



$$G_t = \sum_{k=t+1}^{\infty} \gamma^{k-t-1} R_k$$

$$p(s', r | s, a) = \Pr[S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a]$$

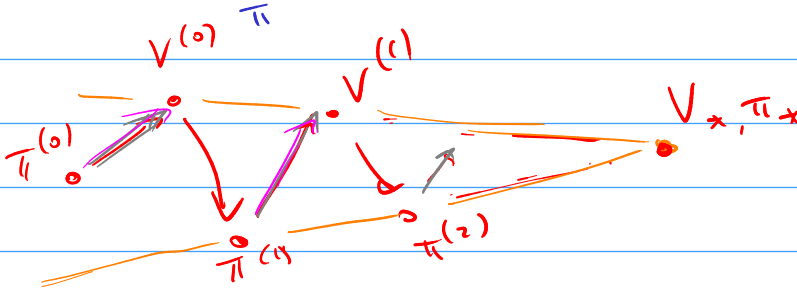
$$\pi: \mathcal{S} \rightarrow \text{Prob}[A(S)], \quad \forall s \quad \pi(a|s) = \Pr[A_t = a | S_t = s]$$

$$V_{\pi}(s) = E_{\pi}[G_t | S_t = s]$$

$$\overline{V_{\pi}(s)} = f(\overline{V_{\pi}(s)})$$

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

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



MC estimation: $V_{\pi}(s) = \text{Avg}[\text{Returns}(s)]$
 $Q_{\pi}(s, a) = \text{Avg}[\text{Ret}(s, a)]$

on-policy

MC control:

$$Q(s_t, a_t) := \text{Avg}[\text{ret}(s, a)]$$

ϵ -soft $\forall a \quad \pi(a|s_t) = \begin{cases} 1 - \epsilon + \epsilon / |A(s_t)|, & a = a^* \\ \epsilon / |A(s_t)|, & \forall a \neq a^* \end{cases}$

Off-policy MC control:

- no policy analyzer no sp. $b(a|s)$ (behavior)
- off-policy $Q_{\pi}(s, a)$, $a \neq Q_b(s, a)$

Importance sampling $E_p[f] = \int p(\bar{x}) f(\bar{x}) d\bar{x} =$

$$\bar{x} \sim q = \int q(\bar{x}) \frac{p(\bar{x})}{q(\bar{x})} f(\bar{x}) d\bar{x} = E_q \left[\frac{p}{q} f \right]$$

importance weights

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] = \mathbb{E}_b \left[G_t \cdot \frac{\Pr[\text{Traj} | \pi]}{\Pr[\text{Traj} | b]} | S_t = s \right] =$$

$S_t, A_t, S_{t+1}, R_{t+1}, A_{t+1}, S_{t+2}, \dots$

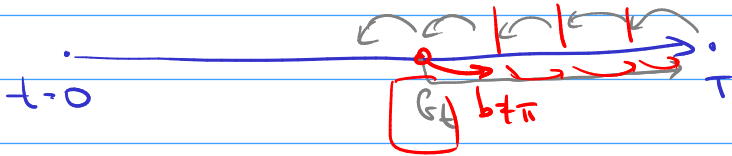
$$\Pr[\text{Traj} | \pi, S_t = s] = \pi(A_t | S_t) p(S_{t+1}, R_{t+1} | S_t, A_t) \pi(A_{t+1} | S_{t+1}) \dots$$

$$= \mathbb{E}_b \left[G_t \cdot \frac{\prod_{k=t}^T \pi(A_k | S_k)}{\prod_{k=t}^T b(A_k | S_k)} | S_t = s \right]$$

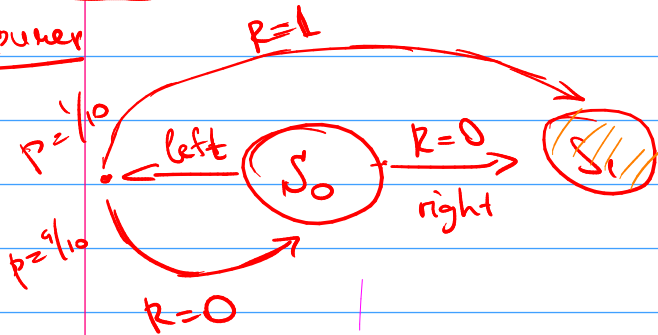
ynobue: even $\pi(a|s) > 0$,
to $b(a|s) > 0$

$$\underline{V_{\pi}(s)} = \text{Avg}(\text{Ret}(s))$$

$G_t \cdot \beta_{t:T}$



Pruner



$$\pi(\text{left} | S_0) = 1$$

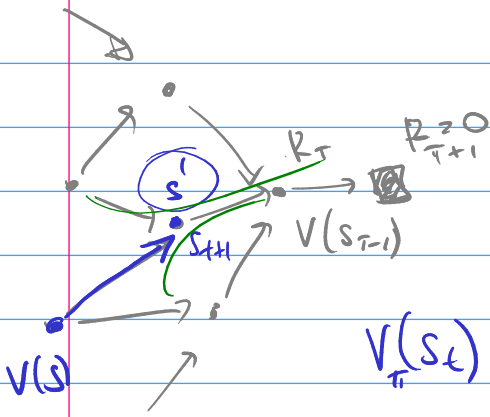
$$b(\text{left} | S_0) = b(\text{right} | S_0) = \frac{1}{2}$$

$$\mathbb{E}_b \left[\left(G \cdot \frac{\prod_k \pi(a_k | S_k)}{\prod_k b(a_k | S_k)} \right)^2 \right] = \left(\frac{1}{2} \cdot \frac{1}{10} \right) \cdot \left(\frac{1}{4} \right)^2 +$$

$$+ \left(\frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \frac{1}{10} \right) \left(\frac{1}{4} \right)^2 + \dots =$$

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

TD-learning (temporal difference)



TD(0)

$$V(s_t) := V(s_t) + \alpha [R_{t+1} + \gamma \cdot V(s_{t+1}) - V(s_t)]$$

$$\text{New} = \text{Old} + \text{Learning Rate} [\text{Target} - \text{Old}]$$

$$V_{\pi}(s_t) = E_{\pi} [G_t | s_t] \approx \text{Arg}(G_t)$$

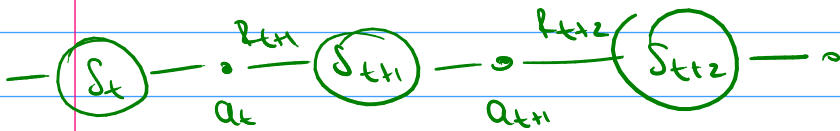
$$V_{\pi}(s_t) = V_{\pi}(s_t) + \frac{1}{n+1} (G_t - V_{\pi}(s_t))$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots =$$

$$= R_{t+1} + \gamma \cdot G_{t+1} \approx R_{t+1} + \gamma \cdot V_{\pi}(s_{t+1})$$

On-policy TD-control

Sarsa



(s, a, r, s', a')

- init

- loop no break $t = 0, \dots, T$:

- (s, a, r, s')

- bootstrap a' by s' no $\pi(s') = \epsilon$ -stoch. exp. no $Q(s, a)$

$$Q(s, a) := Q(s, a) + \alpha (\underbrace{r + \gamma \cdot Q(s', a')}_{\text{target}} - Q(s, a))$$

Off-policy TD-control

Q-learning

1989

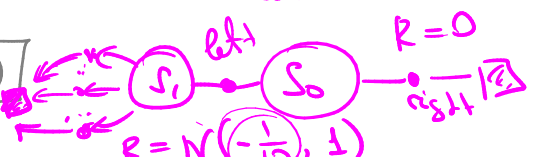
$$Q(s, a) := Q(s, a) + \alpha (r + \gamma \cdot \max_{a'} Q(s', a') - Q(s, a))$$

Q_*

$$\pi_*(s) = \underset{a}{\text{argmax}} Q_*(s, a) \quad \text{also, best}$$

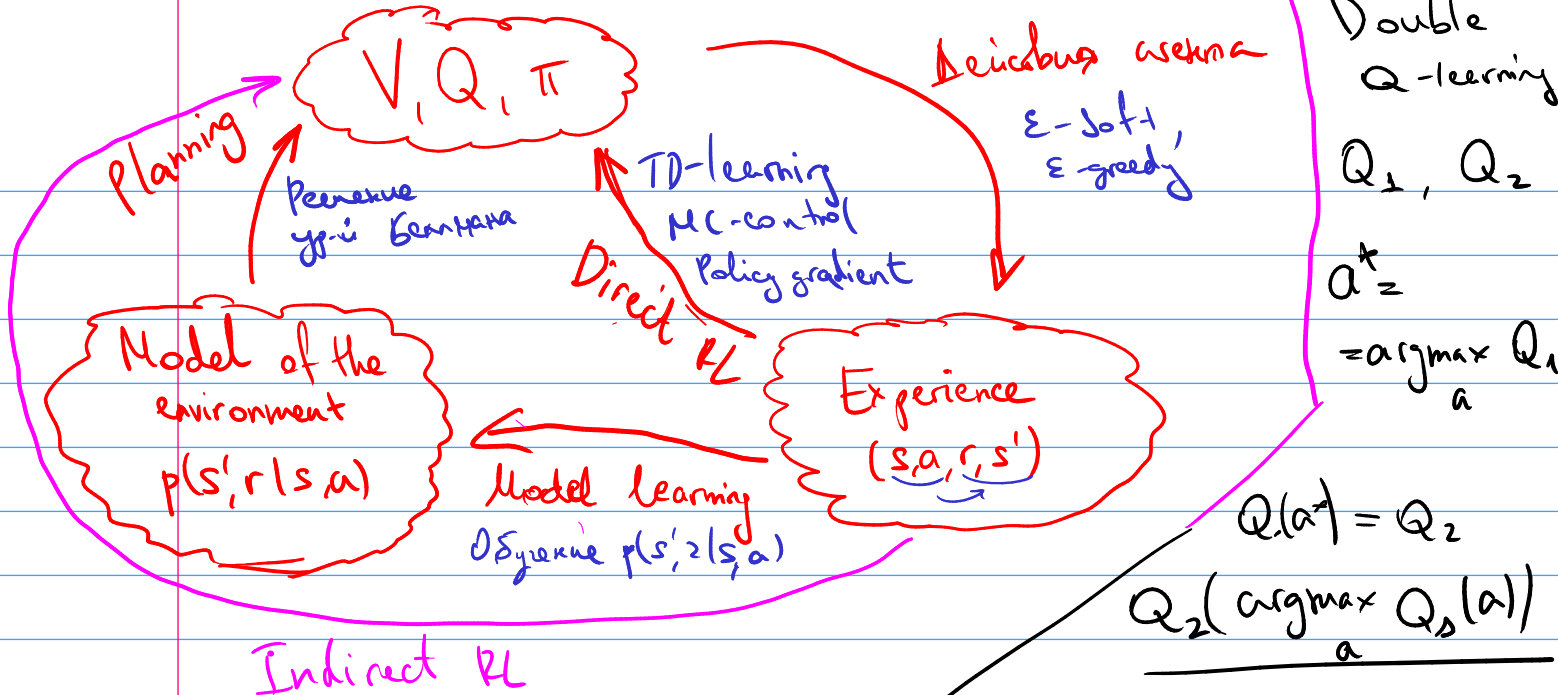


$$N(x | \mu, \sigma^2)$$

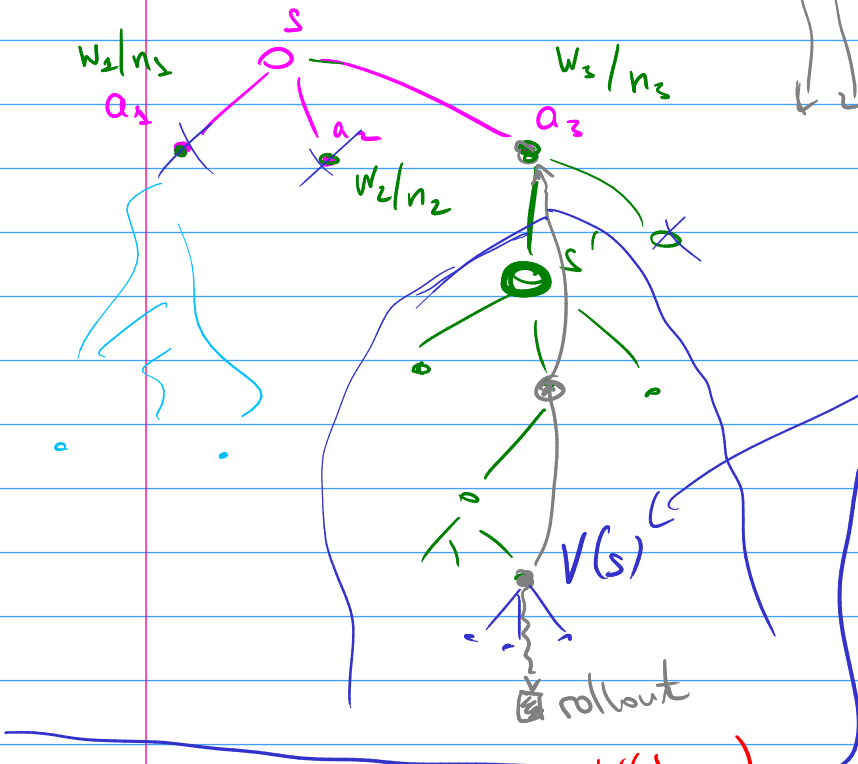


Problems:

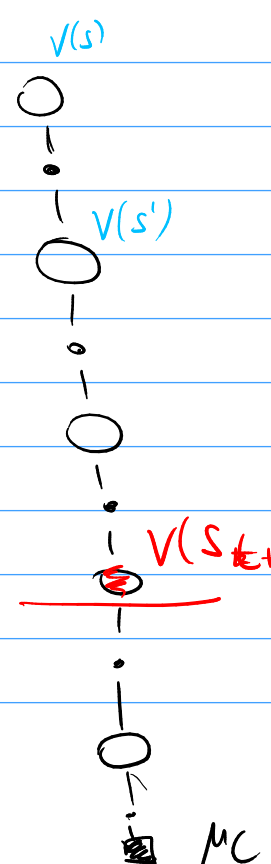
winner's curse
maximization bias



$$V'(Q) \approx \alpha V(s) + (1-\alpha) \cdot \text{Avg}(\text{Rollout})$$



$$\begin{aligned}
 G_t &= R_{t+1} + \gamma(G_{t+1}) \approx V(s_{t+1}) \\
 &= R_{t+2} + \gamma R_{t+2} + \gamma^2(G_{t+2}) = \\
 &= \sum_{i=1}^{k-1} \gamma^i R_{t+i} + \gamma^k(G_{t+k}) \approx V(s_{t+k})
 \end{aligned}$$



Approximate RL

$$Q(s,a) \approx \hat{Q}(s,a;\bar{\theta})$$

$$V(s) \approx \hat{V}(s;\bar{\theta})$$

$$L(\theta) = \sum_s p(s) \left(\underbrace{V(s)}_{s_t} - \hat{V}(s;\bar{\theta}) \right)^2 \xrightarrow{\bar{\theta}} \min$$

SGD

s_t

$$\approx \nabla_{\bar{\theta}} L(\bar{\theta})$$

$$\bar{\theta} := \bar{\theta} + \alpha \left(V(s_t) - \hat{V}(s_t;\bar{\theta}) \right) \cdot \nabla_{\bar{\theta}} \hat{V}(s_t;\bar{\theta})$$

Gradient MC: $V(s_t) \approx G_t$

$$\bar{\theta} := \bar{\theta} + \alpha \left(G_t - \hat{V}(s_t;\bar{\theta}) \right) \cdot \nabla_{\bar{\theta}} \hat{V}(s_t;\bar{\theta})$$

Semi-gradient TD(0):

$$\bar{\theta} := \bar{\theta} + \alpha \left(R_{t+1} + \gamma \hat{V}(s_{t+1};\bar{\theta}) - \hat{V}(s_t;\bar{\theta}) \right) \nabla_{\bar{\theta}} \hat{V}(s_t;\bar{\theta})$$

-/- Q-Learning

$$R_{t+1} + \gamma \max_a \hat{Q}(s_{t+1}, a; \bar{\theta})$$

