

30.11.18 Reinforcement Learning

ML $\{x_i, y_i\}_{i=1}^N, f(x, \theta)$

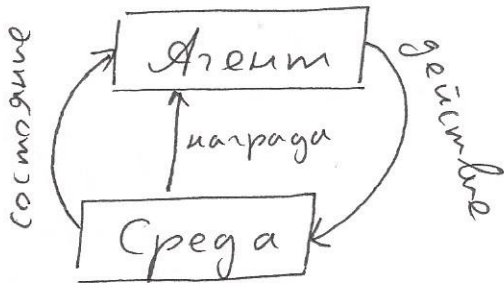
обучаем параметрами θ

$$F(\theta) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i, \theta)) + \lambda R(\theta) \rightarrow \min_{\theta}$$

отложенное вознаграждение?

взаимодействие со средой?

RL



$s \in S$ - сост. среды

$a \in A$ - действие агента

$\pi(a|s)$ - политика агента

$p(s'|s, a)$ - вер-ти переходов

$r(s, a)$ - награда

MDP, Markov Decision Process

$$\{s_0, a_0, r_0, s_1, a_1, r_1, \dots\}$$

суммарный дисконтированный возврат

$$R_{\infty} = r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) + \dots, \gamma \in (0, 1)$$

$$R_t = r(s_t, a_t) + \gamma R_{t+1} \rightarrow \max_{\pi}$$

Value Iteration

$$V^{\pi}(s) = \mathbb{E}[R_{\infty} | s_t = s]$$

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

$$\pi^*(s) = \arg \max_{\pi} V^{\pi}(s)$$

$$R_0 = r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) + \dots = r(s_0, a_0) + \gamma$$

$$V^*(s) = \max_{\pi} \mathbb{E}[R_0 | s_0 = s] = \max_{a_0, a_1, a_2, \dots} \mathbb{E}[R_0 | s_0 = s] =$$

марковский процесс, отсюда равенство

$$= \max_{a_0, a_1, a_2, \dots} [r(s_0, a_0) + \gamma \mathbb{E}_{p(s'|s_0, a_0)} \mathbb{E}[R_1 | s_1 = s']] =$$

$$= \max_{a_0} [r(s_0, a_0) + \gamma \mathbb{E}_{p(s'|s_0, a_0)} \underbrace{\max_{a_1, a_2, \dots} \mathbb{E}[R_1 | s_1 = s']}_{V^*(s')}]$$

$$V^*(s) = \max_a (r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} V^*(s'))$$

оптимальное уравнение Беллмана, $x = f(x)$

fixed point iteration, оптимальное значение $x_{k+1} = f(x_k)$

$$V_{\text{new}}(s) = \max_a (r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} V_{\text{old}}(s')) \quad \forall s$$

$$\pi^*(s) = \arg \max_a (r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} V^*(s'))$$

Пример

мерн. соим.

$$r:$$

0	0	0	1
0	1	0	-1
0	0	0	0

$\rightarrow V_1:$

0	0	0	1
0	1	0	-1
0	0	0	0

\rightarrow

$$V_2:$$

0	0	0	1
0	1	0	-1
0	0	0	0

$\rightarrow V_5:$

γ^3	γ^2	γ	1
γ^4	γ^3	γ^2	-1
γ^5	γ^4	γ^3	γ^4

Q-learning

$$Q^\pi(s, a) = \mathbb{E}[R_t | s_t = s, a_t = a]$$

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

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

$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = \max_{a_1, a_2, a_3, \dots} \mathbb{E}[R_0 | s_0 = s, a_0 = a] =$$

$$= \max_{a_1, a_2, a_3, \dots} [\gamma(s_0, a_0) + \delta \mathbb{E}_{p(s'|s_0, a_0)} \mathbb{E}[R_1 | s_1 = s', a_1 = a_1]] =$$

$$= \gamma(s_0, a_0) + \delta \mathbb{E}_{p(s'|s_0, a_0)} \underbrace{\max_{a_2} \max_{a_2, a_3, \dots} \mathbb{E}[R_1 | s_1 = s', a_1 = a_1]}_{Q^*(s', a_1)}$$

$$Q^*(s, a) = \gamma(s, a) + \delta \mathbb{E}_{p(s'|s, a)} \max_{a'} Q^*(s', a')$$

$$\forall s, a$$

$$F(Q^*) = \frac{1}{|S| |A|} \sum_{s, a} (Q^*(s, a) - \gamma(s, a) - \delta \mathbb{E}_{p(s'|s, a)} \max_{a'} Q^*(s', a'))^2 \rightarrow \min_{Q^*}$$

1) Обновление (s, a, r, s')

$$2) Q_{\text{new}}^*(s, a) = Q_{\text{old}}^*(s, a) - \alpha_{s, a} \cdot 2 (Q_{\text{old}}^*(s, a) - \gamma(s, a) - \delta \max_{a'} Q_{\text{old}}^*(s', a'))$$

Exploration - Exploitation dilemma

Локально-оптимальные решения?

ϵ - малая переменная:

$$\pi(s) = \begin{cases} \operatorname{argmax}_a Q(s, a), & \text{с вер-ю } 1-\epsilon \\ \sim R, & \text{с вер-ю } \epsilon \end{cases}$$

Softmax - функция

$$\pi(a|s) \sim \text{Softmax}(Q(s,a)/T)$$

сэмплинг, $T \rightarrow \infty$ \mathcal{U}
 $T \rightarrow 0$ argmax

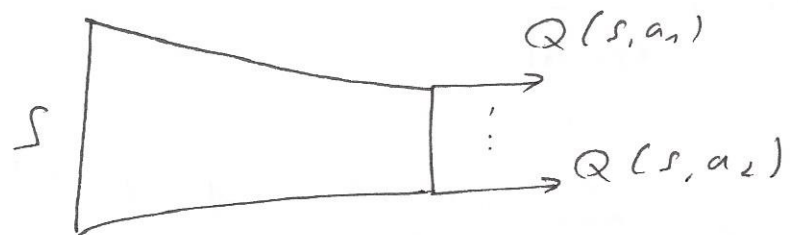
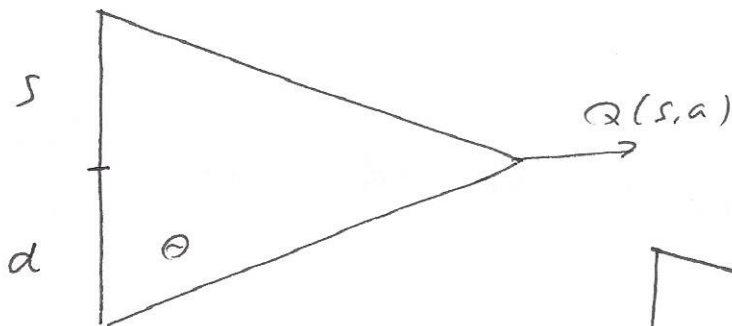
$$Q^*(s,a) \approx Q^*(s,a|\theta)$$

параметризация функции

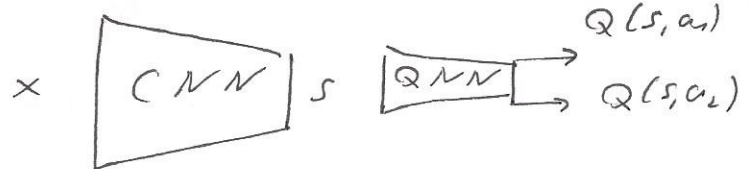
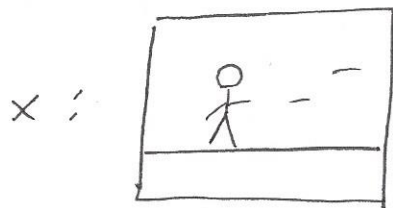
$$F(\theta) = \frac{1}{|S||A|} \sum_{s,a} (Q(s,a|\theta) - r(s,a) -$$

$$- \gamma \mathbb{E}_{p(s'|s,a)} \max_{a'} Q(s',a'|\theta))^2 \rightarrow \min_{\theta}$$

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha \cdot 2 \cdot (Q(s,a|\theta_{\text{old}}) - \dots) \nabla_{\theta} Q(s,a|\theta_{\text{old}})$$



DQN



Инициализация θ, M

...

Иници. θ, M

Для эпизодов $1, \dots, p$

инициализация x_1

$$s_1 = \text{CNN}(x_1)$$

для $t = 1 \dots T$:

$$a_t = \begin{cases} \arg \max Q(s_t, a; \theta) & \text{с вер-ю } 1-\epsilon \\ \sim R, \epsilon & \end{cases}$$

$$(s_t, a_t, \hat{c}_t, x_{t+1} \rightarrow s_{t+1}) \rightarrow M \quad \left. \begin{array}{l} \text{experience} \\ \text{replay, mini-batch} \end{array} \right\}$$

$$(s_j, a_j, r_j, s_{j+1}) \leftarrow M$$

$$y_j = \begin{cases} r_j + \gamma \max_{a'} Q(s_j, a'; \theta), & s_j - \text{не мерн.} \\ r_j, & \text{если } s_j - \text{мерн.} \end{cases}$$

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha_2 (Q(s_j, a_j; \theta) - y_j) \nabla_{\theta} Q(s_j, a_j; \theta)$$