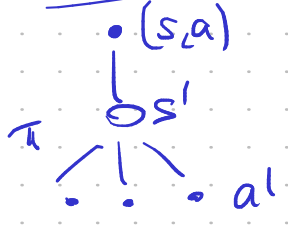
- bootstrap a no epa.
- $Q(s, a) := Q(s, a) + \alpha (r + \gamma \cdot Q(s', a') - Q(s, a))$

$$\pi(a|s) = \begin{cases} \arg \max_a Q(s, a), & 1-\epsilon \\ \text{unit}, & \epsilon \end{cases}$$

Expected  
Sarsa

$(s_t, a_t, r_t, s_{t+1}, a_{t+1})$

$$Q(s_t, a_t) := Q(s_t, a_t) + \alpha (r_t + \gamma \cdot \sum_{a'} \pi(a'|s_{t+1}) Q(s_{t+1}, a') - Q(s_t, a_t))$$



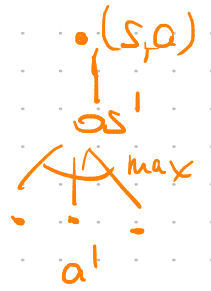
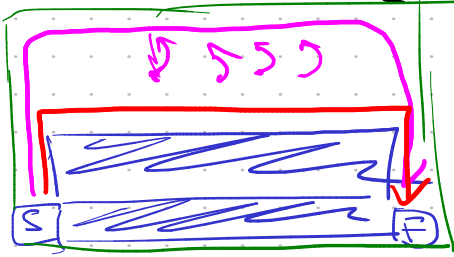
$$r_t + \gamma \cdot \sum_{a_1 \neq a'} \pi(a_1|s_{t+1}) \cdot Q(s_{t+1}, a_1) +$$

$$+ \gamma \cdot \pi(a'|s_{t+1}) \cdot (r_{t+1} + \gamma \cdot \sum_{a_2} \pi(a_2|s_{t+2}) Q(s_{t+2}, a_2))$$



② Q-learning (1989, Watkins) - off policy TD control

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



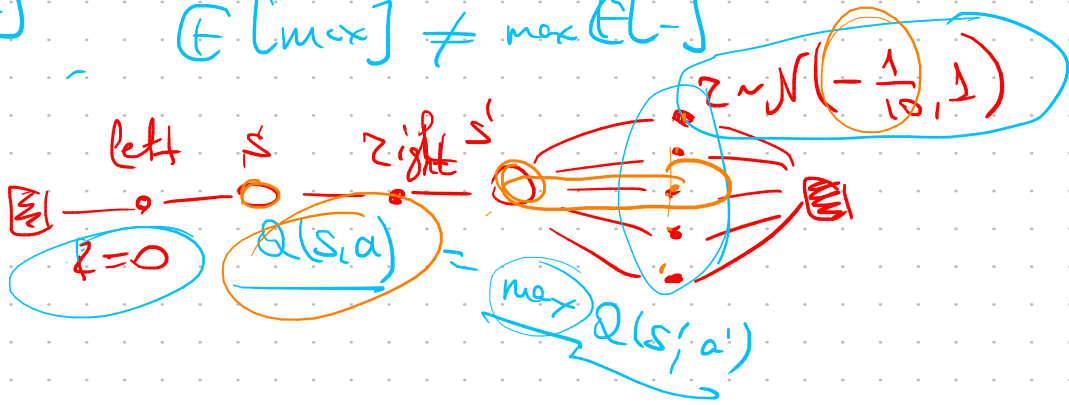
$Q(s,a) \rightarrow r + \gamma \cdot \max_{a'} Q(s',a')$

$E[Q_1(s, a) \max_{a'} Q_2]$

$\pi(s) = \operatorname{argmax}_a Q(s,a)$

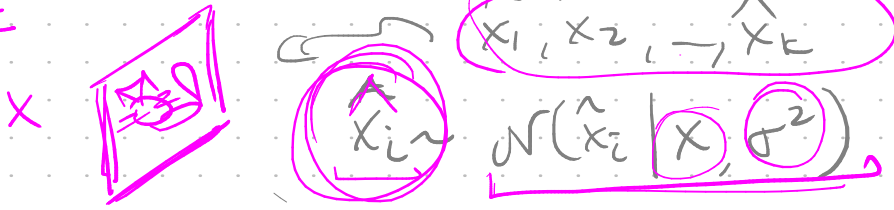
$E_{\pi}[G(s,a)]$

$E[\max] \neq \max E[\ ]$



Winner's curse

Mechanism design



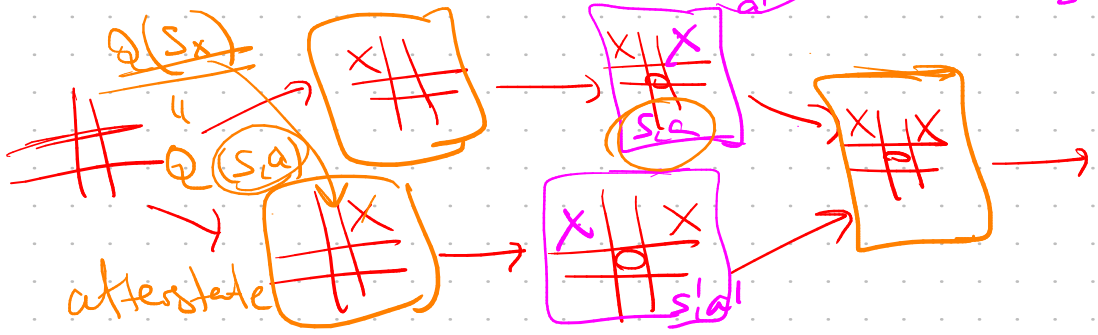
Double Q-learning (double DQN)

$Q_1(s,a) \quad Q_2(s,a)$

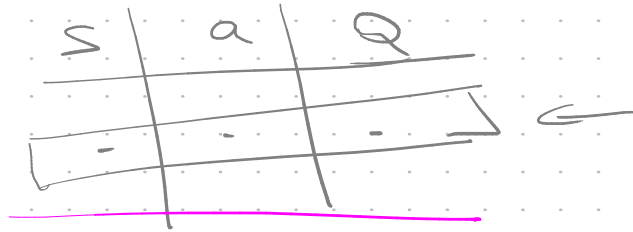
(s,a, z, s'): -  $p = \frac{1}{2}$ :  $Q_1(s,a) := Q_1(s,a) + \alpha [r + \gamma \cdot Q_2(s', \operatorname{argmax}_{a'} Q_2(s', a')) - Q_1(s,a)]$

-  $p = \frac{1}{2}$ :  $Q_2(s,a) := Q_2(s,a) + \alpha [r + \gamma \cdot Q_1(s', \operatorname{argmax}_{a'} Q_1(s', a')) - Q_2(s,a)]$

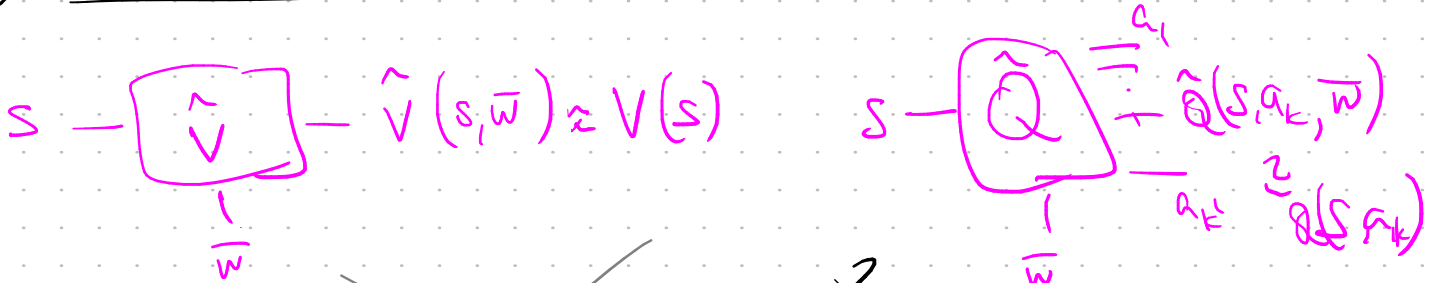
③ Afterstates



Tabular RL



④ Approximate RL



$$L(\bar{w}) = \sum_{s \in \mathcal{S}} (V_{\pi}(s) - \hat{V}_{\pi}(s, \bar{w}))^2 \xrightarrow{\bar{w}} \min$$

GD  
 $F(\bar{x}) \rightarrow \min$

$$F(\bar{x}) = \frac{1}{N} \sum_{n=1}^N \ell(d_n, \bar{x}) = \mathbb{E}_{\text{unif}}[\ell(d, \bar{x})]$$

$$\bar{x} := \bar{x} - \alpha \cdot \nabla_{\bar{x}} F(\bar{x})$$

SGD: mini-batch stochastic gradient descent

$$\approx \frac{1}{M} \sum_{m=1}^M \ell(d_m, \bar{x})$$

$$L(\bar{w}) = \sum_{s \in \mathcal{S}} \mu(s) \cdot (V_{\pi}(s) - \hat{V}_{\pi}(s, \bar{w}))^2 \xrightarrow{\bar{w}} \min$$

$\mu(s)$  = "given by policy  $\pi$  on  $\mathcal{S}$ "

$$\mu(s) = \mathbb{P}_{\pi}[S_t = s] \quad \mathbb{E}_{\pi}[e + \gamma \cdot V_{\pi}(s', \bar{w})]$$

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

$$\bar{w} := \bar{w} - \alpha \nabla_{\bar{w}} \ell(s, \bar{w}) =$$

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

$$G_t(\mu) \leftarrow e + \gamma \cdot \hat{V}_{\pi}(s', \bar{w}) \quad (TD)$$

Gradient MC estimation:

③ Off-policy MC

$$\bar{w} := \bar{w} + \alpha (G_t - \hat{V}_{\pi}(s, \bar{w})) \cdot \nabla_{\bar{w}} V_{\pi}(s, \bar{w})$$

Semi-gradient TD estimation  $\approx V(S_t)$

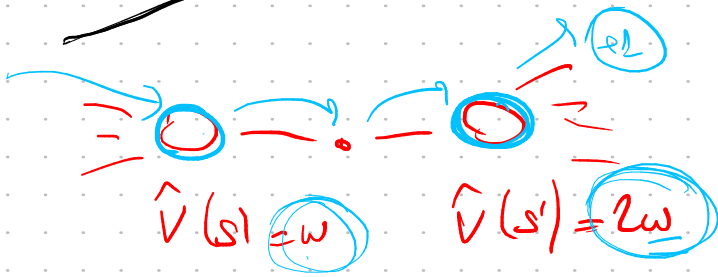
$$\bar{w} := \bar{w} + \alpha (R_{t+1} + \gamma \hat{V}_\pi(S_{t+1}, \bar{w}) - \hat{V}_\pi(S_t, \bar{w})) \cdot \nabla_{\bar{w}} \hat{V}(S_t, \bar{w})$$

Semi-gradient DDPG

*n-step, expected DDPG*

$$\bar{w} := \bar{w} + \alpha (R_{t+1} + \gamma \hat{Q}(S_{t+1}, A_{t+1}, \bar{w}) - \hat{Q}(S_t, A_t, \bar{w})) \cdot \nabla_{\bar{w}} \hat{Q}(S_t, A_t, \bar{w})$$

~~Off-policy semi-gradient TD control~~



Deadly triad

- approximation
  - bootstrapping
  - off-policy
- Q-learning



experience replay

5 Policy gradient

~~$s - Q \equiv Q(s,a) \} \text{argmax}$~~

$s - \pi \equiv \pi(a|s, \theta) = P_r(A_t = a | S_t = s)$

$\theta$

$s - \pi \equiv h(a, s, \theta) - \int_{\text{stochastic}} - \pi(a|s)$

softmax cost

$$\pi(a|s, \theta) = \frac{e^{h(a, s, \theta)}}{\sum_{a'} e^{h(a', s, \theta)}}$$

J( $\bar{\theta}$ ) =  $V_{\pi_{\bar{\theta}}}(s_0) = \mathbb{E}_{\pi_{\bar{\theta}}}[G_t | s_0] \xrightarrow{\bar{\theta}} \text{max}$

$\bar{\theta} := \bar{\theta} + \alpha \cdot \nabla_{\bar{\theta}} J(\bar{\theta}) = \mathbb{E}_{\pi_{\bar{\theta}}}(G_t)$

Policy gradient theorem



$$\begin{aligned} \nabla_{\theta} V_{\pi_{\theta}}(s) &= \nabla_{\theta} \left[ \sum_a \pi_{\theta}(a|s) \cdot Q_{\pi_{\theta}}(s,a) \right] = \\ &= \sum_a \left[ (\nabla_{\theta} \pi_{\theta}(a|s)) \cdot Q_{\pi_{\theta}}(s,a) + \pi_{\theta}(a|s) \cdot \nabla_{\theta} Q_{\pi_{\theta}}(s,a) \right] = \end{aligned}$$

$$\begin{aligned}
&= \sum_a \left[ Q \cdot \nabla_{\pi} (a|s) + \pi_{\theta}(a|s) \cdot \nabla_{\theta} \left[ \sum_{z, s'} p(z, s'|s, a) (z + V_{\pi_{\theta}}(s')) \right] \right] \\
&= \sum_a \left[ Q \cdot \nabla_{\pi} (a|s) + \pi_{\theta}(a|s) \cdot \sum_{z, s'} p(z, s'|s, a) \left( \cancel{z} + \nabla_{\theta} V_{\pi_{\theta}}(s') \right) \right] \\
&= \sum_a \left[ \underbrace{(Q \cdot \nabla_{\pi}) (a|s)} + \pi_{\theta}(a|s) \sum_{z, s'} p(z, s'|s, a) \times \right. \\
&\quad \left. \times \sum_{a'} \left[ \underbrace{(Q \cdot \nabla_{\pi}) (a', s')} + \pi_{\theta}(a'|s') \cdot \nabla_{\theta} Q_{\pi_{\theta}}(s', a') \right] \right] \\
&= \sum_{s'} \left[ \left( \sum_a \underbrace{Q_{\pi_{\theta}}(s', a)} \cdot \nabla_{\pi_{\theta}} (a|s') \right) \cdot \sum_{k=0}^{\infty} p_r \left[ \begin{array}{c} y \text{ } s \rightarrow s' \\ \text{no } \pi \\ \text{no } \pi \end{array} \right] \right]
\end{aligned}$$

2-ü war:  $\left( \frac{s'}{\pi} \right) \times \left( \sum_a \pi_{\theta}(a|s) - \sum_{z'} p(z, s'|s, a) \right)$

$$\nabla_{\theta} J(\bar{\theta}) = \nabla_{\theta} V_{\pi}(s_0) = \sum_s \left( \sum_{t=0}^{\infty} p_r \left[ \begin{array}{c} s_0 \rightarrow s \\ \text{no } \pi \\ \text{no } \pi \end{array} \right] \right) \left( \sum_a \underbrace{Q_{\pi_{\theta}}(s, a)} \nabla_{\theta} \pi_{\theta}(a|s) \right)$$

$E \# [t: s_t = s] \propto \mu(s)$

$$\nabla_{\theta} J(\bar{\theta}) \propto \sum_s \mu(s) \sum_a \underbrace{Q_{\pi_{\theta}}(s, a)} \nabla_{\theta} \pi_{\theta}(a|s)$$

1) All-actions method (actor-critic)

$$\bar{\theta} := \bar{\theta} + \alpha \cdot \sum_a \hat{Q}(s, a | \bar{w}) \cdot \nabla_{\theta} \pi(a|s, \bar{\theta})$$

2) REINFORCE (Williams, 1992)

$$\sum_a \pi_{\theta}(a|s_t)$$

$$\nabla_{\theta} J(\theta) \propto E_{\pi_{\theta}} \left[ \sum_a Q_{\pi_{\theta}}(s_t, a) \nabla_{\theta} \pi_{\theta}(a|s_t) \right] =$$

$$= E_{\pi} \left[ \sum_a \pi_{\theta}(a|s_t) \cdot Q_{\pi_{\theta}}(s_t, a) \cdot \frac{\nabla_{\theta} \pi_{\theta}(a|s_t)}{\pi_{\theta}(a|s_t)} \right] =$$

$$= E_{\pi_{\theta}} \left[ Q_{\pi_{\theta}}(s_t, A_t) \cdot \frac{\nabla_{\theta} \pi_{\theta}(A_t|s_t)}{\pi_{\theta}(A_t|s_t)} \right] =$$

$$= E_{\pi_{\theta}} [G_t | s_t, A_t]$$

$$= E_{\pi_{\theta}} \left[ G_t \cdot \frac{\nabla_{\theta} \pi_{\theta}(A_t|s_t)}{\pi_{\theta}(A_t|s_t)} \right] \quad \checkmark \text{ REINFORCE}$$

$$\bar{\theta} := \bar{\theta} + \alpha \cdot G_t \cdot \nabla_{\theta} [\log \pi_{\theta}(A_t|s_t, \bar{\theta})]$$

3) REINFORCE w/ baselines

$$\nabla J(\bar{\theta}) = \sum_s \mu(s) \sum_a (Q_{\pi}(s, a) - b(s)) \nabla_{\theta} \pi(a|s)$$

$$\sum_s \mu(s) \sum_a (Q - b(s)) \cdot \nabla_{\theta} \pi(a|s) =$$

$$= \sum_s \dots - \sum_s \mu(s) \sum_a b(s) \cdot \nabla_{\theta} \pi(a|s)$$

$$= \sum_a \nabla_{\theta} \pi(a|s) =$$

$$\nabla_b \sum_a \pi(a|s) = \nabla_b 1 = 0$$

Actor-critic

$$\nabla J(\bar{\theta}) \propto \sum_s \mu(s) \sum_a (Q_{\pi}(s, a) - \hat{V}(s, \bar{w})) \cdot \nabla_{\theta} \pi(a|s)$$

- gen.  $S_0, A_0, r_1$  — — — — —  $G=0$

- for  $t = T-1, \dots, 0$

$$- G = R_t + \delta G$$

$$- \bar{w}_t = \bar{w} + \alpha_w (G - \hat{V}(S_t, \bar{w})) \nabla_{\bar{w}} \hat{V}(S_t, \bar{w})$$

$$- \bar{\theta} := \bar{\theta} + \alpha_{\theta} (G - \hat{V}(S_t, \bar{w})) \nabla_{\bar{\theta}} \log \pi(A_t | S_t, \bar{\theta})$$