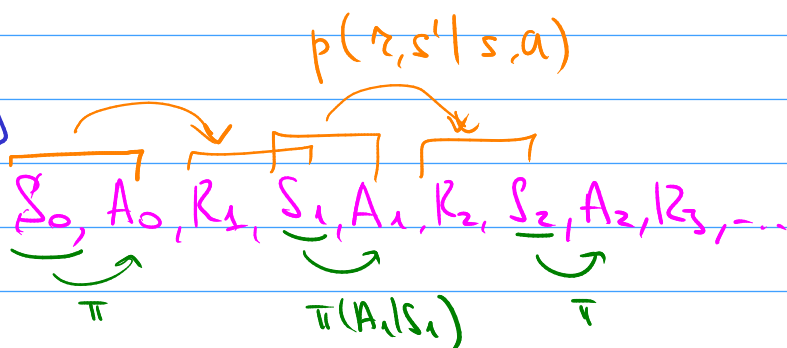
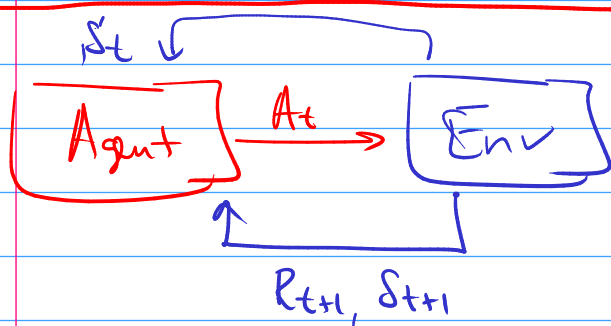


$$\mathbb{E}_{x \sim p(x)} [f(x)] = ?$$

$$x \sim q(x)$$

$$\mathbb{E}_{x \sim q(x)} [g(x)] = \int g(x) q(x) dx$$

$$g(x) = \frac{f(x) p(x)}{q(x)}$$



Return

$$G_t = R_{t+1} + \dots + R_T$$

$$G_t = R_{t+1} + \gamma (R_{t+2} + \gamma^2 R_{t+3} + \dots) = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

State value function

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right]$$

Action value function  $Q$ -function

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi} \left[ R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_t = s, A_t = a \right] = \mathbb{E}_{\pi} \left[ \mathbb{E}_{R_{t+1}, S_{t+1}} \left[ R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | A_{t+1}, R_{t+2}, S_{t+2}, \dots \right] \right]$$

$$= \sum_{s, s'} p(s, s' | s, a) \left( r + \underbrace{\mathbb{E}_{\pi}}_{\substack{A_t, R_{t+1}, S_{t+1}, A_{t+1}, \dots}} \left[ \gamma \cdot \sum \gamma^k R_{t+k+2} \mid S_t = s, A_t = a, \underline{s_{t+1} = s'} \right] \right)$$

$$Q_{\pi}(s, a) = \sum_{s, s'} p(s, s' | s, a) \left( r + \gamma \cdot V_{\pi}(s') \right)$$

$$V_{\pi}(s) = \underbrace{\mathbb{E}_{\pi}}_{\substack{A_t, R_{t+1}, S_{t+1}, A_{t+1}, \dots}} \left[ \sum \gamma^k R_{t+k+1} \mid S_t = s \right] = \mathbb{E}_{a \sim \pi} \left[ \mathbb{E}_{\pi} \left[ G_t \mid S_t = s, A_t = a \right] \right]$$

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

$$\pi' \succeq \pi \Leftrightarrow \forall s \quad V_{\pi'}(s) \geq V_{\pi}(s)$$

$\pi_*$  = ?

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

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

$$\underline{\pi_*(s)} = \operatorname{argmax}_a Q_*(s, a)$$

$$Q_{\pi_*}(s, a) = \mathbb{E}_{\pi_*} [G_t \mid S_t = s, A_t = a] =$$

$$= \sum_{s, s'} p(s, s' | s, a) \left( r + V_{\pi_*}(s') \right)$$

$$V_*(s) = \max_{\pi} V_{\pi}(s) = \max_a \left( \max_{\pi} \sum \pi(a|s) \cdot Q(s, a) \right) =$$

$$= \max_a \max_{\pi} Q(s, a) = \max_a Q_*(s, a)$$

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t \mid S_t = s] = \mathbb{E}_{\pi} [R_{t+1} + \gamma \cdot G_{t+1} \mid S_t = s] =$$

$$= \mathbb{E}_{\pi} [R_{t+1} \mid S_t = s] + \gamma \cdot \mathbb{E}_{\pi} [G_{t+1} \mid S_t = s] =$$

$$= E_{\pi, A_t, R_{t+1}, S_{t+1}} \left[ R_{t+1} + \gamma \underbrace{E_{\pi} [G_{t+1} | S_{t+1}]}_{= V_{\pi}(S_{t+1})} \right]$$

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{z, s'} p(z, s' | s, a) \cdot (z + \gamma V_{\pi}(s')) \quad \text{Bellman equations}$$

Урав. на  $V_{\pi}(s)$ ,  $s \in \mathcal{S}$   $\bar{x} = f(\bar{x})$

$$\begin{aligned} Q_{\pi}(s, a) &= E_{\pi} [G_t | S_t = s, A_t = a] = E_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\ &= E_{\pi, R_{t+1}, S_{t+1}} [R_{t+1} + \gamma E_{\pi} [G_{t+1} | S_{t+1}]] = \\ &= \sum_{z, s'} p(z, s' | s, a) (z + \gamma \cdot \underbrace{V_{\pi}(s')}_{= E_{A_{t+1}} [Q_{\pi}(s, a)]}) \end{aligned}$$

$$Q_{\pi}(s, a) = \sum_{z, s'} p(z, s' | s, a) (z + \gamma \cdot \sum_a' \pi(a' | s') Q_{\pi}(s', a')) \quad \text{Bellman equations}$$

① Переопределить  $\pi$  рекурсивно  $\pi: \mathcal{S} \rightarrow \mathcal{A}$   $|\mathcal{S}|^{|\mathcal{A}|}$

② Мы не знаем  $p(z, s' | s, a)$

③ Все плохо уравнению неизвестно - как  $|\mathcal{S}|$  или  $|\mathcal{S}| \cdot |\mathcal{A}|$

$$\begin{aligned} V_*(s) &= \max_{\pi} V_{\pi}(s) = \max_{\pi} \sum_a \pi(a|s) \sum_{z, s'} p(z, s' | s, a) (z + \gamma V_{\pi}(s')) = \\ &= \max_a \sum_{z, s'} ( \dots, \gamma \max_{\pi} V_{\pi}(s') ) \end{aligned}$$

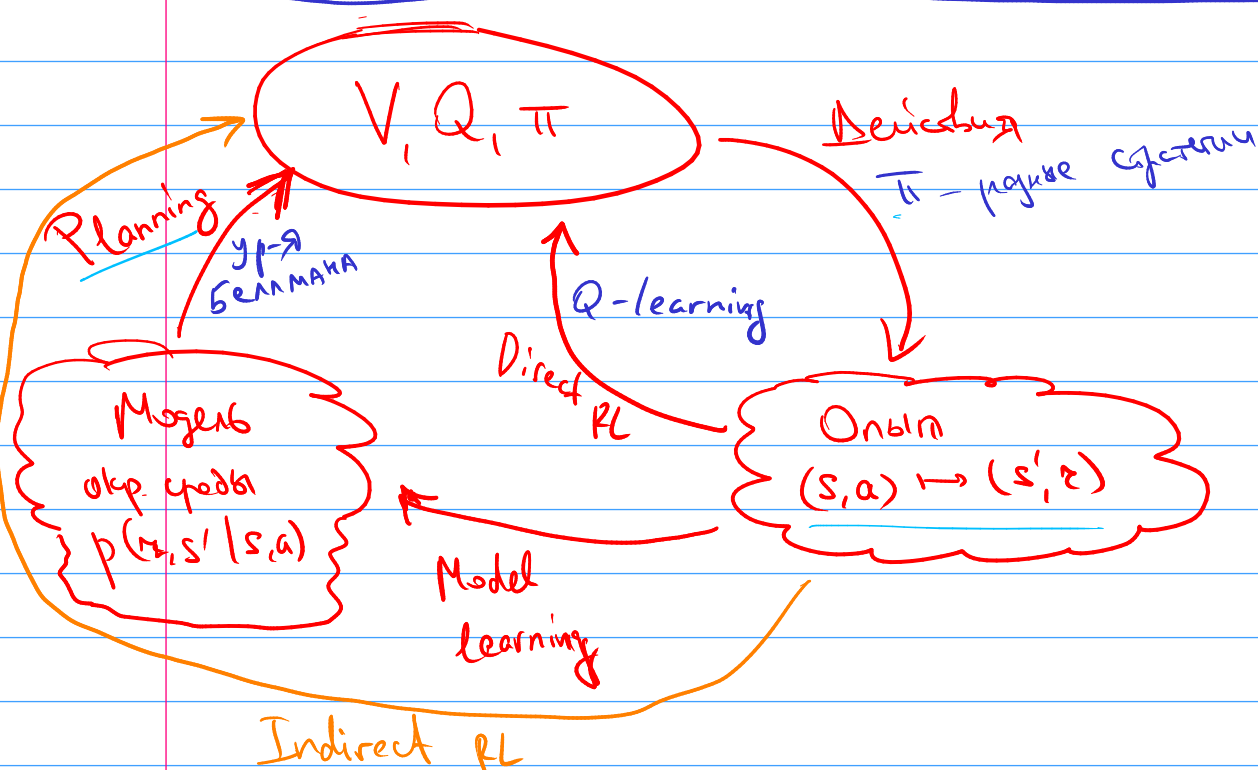
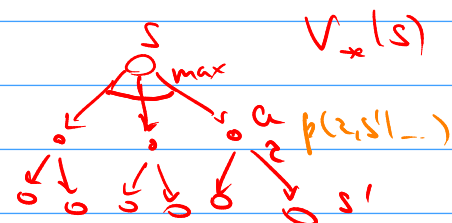
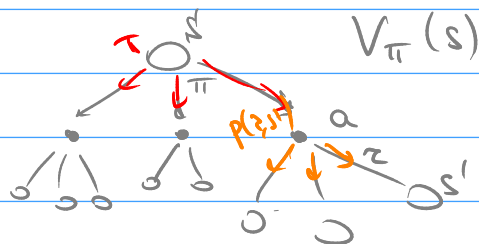
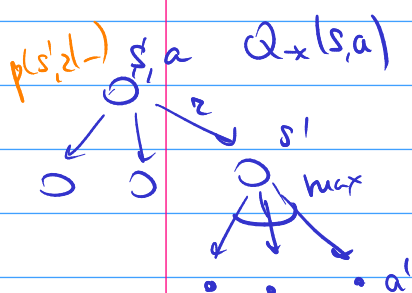
$$V_*(s) = \max_a \sum_{z, s'} p(z, s' | s, a) (z + \gamma V_*(s')) \quad \text{Bellman equations}$$

$$Q_{\pi}(s, a) = \max_{\pi} Q_{\pi}(s, a) = \left( \max_{\pi} \sum_{z, s'} p(z, s' | s, a) \left( z + \gamma \sum_{a'} \pi(a' | s') Q_{\pi}(s', a') \right) \right)$$

$$Q_{\pi}(s, a) = \sum_{z, s'} p(z, s' | s, a) \left( z + \gamma \max_{a'} Q_{\pi}(s', a') \right)$$

$$V^{(0)}(s) \quad V^{(k+1)}(s) := \max_{\pi} ( \dots V^{(k)}(s') )$$

backup diagram



Policy improvement

$\pi \rightsquigarrow \pi'!$

$\pi \rightsquigarrow V_{\pi}(s), Q_{\pi}(s, a)$

Policy impr. thm. Esau der byx geograph.  $\pi, \pi'$

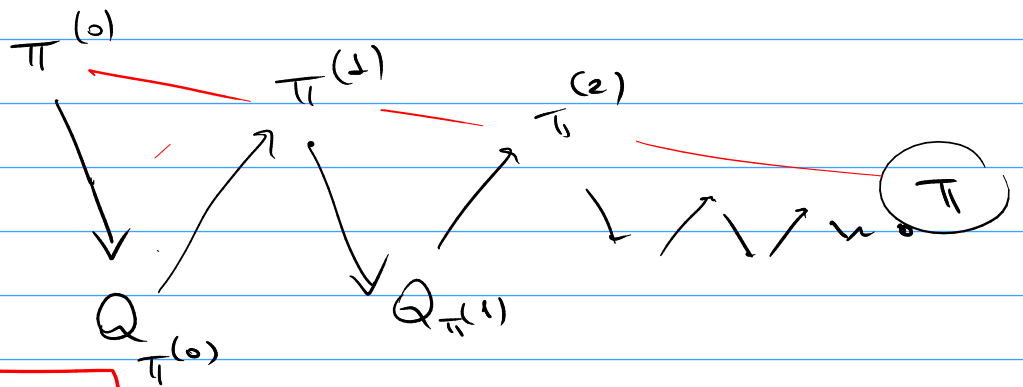
$$\forall s \in \mathcal{S} \quad Q_{\pi'}(s, \pi'(s)) \geq V_{\pi}(s) = Q_{\pi}(s, \pi(s))$$

$$\text{TO } \forall s \in \mathcal{S} \quad V_{\pi'}(s) \geq V_{\pi}(s)$$

$$\begin{aligned} V_{\pi}(s) &\leq Q_{\pi}(s, \pi'(s)) = E_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s, A_t = \pi'(s)] \\ &= E_{\pi'} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s] \leq \\ &\leq E_{\pi'} [R_{t+1} + \gamma Q_{\pi}(s_{t+1}, \pi'(s_{t+1})) | s_t = s] = \\ &\leq E_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{\pi}(s_{t+2}) | s_t = s] \leq \\ &\leq \dots \leq E_{\pi'} [R_{t+1} + \gamma R_{t+2} + \dots | s_t = s] = V_{\pi'}(s) \end{aligned}$$

$$\pi \rightarrow V_{\pi}, Q_{\pi} \quad \boxed{\pi'(s) := \arg \max_a Q_{\pi}(s, a)}$$

Policy Iteration



$$\pi(s) = \arg \max_a Q_{\pi}(s, a)$$

