



Продолжение Q-learning. DQN



- ✓ Double Q-Learning
- ✓ Нейронные сети
- ✓ Deep Q Network (DQN)
- ✓ Continuous actions



Q-обучение

Q – обучение (Q – *learning*):

Есть какая – то $Q(s, a)$. Рассмотрим пример $s_t, a_t \rightarrow r_t \rightarrow s_{t+1}$

Можем посчитать $Q(s_t, a_t)$ и $r_t + \gamma \max_{a'} Q(s_{t+1}, a')$

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

Для Q^* стратегия тривиальна:

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

Что делать для произвольной Q в процессе обучения?

$\pi(a | s)$ = с некоторой вероятностью не оптимальное действие(в смысле Q функции)

Exploration-exploitation dilemma

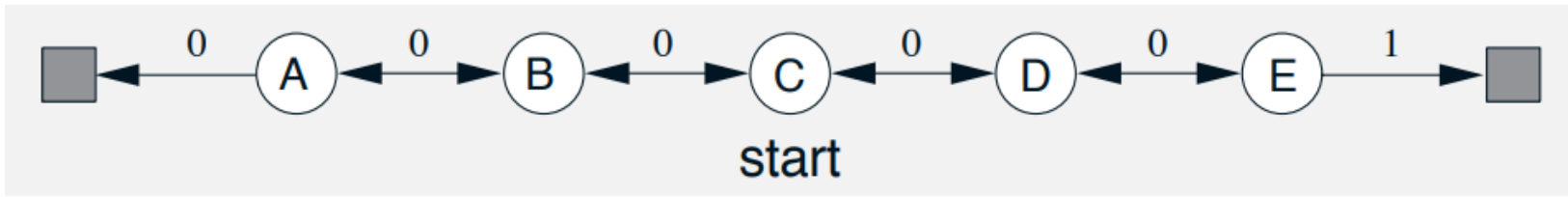


ε – жадная стратегия:

$$\pi(a | s) = \begin{cases} \underset{a}{\operatorname{argmax}} Q(s, a) & \text{с вероятностью } 1 - \varepsilon \\ \forall a \neq \underset{a}{\operatorname{argmax}} Q(s, a) & \text{с вероятностью } \frac{\varepsilon}{|A| - 1} \end{cases}$$



Пример вычисления функции ценности



$$V^*(s) = \max_{\pi} \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s]$$

$$V^*(s) = \max_a \mathbb{E}_{s' \sim P} [r(s, a) + \gamma V^*(s')]$$

$$V^*(C) = \frac{1}{2}$$

$$V^*(A) = 0.5 * 0 + 0.5 * (0 + V^*(B))$$

$$V^*(B) = 0.5 * (0 + V^*(C)) + 0.5 * (0 + V^*(A))$$

$$V^*(C) = \frac{1}{2}$$

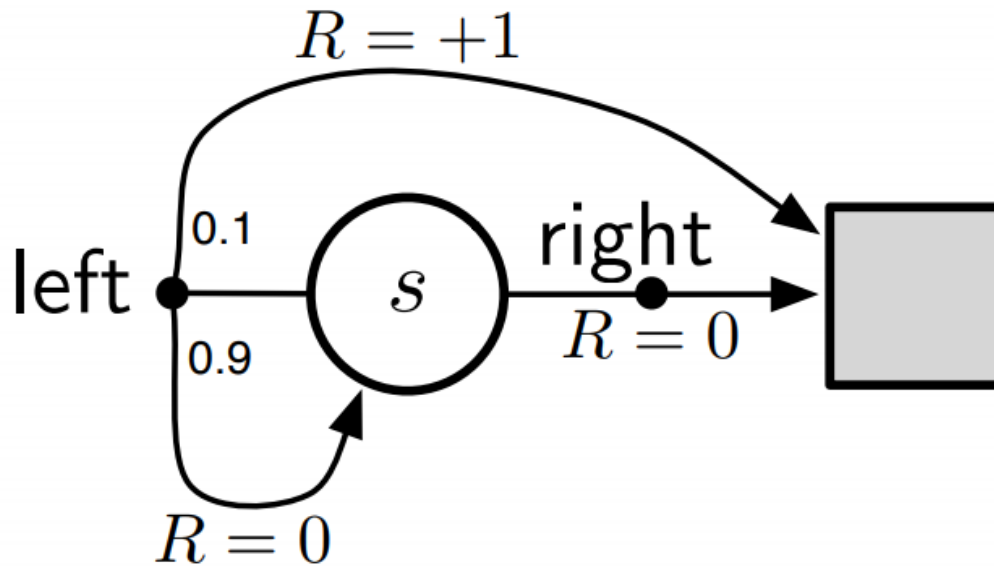
$$2V^*(A) = V^*(B)$$

$$4V^*(B) = 2V^*(A) + 1$$

$$V^*(B) = \frac{1}{3}$$



Ещё пример



Политики:

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

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

Найти $V^\pi(s)$ и $V^b(s)$:

- С помощью уравнения Беллмана
- Через определение



Double Q-Learning

$$\mathbb{E} \max(X, Y) \geq \max(\mathbb{E}X, \mathbb{E}Y)$$

$$\mathbb{E} \max(X, Y) = \int \max(x, y) f(x)f(y) dx dy \geq \int x f(x) f(y) dx dy = \mathbb{E}X$$

$$\mathbb{E} \max(X, Y) = \int \max(x, y) f(x)f(y) dx dy \geq \int y f(x) f(y) dx dy = \mathbb{E}Y$$

Пусть $X, Y \sim U[0,1], i. i. d$

$$\max(\mathbb{E}X, \mathbb{E}Y) = 0.5$$

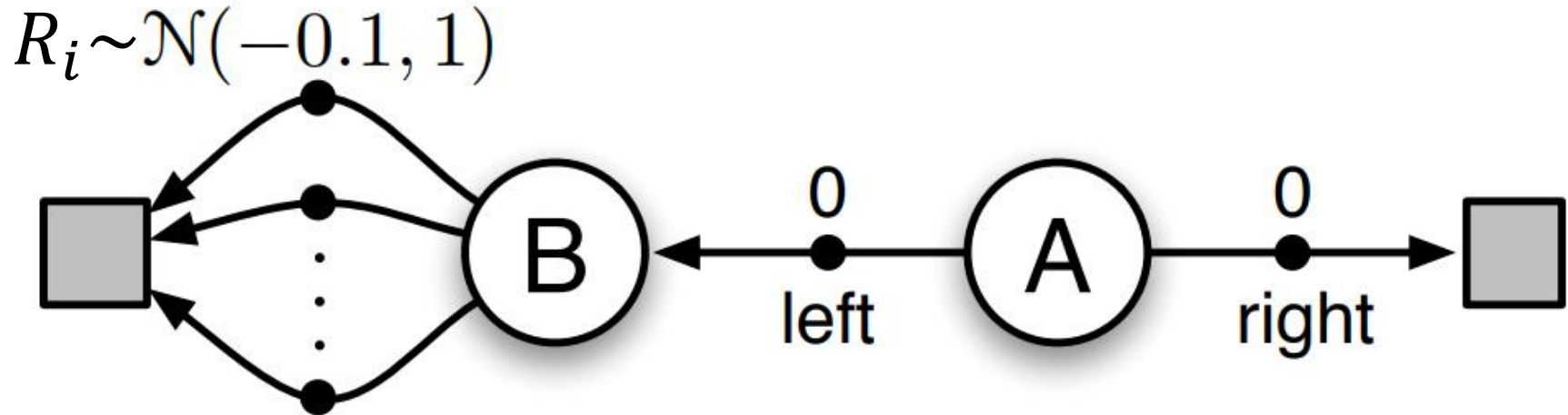
$$\mathbb{E} \max(X, Y) = \int \max(x, y) dx dy = 2 \int x dx dy = 2 \int_0^1 x^2 dx = 2/3$$

\square

\triangle



Double Q-Learning



*Посчитаем Q^**

$$Q^*(A, \text{right}) = 0$$

$$Q^*(A, \text{left}) = \max Q^*(B, a) = \max \mathbb{E}(R_i) = -0.1$$

**Но в начале мы не
располагаем большой
выборкой, поэтому оценка
среднего получается
с большой дисперсией**

Оценка $Q(A, \text{left})$ после первого шага:

$$\max Q(B, a) = \max(R_i)$$

$$\mathbb{E} \max(Q(B, a)) = \mathbb{E} \max(R_i) \geq \max \mathbb{E}(R_i) = -0.1$$

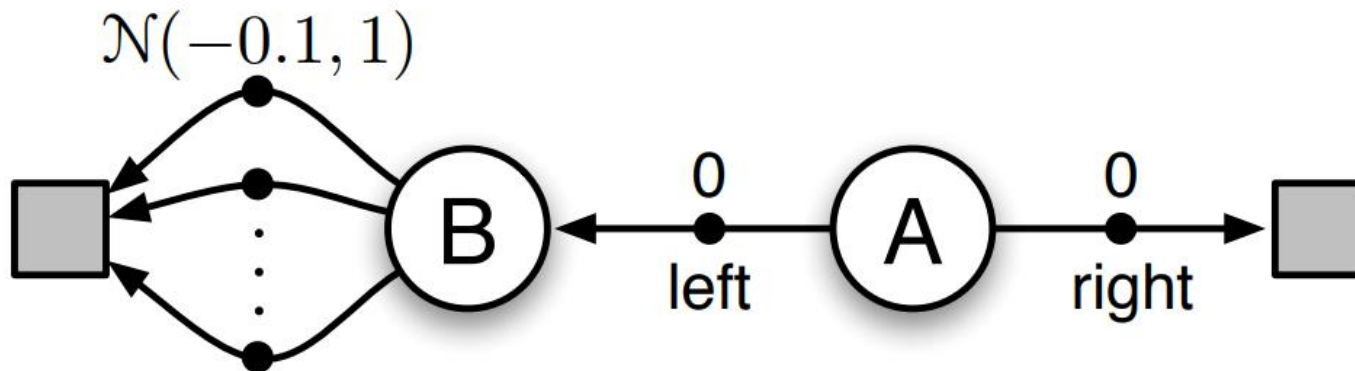


Double Q-Learning

Ну и что с того, что переоценивает?

Нам же важно правильно найти лучшее действие в данном состоянии, а с ранжированием (упорядочиванием) всё в порядке

Посмотрим еще раз



$$\mathbb{E}Q(A, left) = \mathbb{E} \max(Q(B, a) = \mathbb{E} \max(R_i) \geq -0.1$$

$$Q(A, right) = 0 \text{ не сильно больше, чем } -0.1$$

\Rightarrow возможно такое, что в начале $Q(A, left) \geq Q(A, right)$

Но это неверно!



Double Q-Learning

Заметим, что $\max Q(S', a) = Q(S', \operatorname{argmax} Q(S', a))$

Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, such that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

Initialize S

Loop for each step of episode:

Choose A from S using the policy ε -greedy in $Q_1 + Q_2$

Take action A , observe R, S'

With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left(R + \gamma Q_2(S', \operatorname{argmax}_a Q_1(S', a)) - Q_1(S, A) \right)$$

else: Шумы некоррелированы, поэтому переоценивания не происходит

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left(R + \gamma Q_1(S', \operatorname{argmax}_a Q_2(S', a)) - Q_2(S, A) \right)$$

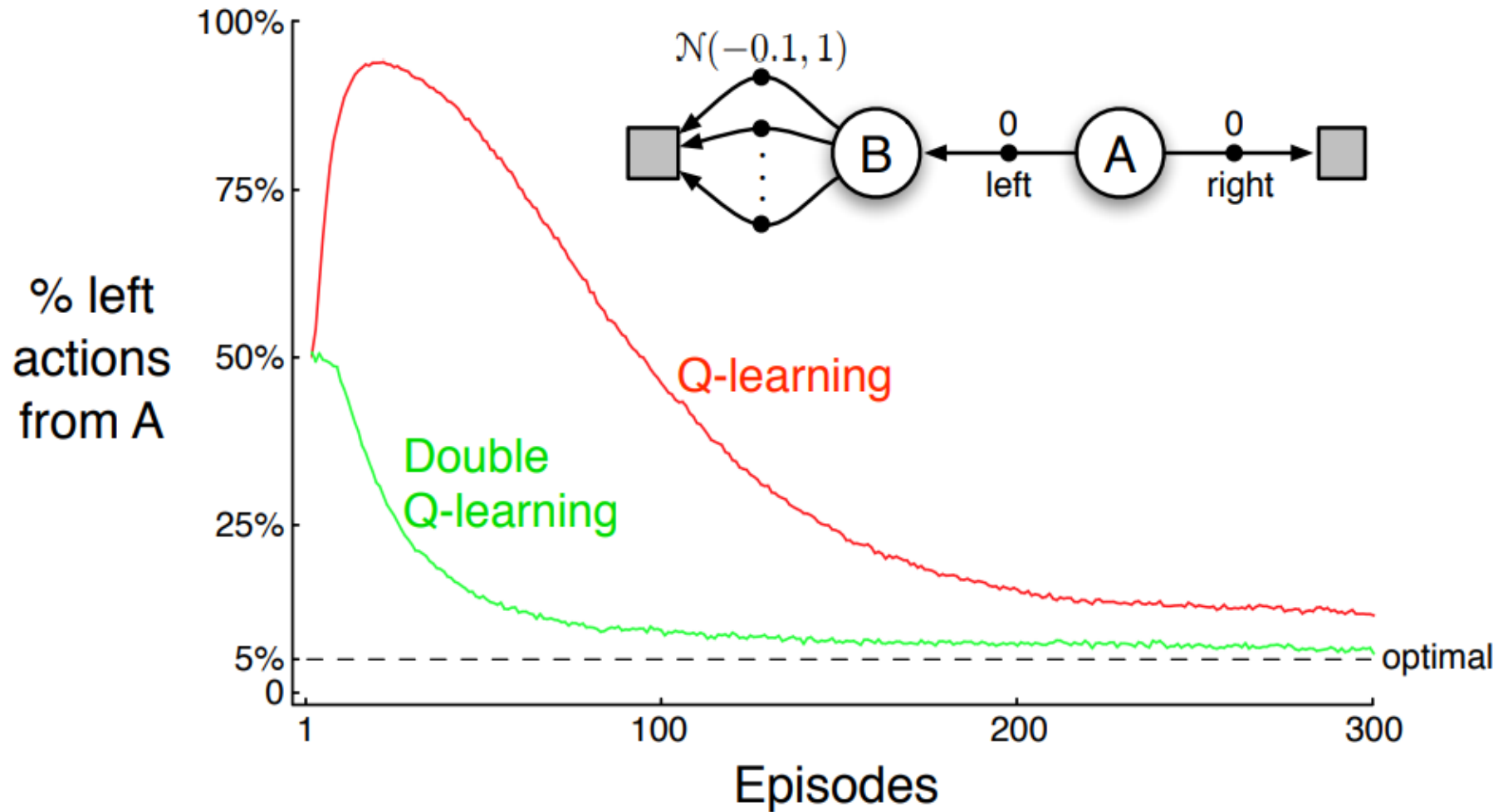
$S \leftarrow S'$

until S is terminal

Что изменилось для предтерминальных состояний?



Double Q-Learning





Что делать, если состояний много?

- Если пространство состояний является непрерывным, то можно его дискретизировать
- Проблема дискретизации: проклятие размерности
Количество состояний экспоненциально растет с ростом числа размерностей
- Если их просто много, но не бесконечно, то получается большая табличка
- Из-за большого количества состояний процесс обучения может затянуться
- Хочется использовать пространственную структуру состояний:
часто близкие состояния имеют близкие значения Q функции
- Не обязательно искать табличное отображение. Можно искать полноценную функцию



Что делать, если состояний много?

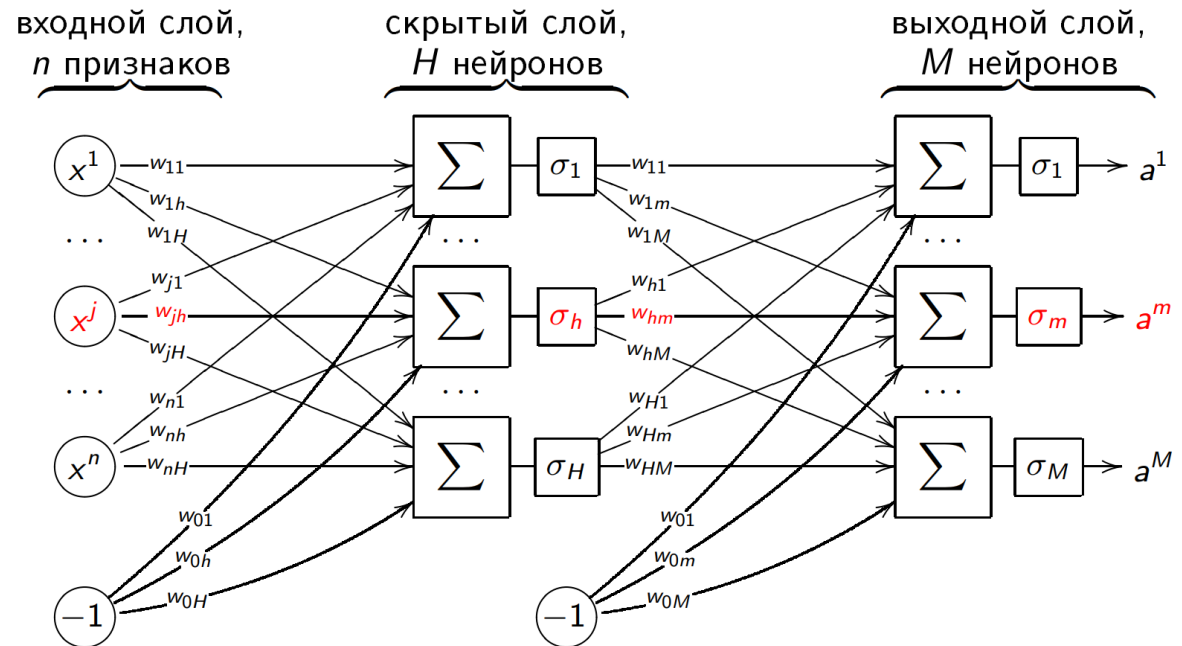
Опять проблема: функций много, непонятно как найти ту самую

Будем аппроксимировать:
выделим параметризованное семейство функций, которым можно хорошо приближать искомую Q-функцию

$$Q(s, a) = Q(s, a \mid \theta)$$

Как аппроксимировать:

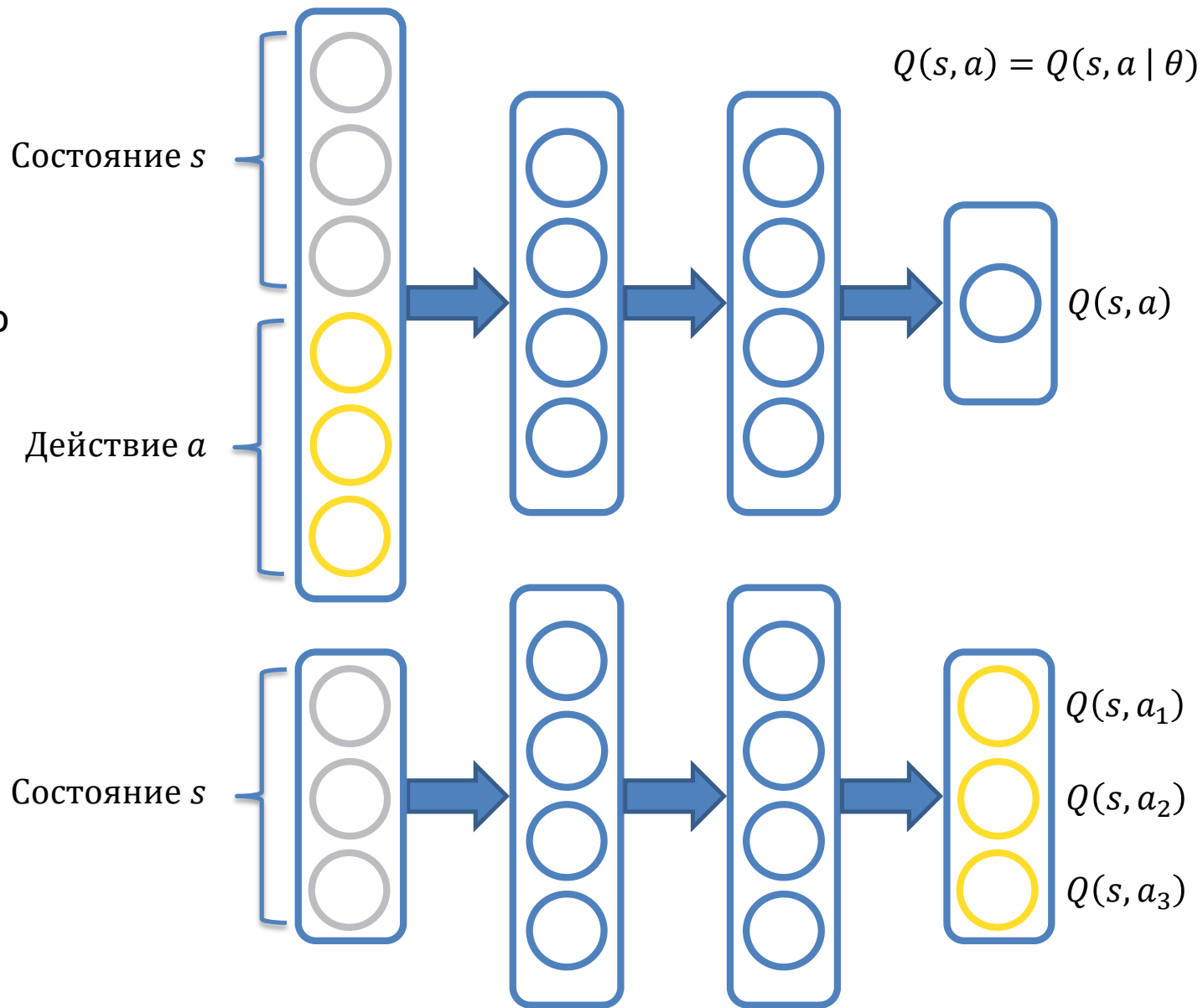
- Линейная регрессия
$$y = \theta_1 x + \theta_0$$
- Полиномы $y = \sum_0^n \theta_k x^k$
- Нейронная сеть



Варианты



Используем
эффективную
аппроксимацию





Посмотрим еще раз на Q-обучение

Табличный метод

$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} \max_{a'} Q^*(s', a')$ – уравнение Беллмана для Q^*

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

$$\begin{aligned} Q_{new}(s, a) &= Q_{old}(s, a) - \alpha \nabla_{Q(s, a)} F(Q) = \\ &= Q_{old}(s, a) - 2\alpha (Q_{old}(s, a) - r(s, a) - \gamma \max_{a'} Q_{old}(s', a')) \end{aligned}$$

Параметризация

$$Q(s, a) = Q(s, a | \theta)$$

$$\begin{aligned} \theta_{new} &= \theta_{old} - \alpha \nabla_{\theta} F(Q) = \\ &= \theta_{old} - 2\alpha (Q(s, a, \theta_{old}) - r(s, a) - \gamma \max_{a'} Q(s', a', \theta_{old})) \nabla_{\theta} Q(s, a, \theta) |_{\theta_{old}} \end{aligned}$$



1. Сделаем некоторое действие и получим (s_t, a_t, r_t, s_{t+1})

$$2. y_t = \begin{cases} r_t, & \text{если } s_t - \text{терминальное} \\ r_t + \gamma \max_a Q(s_{t+1}, a, \theta) \end{cases}$$

$$3. \theta \leftarrow \theta - \alpha \sum (Q(s_t, a_t, \theta) - y_t) \nabla_{\theta} Q(s_t, a_t, \theta)$$

Проблемы:

1. Мы решаем задачу регрессии, но данные коррелированы
2. Целевая переменная изменяется с обновлением θ



Использование памяти для Q-обучения

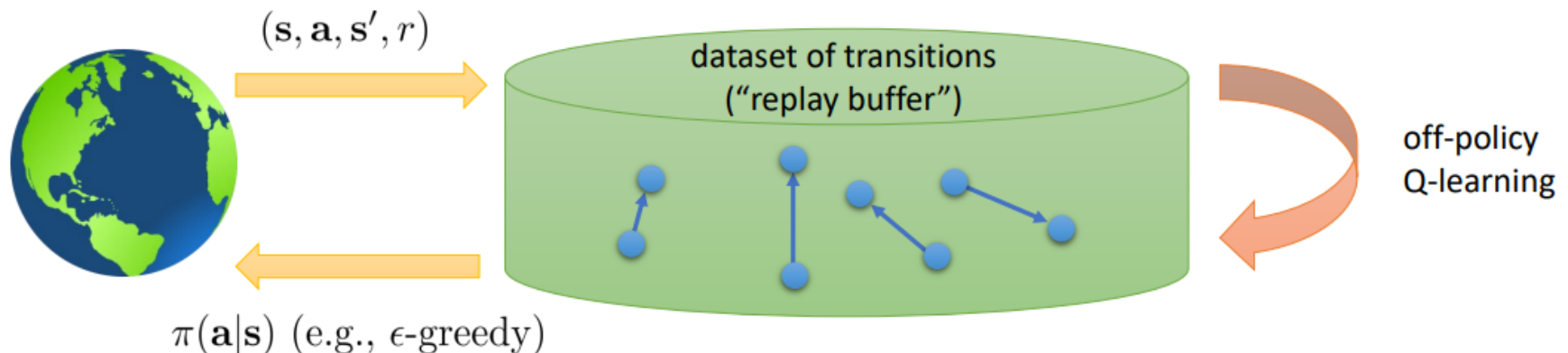
Есть какая-то $Q(s, a | \theta)$.

Действуем ϵ -жадно и получаем примеры $s_t, a_t \rightarrow r_t \rightarrow s_{t+1}$

Experience replay:

Сохраняем опыт $e_t = (s_t, a_t, r_t, s_{t+1})$ в память M

Делаем шаги обучения не только по новому опыту, но и по сохраненному в памяти





Q-обучение'

Инициализируем память M *двусторонней очередью*

Инициализируем функцию Q со случайными весами

for $episode_i$ in $[1 \dots N]$:

Инициализируем s_0

for t in $[0 \dots T]$:

$$a_t = \begin{cases} \operatorname{argmax}_a Q(s_t, a) & \text{с вероятностью } 1 - \varepsilon \\ \forall a \neq \operatorname{argmax}_a Q(s_t, a) & \text{с вероятностью } \frac{\varepsilon}{|A|-1} \end{cases}$$

Получаем s_{t+1}, r_t

Сохраняем $e_t = (s_t, a_t, r_t, s_{t+1})$ в память M

Сэмплируем минибатч (s_j, a_j, r_j, s'_j) из M

$$y_j = \begin{cases} r_j, & \text{если } s_j - \text{терминальное} \\ r_j + \gamma \max_{a'} Q(s'_j, a', \theta) & \end{cases}$$

$$\theta \leftarrow \theta - \alpha \sum (Q(s_j, a_j, \theta) - y_j) \nabla_{\theta} Q(s_j, a_j, \theta)$$

end for

end for



Target network

$$y_j = \begin{cases} r_j, & \text{если } s_j - \text{терминальное} \\ r_j + \gamma \max_{a'} Q(s'_j, a', \theta) \end{cases}$$

Таргет всё ещё двигается

$$\theta \leftarrow \theta - \alpha \sum (Q(s_j, a_j, \theta) - y_j) \nabla_{\theta} Q(s_j, a_j, \theta)$$

Решение: Зафиксируем веса θ' для таргета

$$y_j = \begin{cases} r_j, & \text{если } s_j - \text{терминальное} \\ r_j + \gamma \max_{a'} Q(s'_j, a', \theta') \end{cases} !$$

$$\theta \leftarrow \theta - \alpha \sum (Q(s_j, a_j, \theta) - y_j) \nabla_{\theta} Q(s_j, a_j, \theta)$$

И будем обновлять: $\theta' \rightarrow (1 - \tau)\theta' + \tau\theta$

Например, $\tau = 0.001$



DQN (Deep Q-Network)

$\theta' := \theta$

for $episode_i$ in $[1 \dots N]$:

Инициализируем s_0

for t in $[0 \dots T]$:

$$a_t = \begin{cases} \operatorname{argmax}_a Q(s_t, a) & \text{с вероятностью } 1 - \varepsilon \\ \forall a \neq \operatorname{argmax}_a Q(s_t, a) & \text{с вероятностью } \frac{\varepsilon}{|A|-1} \end{cases}$$

Получаем s_{t+1}, r_t

Сохраняем $e_t = (s_t, a_t, r_t, s_{t+1})$ в память M

Сэмплируем минибатч (s_j, a_j, r_j, s'_j) из M

$$y_j = \begin{cases} r_j, & \text{если } s_j - \text{терминальное} \\ r_j + \gamma \max_{a'} Q(s'_j, a', \theta) & \end{cases}$$

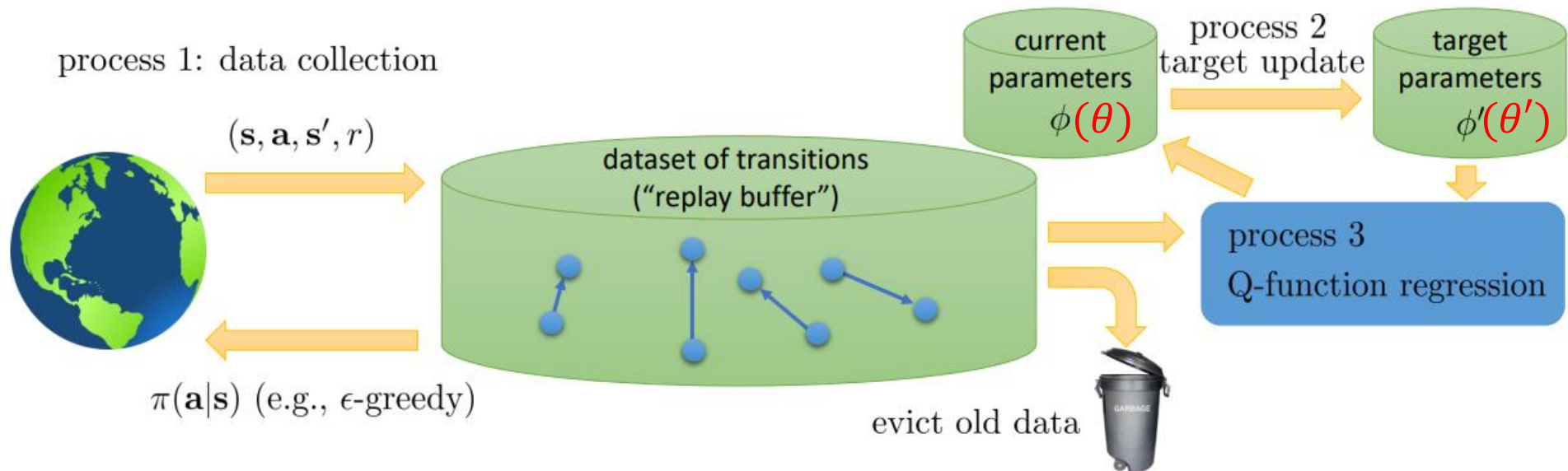
$$\theta \leftarrow \theta - \alpha \sum (Q(s_j, a_j, \theta) - y_j) \nabla_{\theta} Q(s_j, a_j, \theta)$$

$$\theta' \rightarrow (1 - \tau)\theta' + \tau\theta$$

end for

end for

DQN. Схема обучения



Можно ли совместить Double Q-Learning и Target network?



Target Network

$$y_j = \begin{cases} r_j, & \text{если } s_j - \text{терминальное} \\ r_j + \gamma \max_{a'} Q(s'_j, a', \theta') \end{cases}$$
$$\theta \leftarrow \theta - \alpha \sum (Q(s_j, a_j, \theta) - y_j) \nabla_{\theta} Q(s_j, a_j, \theta)$$
$$\theta' \rightarrow (1 - \tau)\theta' + \tau\theta$$

Double Q-Learning

$$y_j = \begin{cases} r_j, & \text{если } s_j - \text{терминальное} \\ r_j + \gamma Q(s'_j, \operatorname{argmax} Q(s'_j, a_j, \theta_1), \theta_2) \end{cases}$$
$$\theta_2 \leftarrow \theta_2 - \alpha \sum (Q(s_j, a_j, \theta_2) - y_j) \nabla_{\theta_2} Q(s_j, a_j, \theta_2)$$

Double Q-Learning with target network

$$y_j = \begin{cases} r_j, & \text{если } s_j - \text{терминальное} \\ r_j + \gamma Q(s'_j, \operatorname{argmax} Q(s'_j, a_j, \theta), \theta') \end{cases}$$
$$\theta \leftarrow \theta - \alpha \sum (Q(s_j, a_j, \theta) - y_j) \nabla_{\theta} Q(s_j, a_j, \theta)$$
$$\theta' \rightarrow (1 - \tau)\theta' + \tau\theta$$

State



Grayscale + scaling down

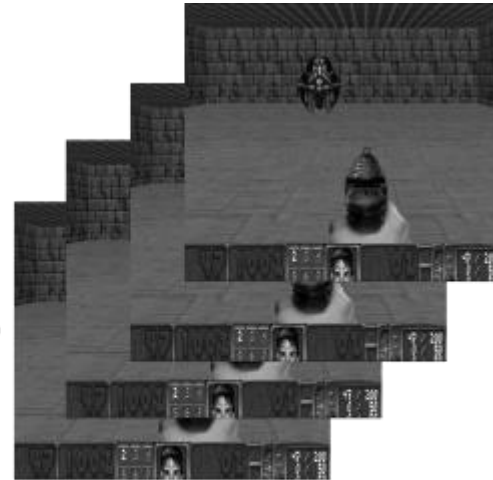


State frame

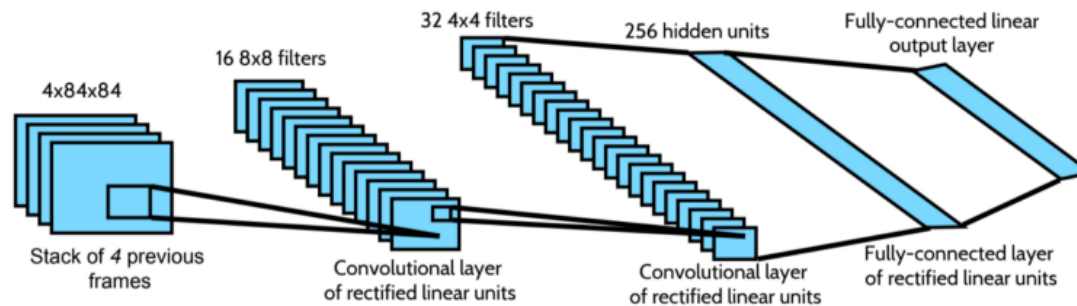


Preprocessed frame

+



Preprocessed stack of frames





Непрерывные действия

Почему это проблема для рассмотренных алгоритмов?

1. Нейронные сети?

- Как изменится архитектура?
- А как искать действие с максимальным Q -значением?
- Численная оптимизация - медленно

2. Дискретизация

- Часто непрерывные действия уже дискретизированы в среде:
Например, в гонках скорость машины регулируется 2-3мя действиями
- Теряем информацию
- Не делаем оптимальные действия

3. Случайно сэмплируем n действий;

Затем выбираем действие с максимальным Q

- Можно параллелить

Сделаем еще одну нейронную сеть $\mu(s|\phi)$ для выбора действия - actor

Как?

Мы хотим добиться $\mu(s|\phi) \cong \operatorname{argmax}_a Q(s, a|\theta)$

Значит задача:

$$\phi \rightarrow \operatorname{argmax}_{\phi} Q(s, \mu(s|\phi)|\theta)$$
$$Q(s, \mu(s|\phi)|\theta) \rightarrow \max_{\phi}$$

Градиентный подъём

$$\frac{dQ}{d\phi} = \frac{dQ}{d\mu} \frac{d\mu}{d\phi}$$

DDPG:

1. take some action \mathbf{a}_i and observe $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$, add it to \mathcal{B}
2. sample mini-batch $\{\mathbf{s}_j, \mathbf{a}_j, \mathbf{s}'_j, r_j\}$ from \mathcal{B} uniformly
3. compute $y_j = r_j + \gamma \max_{\mathbf{a}'} Q_{\phi'}(\mathbf{s}'_j, \mu_{\theta'}(\mathbf{s}'_j))$ using *target* nets $Q_{\phi'}$ and $\mu_{\theta'}$
4. $\phi \leftarrow \phi - \alpha \sum_j \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_j, \mathbf{a}_j)(Q_{\phi}(\mathbf{s}_j, \mathbf{a}_j) - y_j)$
5. $\theta \leftarrow \theta + \beta \sum_j \frac{d\mu}{d\theta}(\mathbf{s}_j) \frac{dQ_{\phi}}{d\mathbf{a}}(\mathbf{s}_j, \mathbf{a})$
6. update ϕ' and θ' (e.g., Polyak averaging)

Summary



- ✓ Обычный Q-learning может переоценивать истинные значения
 - ✓ Double Q-learning
- ✓ Нейронные сети отлично подходят для RL
 - ✓ Хорошо выделяют сложные нелинейные зависимости
 - ✓ Позволяют дообучаться по новым примерам
- ✓ DQN – естественное расширение табличного Q-обучения
 - ✓ Вместо таблиц используем параметризованную нейросеть Q-функцию
 - ✓ Стабилизируем обучение использованием памяти
 - ✓ Target network, чтобы целевая переменная была менее чувствительна к изменениям весов сети
- ✓ DDPG для работы с непрерывными действиями