

# Q-learning

(\*и разные трюки)

Подготовил: Грузицкий Андрей  
БПМИ211

# Вспомним

\*дальше будет меньше формул, не пугайтесь

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}. \end{aligned}$$

$$v_{\pi}(s) \triangleq \mathbb{E}_{\pi} [G_t \mid S_t = s]$$

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t \mid S_t = s, A_t = a] = \sum_{r, s'} p(r, s' \mid s, a) [r + \gamma v_{\pi}(s')]$$

**$Q^*(s, a)$**  - сейчас совершили действие  $a$ , а потом выбираем действия с наибольшим  $q$

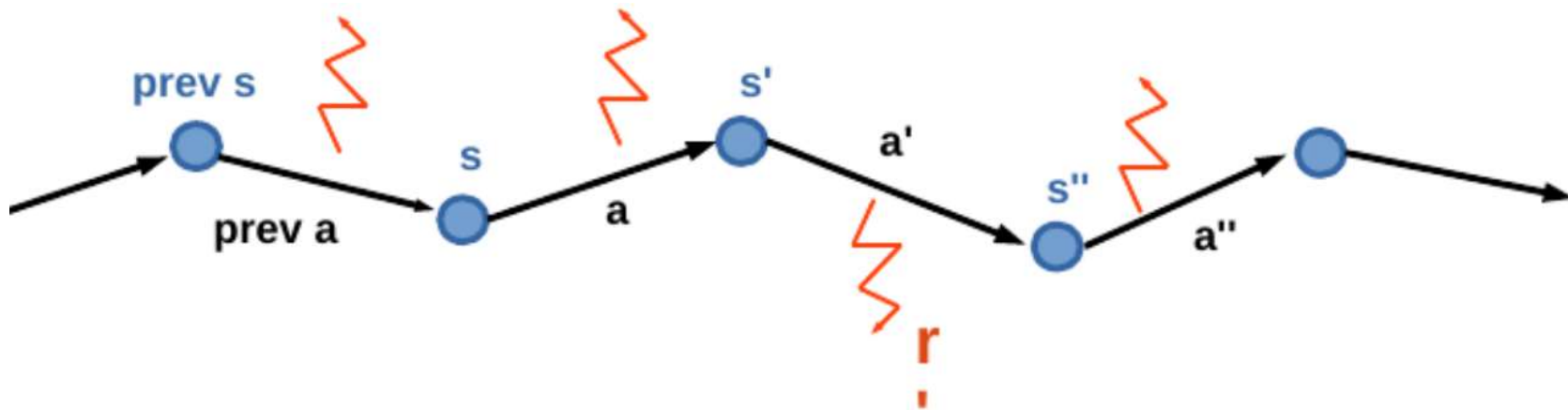
# Model-free reinforcement learning

Теперь у нас нет  $P(s', r | s, a)$  !!!

Model based: можем планировать наперёд

Model free: только пробовать

## Учимся на траекториях

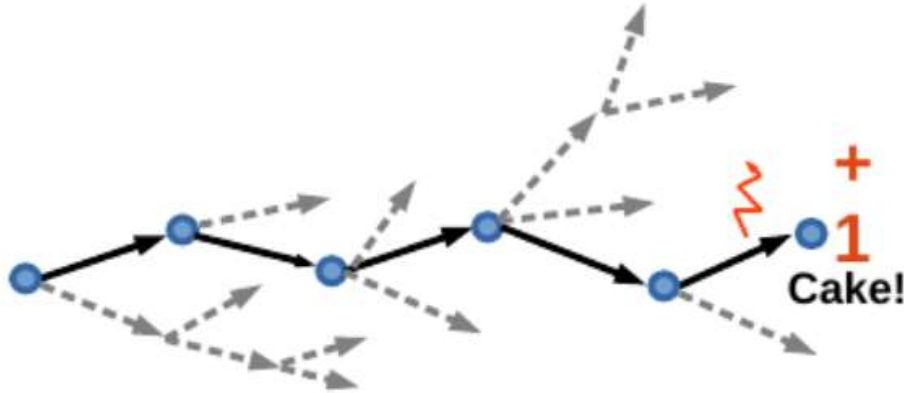


Теперь  $V(s)$  - бесполезно, т.к. не знаем  $P(s'|s, a)$

Учим  $Q(s, a)$  - берём  $\arg\max$  по  $a \Rightarrow$  победа

## Monte-carlo

- Averages  $Q$  over sampled paths



Требует очень много сессий!

## Temporal Difference

- Uses recurrent formula for  $Q$



Меньше сэмплов. Обучаемся даже не доиграв одной игры.

# Temporal difference

Можем обновлять значения Q итеративно:

$$Q(s_t, a_t) \leftarrow E_{r_t, s_{t+1}} r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')$$

# Temporal difference

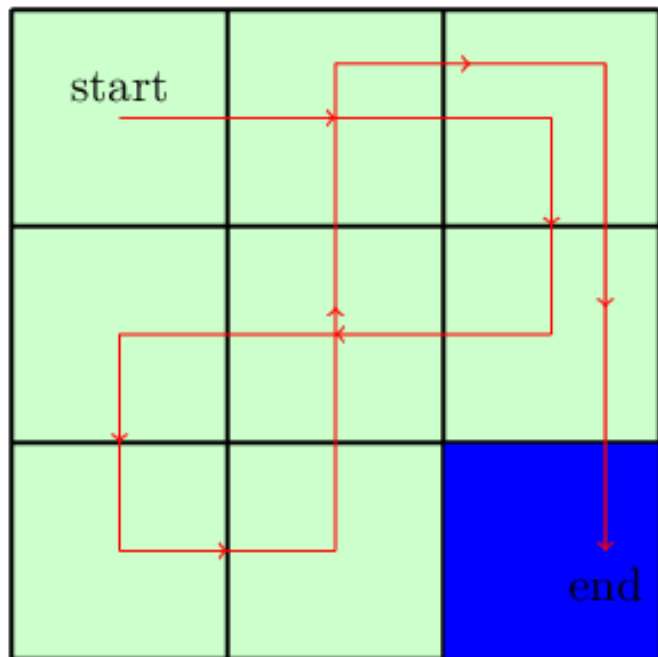
- Приближаем сэмплами:

$$E_{r_t, s_{t+1}} r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a') \approx \frac{1}{N} \sum_i r_i + \gamma \cdot \max_{a'} Q(s_i^{\text{next}}, a')$$

- Пользуемся скользящим средним. На одном шаге используем только один объект!!!

$$Q(s_t, a_t) \leftarrow \alpha \cdot (r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')) + (1 - \alpha) Q(s_t, a_t)$$

# Q-learning



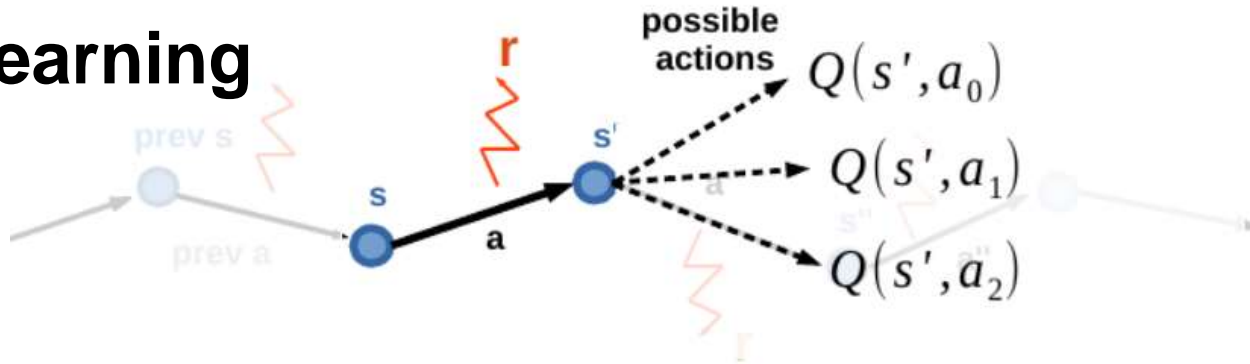
-0.1 →	0 ↑	-0.1 →	0 ↑	0 →	0 ↑
0 ←	0 ↓	0 ←	0 ↓	0 ←	-0.1 ↓
0 →	0 ↑	0 →	-0.1 ↑	0 →	0 ↑
0 ←	-0.1 ↓	-0.1 ←	0 ↓	-0.1 ←	0 ↓
-0.1 →	0 ↑	0 →	-0.1 ↑		
0 ←	0 ↓	0 ←	0 ↓		



# Q-learning

-3.0 →	-3.6 ↑	-2.0 →	-2.7 ↑	-1.2 →	-1.2 ↑
-3.8 ←	-3.0 ↓	-3.0 ←	-2.0 ↓	-1.6 ←	-1.0 ↓
-2.0 →	-2.9 ↑	-1.0 →	-1.7 ↑	-0.6 →	-0.8 ↑
-2.7 ←	-2.0 ↓	-1.8 ←	-1.0 ↓	-1.4 ←	0.0 ↓
-1.0 →	-1.4 ↑	0.0 →	-0.4 ↑		
-1.4 ←	-1.5 ↓	-1.0 ←	-0.7 ↓		

# Q-learning



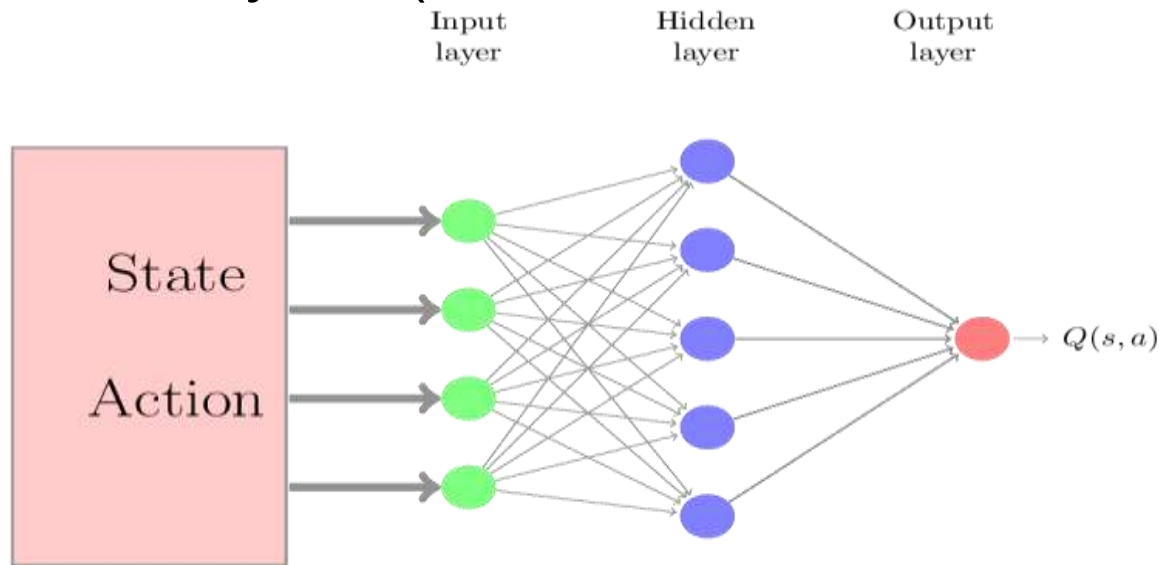
Инициализируем  $Q(s, a)$  нулями

- Loop:

- Из среды получаем значения  $\langle s, a, r, s' \rangle$
- Считаем  $\hat{Q}(s, a) = r(s, a) + \gamma \max_i Q(s', a_i)$
- Обновляем  $Q(s, a) \leftarrow \alpha \cdot \hat{Q}(s, a) + (1 - \alpha) Q(s, a)$

Непрерывный случай (таблица со значениями Q не спасёт)

**DQN:**



$$w = \operatorname{argmin} \left( f^w(s, a) - \left( r + \max_b (f^w(s', b)) \right) \right)^2$$

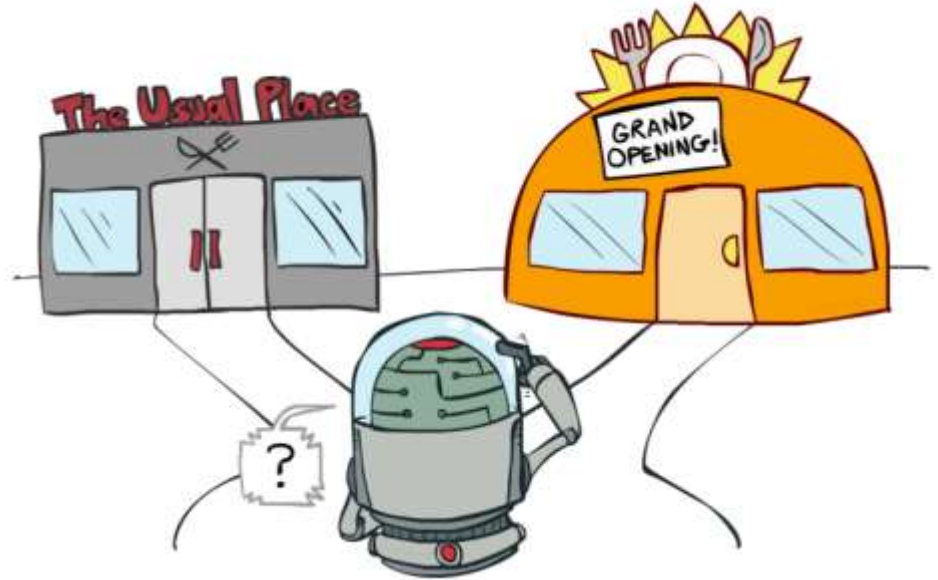
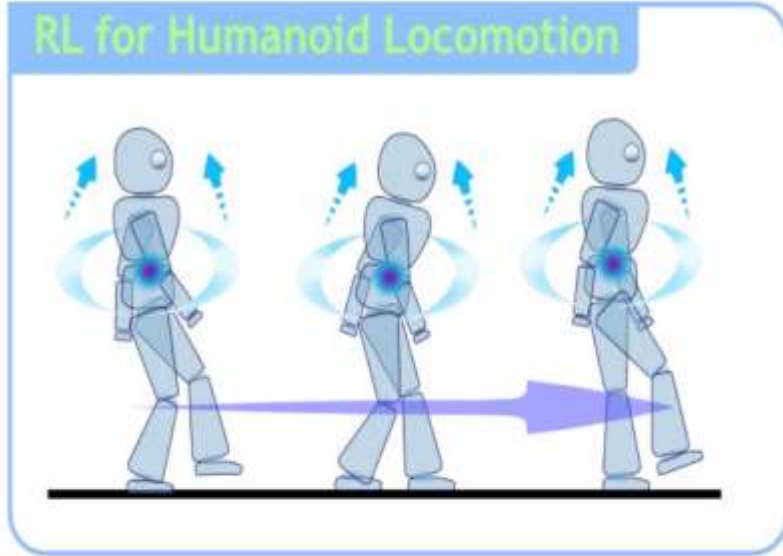
# Непрерывный случай

Можем использовать MSE:

$$w = \operatorname{argmin} \left( f^w(s, a) - \left( r + \max_b (f^w(s', b)) \right) \right)^2$$

$$L = (Q(s_t, a_t) - [r + \gamma \cdot \max_{a'} Q(s_{t+1}, a')])^2$$

# Exploration Vs Exploitation:



# Exploration Vs Exploitation

- $\epsilon$ -greedy

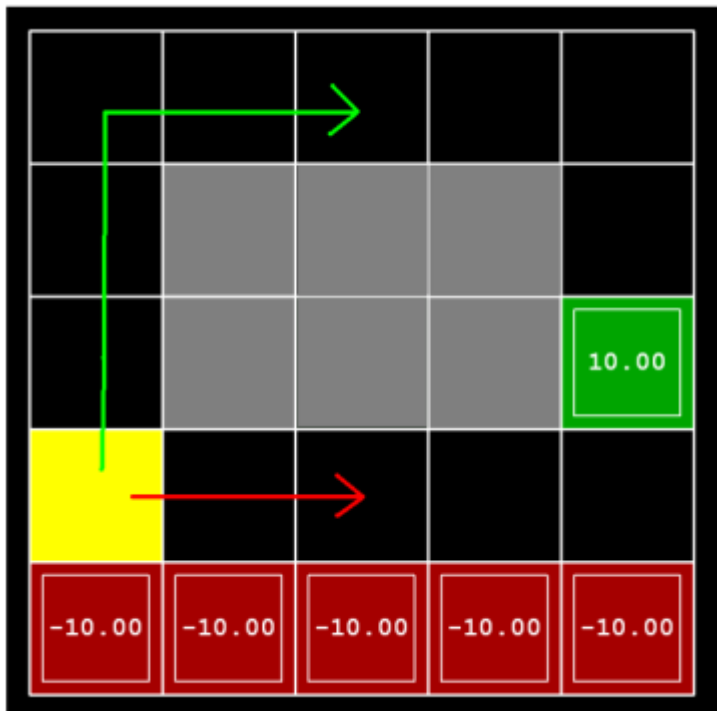
С вероятностью  $\epsilon$  выбираем случайное действие, иначе берём оптимальное

- Softmax

Выбираем действие пропорционально значению softmax Q-value нормированного на  $\tau$

$$\pi(a|s) = \text{softmax}\left(\frac{Q(s, a)}{\tau}\right)$$

## Cliff problem:



- Q-learning

$$\gamma = 0.99 \quad \epsilon = 0.1$$

- не хотим падать в лаву!!!

## Cliff problem:

$$Q(s_t, a_t) \leftarrow \alpha \cdot \hat{Q}(s_t, a_t) + (1 - \alpha) Q(s_t, a_t)$$

Q-learning

“better  $Q(s,a)$ ”



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

SARSA

$$\hat{Q}(s, a) = r(s, a) + \gamma \cdot E_{a' \sim \pi(a'|s')} Q(s', a')$$



# On-policy Vs Off-policy



# Experience replay

**Idea:** store several past interactions

$\langle s, a, r, s' \rangle$

Train on random subsamples

**Training curriculum:**

- play 1 step and record it
- pick N random transitions to train

**Profit:** you don't need to re-visit same  $(s, a)$  many times to learn it.

**Only works with  
off-policy algorithms!**

**Old  $(s, a, r, s')$  are  
from older/weaker  
version of policy!**

