$$\mathcal{S}_t, A_t$$

$$p(s', z | s, a) = \Pr[R_t = z, \mathcal{S}_t = s' | S_{t-1} = s, A_{t-1} = a]$$

over it: $p(z|s,a)$   $p(s'|s,a)$

$$\pi: \mathcal{S} \to \Delta[A(s)] \quad, \quad \pi(a|s)$$

$$G_t = \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k = R_{t+1} + \gamma G_{t+1}$$

$$V_\pi(s) = \mathbb{E}_\pi[G_t | \mathcal{S}_t = s]$$

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

Bellman equations

$$V_*(s) = \max_\pi V_\pi(s) \quad ; \quad Q_*(s,a) = \max_\pi Q_\pi(s,a)$$

Value iteration  $V_*(s)$

Policy iteration



$V = V_\pi$

$\pi, V$

$V_*, \pi_*$

$\pi := Greedy(V)$

---

① Monte-Carlo estimation        $\pi, \boxed{V_\pi} \; Q_\pi \quad p(z, s' | s, a)$

— init $\underline{\pi}$, $V_\pi(s)$, Returns(s) := []

— loop:

    — поройм. эпизод по $\pi$ , $\quad G := 0$        $v(s)$

    — $\forall t = T-1, T-2, \ldots, 0$

    — $G := \gamma G + R_{t+1}$            first-visit MC

—init random  
$\boxed{(S_0, a_0)}$  $S_0$

Exploring starts

    — если надо:        every-visit MC

    — Returns($s_t$).append( $G$ )         $V_\pi(S_t), \; N(S_t)$

    — $V_\pi(S_t) := Avg(Returns(s_t))$     $V_\pi(S_t) =$

    $= V_\pi(S_t) + \frac{1}{N}(G - V_\pi(S_t))$

$Q_\pi(s,a)$    сохранять  
    — exploration

---

② On-policy MC control       $\forall a, s$   $\pi(a|s) \geq \varepsilon / |A(s)|$

— —"—   $\pi$ — $\varepsilon$-мягкой ,

— loop:

    — —"— ,  $G := 0$       $\pi: s \to$ распр. на $A(s)$

    — $\forall t = T-1, T-2, \ldots, 0$.

    — $G = \gamma G + R_{t+1}$     $Q_*(S, A) \to \mathbb{R}$

— если надо:
— Returns$(S_t, a_t)$.append$(G)$
— $\boxed{Q(S_t, a_t)} := Avg(Returns(S_t, a_t))$
— $\forall a: \quad \pi(a|S_t) = \begin{cases} 1 - \varepsilon + \varepsilon/|A(S_t)|, & a = a^* \\ \varepsilon/|A(S_t)|, & a \neq a^* \end{cases}$

$\underset{a}{argmax}\, Q(S_t, a)$
$\|$
$a = a^*$

③ Off-policy MC control

(?) $\left[ \begin{array}{l} - \text{порождать эпизоды по } b \quad (behaviour) \\ - \text{обучать не } V_b(s), Q_b(S, c), \text{ a } \quad V_\pi(s), Q_\pi(S, a), \quad \pi \neq b \end{array} \right]$

Importance sampling

$$V_\pi(s) = \mathbb{E}_\pi [G_t | S_t = s]$$

$G_t \sim$ сбор $b$

$\underline{S_t}: \underbrace{A_t, S_{t+1}, A_{t+1}, \ldots, A_{T-1}, S_T}_{Traj}$

$Pr[Traj \,(\pi), S_t] = \pi(A_t | S_t) p(S_{t+1} | S_t, A_t)$
$\qquad \pi(A_{t+1} | S_{t+1}) p(S_{t+2} | \_\_)$
$\qquad \_\_ \_ p(S_T | S_{T-1}, A_{T-1})$

$$= \prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)$$

$\mathbb{E}_{p(x)}[f(x)], \qquad x_n \sim q(x)$
$\|$
$\int f(x) p(x) dx$ $\qquad \int \ldots q(x) dx$
$\|$
$\int \left( f(x) \cdot \dfrac{p(x)}{q(x)} \right) \cdot q(x) dx =$
$= \mathbb{E}_{q(x)} \left[ f \cdot \dfrac{p}{q} \right]$

$p \qquad q$

$\uparrow$ imp. weights

Th-f:
$\forall x \quad p(x) > 0 \Rightarrow q(x) > 0$
$\forall x \quad q(x) = 0 \Rightarrow p(x) = 0$

$\rho_{t:T-1} = \dfrac{Pr[Traj | \pi, S_t]}{Pr[Traj | b, S_t]} = \dfrac{\prod_k \pi(A_k|S_k) p(S_{k+1}|S_k, A_k)}{\prod_k b(A_k|S_k) p(S_{k+1}|S_k, A_k)}$

$\rho_{t:T-1} = \dfrac{\prod_{k=t}^{T-1} \pi(A_k|S_k)}{\prod_{k=t}^{T-1} b(A_k|S_k)} = \prod_{k=t}^{T-1} \dfrac{\pi_k}{b_k}$

Coverage: $\forall s, a \quad \pi(a|s) > 0 \Rightarrow$
$\Rightarrow b(a|s) > 0$

$$V_\pi(s) = \boxed{\mathbb{E}_b \left[ G_t \cdot \rho_{t:T-1} | S_t = s \right]}$$

$\boxed{\rho_T = 1 \quad \rho_t = \rho_{t+1} \left( \dfrac{\pi_t}{b_t} \right)}$

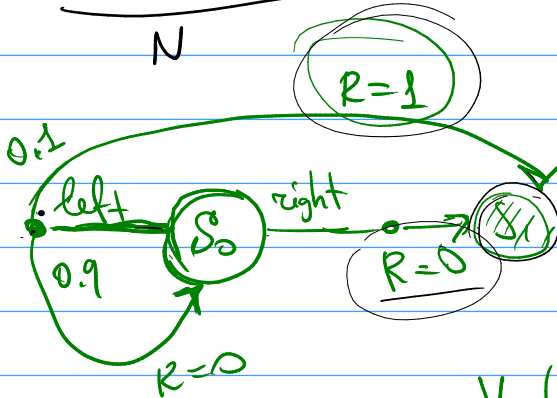$\ldots \qquad t, G, \text{Returns}(S_t).\text{append}\left( G \cdot \mathcal{P}_{t:T-1} \right)$

imp sampling

weighted I.S.

$$V_\pi(s) = \frac{\sum G_i W_i}{N}$$

$$V_\pi(s) = \frac{\sum G_i W_i}{\sum W_i}$$

$R=1$

$\boxed{\gamma = 1}$

$0.1$ left, $0.9$, right, $R=0$, $R=0$

$\pi(\text{left}|s_0) = 1$

$b(\text{left}|s_0) = b(\text{right}|s_0) = \frac{1}{2}$

$V_\pi(s_0) = 1$

$$V_\pi(s_0) = \mathbb{E}_b\left[ G \cdot \frac{\prod_k \pi(a_k|s_k)}{\prod_k b(a_k|s_k)} \right]$$

$Var(x) = \mathbb{E}[x^2] - (\mathbb{E}x)^2$

$$\mathbb{E}_b\left[\left( G \cdot \prod_{k=0}^{T-1} \frac{\pi(a_k|s_k)}{b(a_k|s_k)} \right)^2\right] = \left(\frac{1}{2} \cdot \frac{1}{10}\right) \cdot \left(\frac{1}{1/2}\right)^2 +$$

$$+ \left(\frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \frac{1}{10}\right) \cdot \left(\frac{1}{(1/2)(1/2)}\right)^2 +$$

$$+ \left(\frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \frac{1}{10}\right) \cdot \left(\frac{1}{(1/2)^3}\right)^2 + \ldots$$

$$= \frac{1}{10} \cdot \sum_{k=0}^{\infty} \left(\frac{9}{10}\right)^k \cdot 2^{\gamma(k+1)} \cdot \frac{1}{2^{k+1}} = \frac{1}{10} \cdot \sum_{k=0}^{\infty} \left(\frac{9}{5}\right)^k \to \infty$$

$$\left(\frac{9}{10}\right)^k \cdot 2^k$$

$s_0 \quad \ldots \quad s_T$

$\circledast$ — $G = \gamma G + R_{t+1}$

— $C(s_t, a_t) = C(s_t, a_t) + \rho$

— $Q(s_t, a_t) := Q(s_t, a_t) + \frac{\rho}{C(s_t, a_t)} \cdot (G - Q(s_t, a_t))$

— $\pi(s_t) := \arg\max_a Q(s_t, a)$

— $W := W \cdot \frac{\pi(a_t|s_t)}{b(a_t|s_t)}$

Off policy MC control:

— init, — —

— loop:
   — эпизод из $b$, $G := 0$, $\rho := 1$
   — $\forall t = T-1, \ldots, 0$
      $\hookrightarrow \circledast$

$\pi(s_t) = a_t$?
— если $\pi(s_t) \neq a_t$, to break

$s_{T-2} \ s_{T-1}$

# TD-learning (temporal difference) — bootstrapping

TD(1) $\quad V(S_t) = V(S_t) + \alpha [R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2}) - V(S_t)]$

TD(0): $\quad V(S_t) := V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$

$$V(S_t) \approx G_t = (R_{t+1} + \gamma (R_{t+2} + (\gamma R_{t+3} + \ldots + \gamma^{T-t-2} R_T)))$$
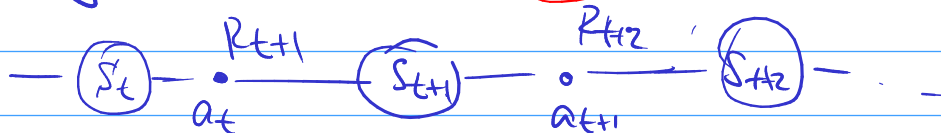
$$" \quad G_{t+1} \approx V(S_{t+1})$$



$V(S_{t+1})$ $R_{t+1}$ $V(S_{T-2})$ $R_{T-1}$ $R_T = 0$ $V(S_{T-1}) = 1$ $V(S_t)$

MC: $V(S_1) = \frac{2}{3}$
$\quad\quad V(S_0) = 0$

TD: $V(S_1) = \frac{2}{3}$
$\quad\quad V(S_0) = \frac{2}{3}$

1) $S_1 \to R = 1$
2) $\bar{S_1} \to R = 0$
3) $\bar{S_1} \to R = 1$
4) $S_1 \to R = 1$
5) $S_0 \to S_1 \to R = 0$
6) $S_1 \to R = 1$

---

## On-policy TD control — Sarsa

$Q_*(s,a)$
$\pi = \pi(Q)$



$S_t \xrightarrow{R_{t+1}}_{a_t} S_{t+1} \xrightarrow{R_{t+2}}_{a_{t+1}} S_{t+2} - \ldots$

— init
— loop. по шагам эпизода:
  — $S, a, s', R$ $\quad\quad (s, a, R, s', a')$
  — выберем $a'$ в $s'$ по $\pi(s') = \varepsilon$-жадн. стр. по $Q(s',a)$
  — $Q(s,a) := Q(s,a) + \alpha [R + \gamma Q(s',a') - Q(s,a)]$
  — $s = s', a = a'$

---

## Off-policy TD control — Q-learning $\boxed{1989}$

$$— Q(s,a) := Q(s,a) + \alpha [R + \gamma \cdot \max_a Q(s',a') - Q(s,a)]$$

$Q_*(s,a)$ — $\pi_*$



start $\to$ finish $R=1$