**Planning** — minimax-tree, rollouts, MCTS

Bellman equations, Experience replay

$V, Q, \pi$

**Direct RL** — MC-methods, TD-обучение

**Model**

Model learning $p(r, s' | s, a)$

Experience $(s, a, r, s')$ , $a'$ → $(x, y)$

**Bellman eq:**

$s$ — $\hat{V}$ — $\hat{V}(s)$ , $\bar{w}$

---

## Approximate RL

без ограничения потери состояний

$$L(\bar{w}) = \sum_{s \in S} \mu(s) \cdot \left(V(s) - \hat{V}(s, \bar{w})\right)^2$$

$$\nabla_{\bar{w}} L = -2 \sum_{s} \mu(s) \left(V(s) - \hat{V}(s, \bar{w})\right) \nabla_{\bar{w}} \hat{V}(s, \bar{w})$$

$s_t:$
$$\bar{w} := \bar{w} + \alpha \left(V(s_t) - \hat{V}(s_t, \bar{w})\right) \nabla_{\bar{w}} \hat{V}(s_t, \bar{w}) \quad \text{Stochastic GD}$$

$$\approx U_t$$

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

$$E[U_t | S_t = s] = V(s)$$

$U_t = G_t \implies$ **Gradient MC**

evaluation ← → control

on-policy →→ off-policy

$\pi, \quad V_\pi(s) = ?$

– generate $S_0, A_0, R_1, S_1 \cdots, S_T$

– $\forall t = T-1, \cdots, 0$

  – $\bar{w} := \bar{w} + \alpha \left(G_t - \hat{V}(S_t, \bar{w})\right) \nabla_{\bar{w}} \hat{V}(S_t, \bar{w})$
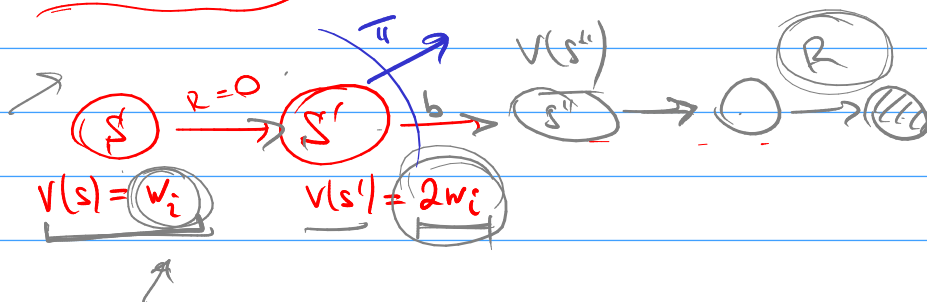
$$\rho_t = \frac{\pi(A_t | S_t)}{b(A_t | S_t)}$$

**Approximate TD**

$$\overline{w} := \overline{w} + \alpha \left( \overbrace{R + \gamma \hat{V}(s', \overline{w}) - \hat{V}(s, \overline{w})}^{U_t} \right) \cdot \nabla_{\overline{w}} \hat{V}(s, \overline{w})$$

Semi-gradient TD

~~Gradient TD~~

off-policy TD



$V(s) = w_i$    $V(s') = 2w_i$

The deadly triad:
1) approximation
2) bootstrapping
3) off-policy

---

Eligibility traces

$\lambda = 0$         $\lambda = 1$



TD(0)

$G_{t:t+3}$

$G_{t:T}$

$G_t$
MC

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \_$$
$$\_ + \gamma^{n-1} R_{t+n} +$$
$$+ \gamma^n V(\hat{S}_{t+n})$$

$$G_t^\lambda = (1-\lambda) G_{t:t+1} + (1-\lambda)\lambda G_{t:t+2} + (1-\lambda)\lambda^2 G_{t:t+3} + \_ \_ \_$$
$$=$$

$$\bar{w} := \bar{w} + \alpha \left( \underline{G_t^\lambda} - \hat{V}(S_t, \bar{w}) \right) \nabla_{\bar{w}} \hat{V}(S_t, \bar{w})$$

## Semi-gradient TD($\lambda$):

- init $S_0$, $\bar{z}_{t=0}$
- $\forall t = 0, \dots T:$
  - b.o.Sup $A$, nerve $R, S'$ & $\hat{S}$
  - $\bar{z} := \gamma \cdot \lambda \bar{z} + \nabla_{\bar{w}} \hat{V}(S_t, \bar{w})$
  - $\bar{w} := \bar{w} + \alpha \left( R + \gamma \hat{V}(S', \bar{w}) - \hat{V}(S, \bar{w}) \right) \bar{z}$

$\lambda = 0 \Rightarrow$ TD
$\lambda = 1 \Rightarrow$ MC

---

# Policy gradient

$S$ — $\boxed{\Pi}$ — $\vdots$     $\pi(a|s, \theta) = \Pr\left[ A_t = a \mid S_t = s, \bar{\theta} \right]$

$\bar{\theta}$

$Q, V$ ~~state-value (action-value) methods~~

$J(\theta) = ?$

$\bar{\theta} := \bar{\theta} + \alpha \nabla_{\bar{\theta}} J(\bar{\theta})$

$S$ — $\boxed{\Pi}$ — $\pi(a|s, \theta)$
$\theta$

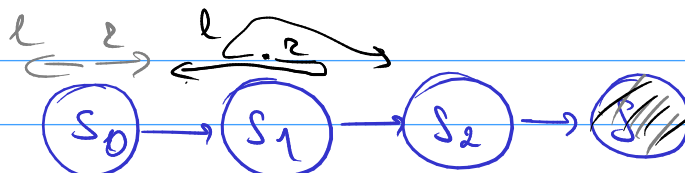## Actor-critic methods : $\Pi(a|s, \bar{\theta}) + V(s, \bar{w})$

$S$ — $\boxed{\Pi} \vdots h(a, s, \theta)$ — $\boxed{\text{softmax}}$ — $\pi(a|s, \theta) = \dfrac{e^{h(a,s,\theta)}}{\sum\limits_{a'} e^{h(a', s, \theta)}}$
$\theta$

$S$ — $\boxed{Q_*} = Q(s,a)_* \Rightarrow \pi := \arg\max\limits_a Q_*(s,a)$



$S_0 \rightarrow S_1 \rightarrow S_2 \rightarrow \bigcirc$

$$J(\bar{\theta}) = V_{\pi_{\bar{\theta}}}(s_0) \qquad \nabla_{\bar{\theta}} J(\bar{\theta}) = ?$$

$$V_{\pi_\theta}(s) = \sum_a \pi(a|s,\theta) \, Q_{\pi_\theta}(s,a)$$

$$\nabla_\theta V_{\pi_\theta}(s) = \sum_a \left( \left( \nabla_\theta \pi(a|s,\theta) \right) \cdot Q_{\pi_\theta}(s,a) + \pi(a|s,\theta) \cdot \nabla_\theta Q_{\pi_\theta}(s,a) \right) =$$

$$= \sum_a \left[ \nabla_\theta \pi \cdot Q_\pi(s,a) + \pi(a|s,\theta) \cdot \nabla_\theta \sum_{s',z} p(s',z|s,a)(z + \gamma V_{\pi_\theta}(s')) \right] =$$

$$= \sum_a \left( \nabla_\theta \pi \cdot Q_\pi(s,a) + \pi(a|s,\theta) \cdot \sum_{s',z} p(s|z|s,a) \left( z + \underset{\gamma<1}{\nabla_\theta V_{\pi_\theta}(s')} \right) \right) =$$

$$= \sum_a \left( Q_\pi(s,a) \cdot \nabla_\theta \pi(a|s) \right) + \pi(a|s,\theta) \cdot \sum_{s',z} p(s',z|s,a) \left( z + \right.$$

$$+ \sum_{a'} \left( \nabla_\theta \pi(a'|s') \cdot Q_\pi(s',a') \right) + \pi(a'|s',\theta) \cdot$$

$$\left. \cdot \nabla_\theta Q_\pi(s',a') \right) =$$

$$= \sum_{s'' \in S} \sum_{k=0}^{\infty} Pr\left[ s \to s' \text{ за } k \text{ шагов} \atop \text{по } \pi \right] \cdot \left( \sum_a \nabla_\theta \pi(a|s') \, Q_\pi(s',a) \right)$$

$$\nabla_\theta J(\theta) = \nabla_\theta V_\pi(s_0) = \sum_s \left( \sum_{t=0}^{\infty} Pr\left[ s_0 \to s' \text{ за } t \text{ шагов} \atop \text{по } \pi \right] \right) \cdot \sum_a \nabla \pi \cdot Q_\pi$$

$$= \sum_s \boxed{\eta(s)} \cdot \sum_a \nabla_\theta \pi(a|s) \cdot Q_\pi(s,a) \qquad \propto$$

$$= \mathbb{E}\left[ \# \text{ попад. в } s \right]$$

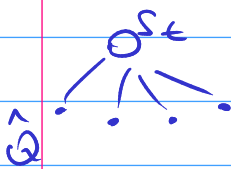$$\propto \sum_s \boxed{\mu(s)} \cdot \sum_a \nabla_\theta \pi(a|s) \, Q_\pi(s,a)$$

доля попад. состоян. $s$

Policy gradient theorem

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

$$= \mathbb{E}_\pi \left[ \sum_a Q_\pi(s,a) \nabla_\theta \pi(a|s,\theta) \right]$$

① **All-actions**

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



② **REINFORCE**

$$\nabla J(\bar{\theta}) \propto \mathbb{E}_\pi \left[ \sum_a Q_\pi(S_t, a) \nabla_\theta \pi(a|S_t,\theta) \right] =$$

$$= \mathbb{E}_\pi \left[ \sum_a \pi(a|s) \cdot \frac{1}{\pi(a|s)} Q_\pi(S_t, a) \nabla_\theta \pi(a|S_t,\theta) \right] =$$

$$[A_t \sim \pi] = \mathbb{E}_\pi \left[ Q_\pi(S_t, A_t) \cdot \frac{\nabla_\theta \pi(A_t|S_t,\theta)}{\pi(A_t|S_t,\theta)} \right] =$$

$$Q_\pi(S_t, A_t) = \mathbb{E}_\pi [G_t | S_t, A_t]$$

$$J(\theta) \propto \mathbb{E}_\pi \left[ G_t \cdot \frac{\nabla_\theta \pi(A_t|S_t,\theta)}{\pi(A_t|S_t,\theta)} \right]$$

$$\bar{\theta} := \bar{\theta} + \alpha G_t \cdot \frac{\nabla_\theta \pi(A_t|S_t,\theta)}{\pi(A_t|S_t,\theta)} = \bar{\theta} + \alpha G_t \nabla_\theta \ln \pi(A_t|S_t,\theta)$$

**REINFORCE:**
- loop
  - gen. ep. $S_0, A_0, R_1, S_1, \dots, S_T$
  - & $t = 0, \dots, T-1$:
    - $G := R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T$

$$- \theta := \theta + \alpha \, \gamma^t \, G \, \nabla \ln \pi (A_t | S_t, \theta)$$

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

$b(s)$ — baseline

$$\sum_s \mu(s) \sum_a \left( Q_\pi(s, a) - b(s) \right) \nabla_\theta \pi(a | s, \theta) =$$

$$= \sum_s \mu(s) \left( \sum_a Q \nabla_\pi - \sum_a b(s) \, \nabla_\theta \pi(a | s, \theta) \right)$$

$$b(s) \cdot \nabla_\theta \left( \underbrace{\sum_a \pi(a | s, \theta)}_{\sim 1} \right)$$

**REINFORCE w/ baseline:**

$$- \theta := \theta + \alpha \left( G_t - \underbrace{b(S_t)}_{\text{``} \hat{V}(S_t, \bar{w})} \right) \cdot \nabla_\theta \ln \pi (A_t | S_t, \theta)$$

$$- \bar{w} := \bar{w} + \alpha' \left( \underbrace{G_t}_{R_{t+1} + \gamma \hat{V}(S_{t+1}, \bar{w})} - \hat{V}(S_t, \bar{w}) \right) \nabla_{\bar{w}} \hat{V}(S_t, \bar{w})$$