

# Q-learning

Васильевых Павел, БПМИ213

# План

1. Условие задачи.
2. Алгоритм.
3. Улучшения алгоритма.

# Какую задачу решает?

$s$  — некоторое состояние

$a$  — некоторое действие

$r(s, a)$  — награда за заданное действие в заданном состоянии

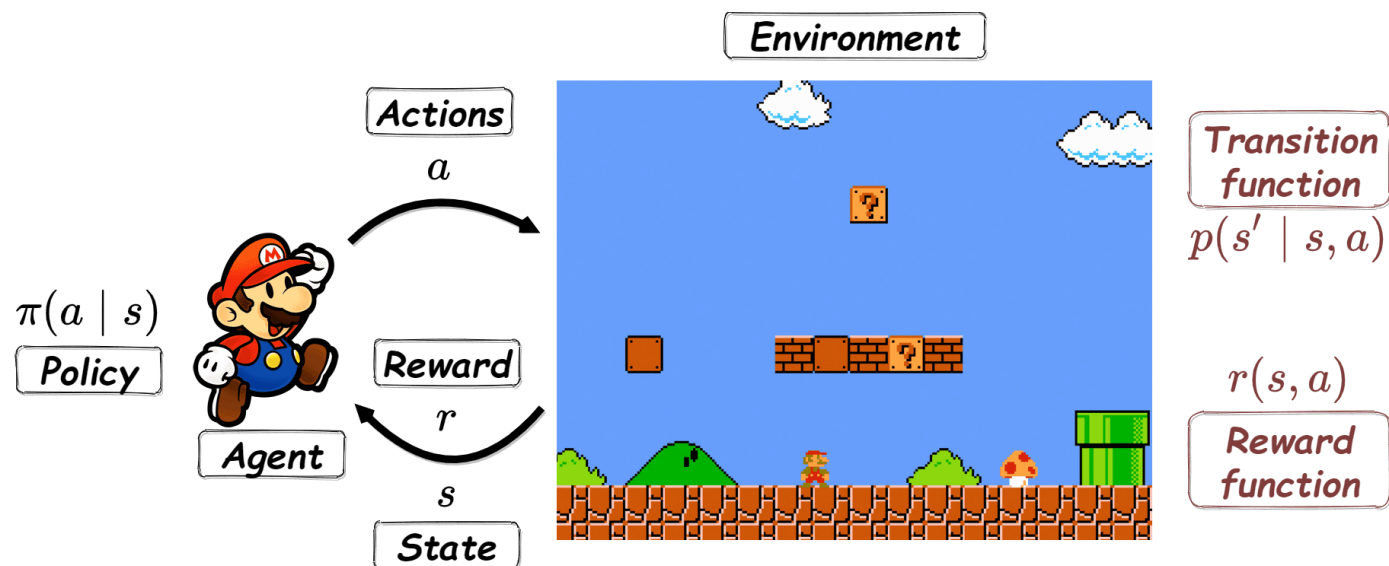
$\pi(a|s)$  — распределение действий агента в заданном состоянии (политика)

$\mathcal{T}$  — траектория (последовательность из реализованных состояний, действий и наград)

$\gamma$  — коэффициент убывания награды

$$Q(s, a) = \max_{\pi} \mathbb{E}_{\mathcal{T} \sim \pi | s_0=s, a_0=a} \sum_{t \geq 0} \gamma^t r(s_t, a_t)$$

Q-learning подбирает  $Q(s, a)$ , по которому можно определить  $\pi(a|s)$ .



# Как упростить задачу?

Во-первых, можно вывести уравнение оптимальности Беллмана:

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

Во-вторых, можно простым образом выразить оптимальную политику через  $Q$ :

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

# Как сделать задачу более практической?

1. Как-то инициализировать  $Q(s, a)$ .
2. Заметить, что можно забыть про  $p(s'|s, a)$  и моделировать переходы экспериментально, совершая действия  $\operatorname{argmax}_a Q(s, a)$ . В результате можно получить переходы вида  $(s, a, r, s')^a$ .

3. Вывести из уравнения Беллмана функцию перехода:

$$Q(s, a) \leftarrow r(s, a) + \gamma \max_{a'} Q(s', a')$$

4. Сгладить функцию перехода:

$$Q(s, a) \leftarrow (1 - \alpha) \cdot Q(s, a) + \alpha \cdot \left( r(s, a) + \gamma \max_{a'} Q(s', a') \right)$$

# СХОДИМОСТЬ

Пусть  $Q(s, a)$  и  $r(s, a)$  инициализированы,  $\gamma \in [0, 1)$ . Доказать, что процесс переходов  $Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(s'|s, a)} \max_{a'} Q(s', a')$  по всему множеству состояний и действий сходится.

Обозначим функционал перехода за  $R$ . Рассмотрим неизвестную истинную функцию  $Q^*(s, a)$ , для неё верно  $R \circ Q^* = Q^*$ . Заметим, что

$$\|R \circ Q - Q^*\|_\infty = \|R \circ Q - R \circ Q^*\|_\infty = \gamma \cdot \left\| \mathbb{E}_{s' \sim p(s'|s, a)} \left[ \max_{a'} Q(s', a') - \max_{a'} Q^*(s', a') \right] \right\|_\infty \leq \gamma \cdot \left\| \mathbb{E}_{s' \sim p(s'|s, a)} [Q(s', a(s')) - Q^*(s', a(s'))] \right\|_\infty \leq \gamma \cdot \|Q - Q^*\|_\infty,$$

где  $a(s') = \operatorname{argmax}_a Q(s', a)$ . Значит,  $0 \leq \lim_{n \rightarrow \infty} \|R^n \circ Q - Q^*\|_\infty \leq \lim_{n \rightarrow \infty} (\gamma^n \cdot \|Q - Q^*\|_\infty) = 0$ .

# Дилемма exploration-exploitation

Из-за детерминированности  $\operatorname{argmax}_a Q(s, a)$  вероятность попасть в некоторые состояния может быть равна 0, нужно добавить случайность, чтобы агент исследовал среду. Но она должна быть незначительной, чтобы агент использовал накопленный опыт.

$$\pi(s) = \begin{cases} \text{случайное действие, вероятность } \varepsilon \\ \operatorname{argmax}_a Q(s, a), \text{ вероятность } 1 - \varepsilon \end{cases}$$

# Что получилось?

Алгоритм, подходящий для  $Q(s, a)$ , которые можно задать таблицей:

1. Инициализировать  $Q(s, a)$ .
2. Получить из среды начальное состояние  $s_0$ .
3. Повторить для  $k = 0 \dots n$ :
  - 1) выбрать случайное действие  $a_k$  с вероятностью  $\varepsilon$ , иначе  $a_k = \underset{a}{\operatorname{argmax}} Q(s_k, a)$ ;
  - 2) получить награду  $r_k$  и следующее состояние  $s_{k+1}$ ;
  - 3) обновить  $Q(s_k, a_k) \leftarrow (1 - \alpha) \cdot Q(s_k, a_k) + \alpha \cdot \left( r_k + \gamma \max_{a'} Q(s_{k+1}, a') \right)$ .



# Таргет-сеть

В случае, когда состояний много, можно задать  $Q(s, a)$  нейросетью с обновляемыми таргетами  $r_k + \gamma \max_{a'} Q(s_{k+1}, a')$ . Но в такой модификации у алгоритма возникает ещё один недостаток — нестабильность сети из-за зависимости таргетов от предыдущих выходов.

Чтобы сделать сеть стабильной, нужна таргет-сеть. Создадим копию имеющейся сети и будем использовать её для обновления таргетов, заменяя её веса параметрами основной сети раз в  $t$  шагов обучения.

# Experience replay

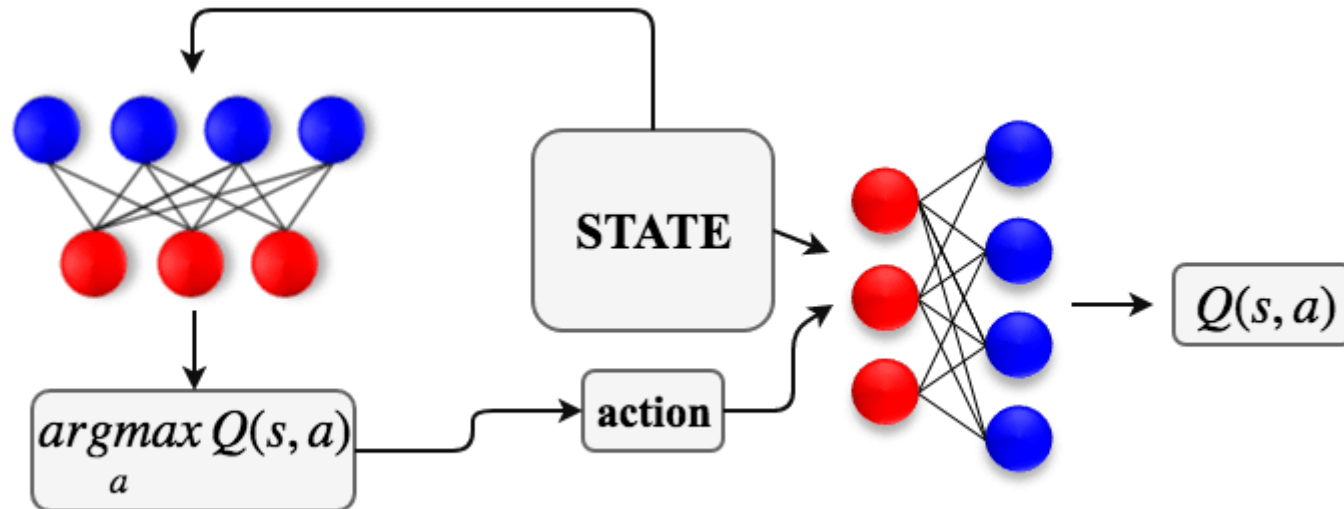
Для качественного обучения нейросети обычно нужно усреднять градиент по батчу. Если мы будем следить за последовательными действиями одного агента, сеть будет забывать предыдущий опыт.

Решения:

- 1) Несколько агентов одновременно;
- 2) experience replay, то есть запись всех переходов в буфер и семплирование случайных батчей из него.

# А если действий много?

Сделаем ещё одну нейросеть  $\pi(s)$  для приближения  $\operatorname{argmax}_a Q(s, a)$ . Её можно обучать параллельно с основной: на каждом шаге алгоритма будем делать шаг градиентного подъёма для батча функций  $Q(s, \pi(s))$ , где параметры сети  $Q(s, a)$  зафиксированы.



# Deep Q-learning (Actor-Critic)

1. Одинаково инициализировать  $Q(s, a)$  и таргет-сеть.
2. Инициализировать  $\pi(s)$ .
3. Получить из среды начальное состояние  $s_0$ .
4. Повторить для  $k = 0 \dots n$ :
  - 1) выбрать случайное действие  $a_k$  с вероятностью  $\varepsilon$ , иначе  $a_k = \underset{a}{\operatorname{argmax}} Q(s_k, a)$ ;
  - 2) получить награду  $r_k$  и следующее состояние  $s_{k+1}$ ;
  - 3) добавить переход  $(s_k, a_k, r_k, s_{k+1})$  в буфер;
  - 4) сэмплировать батч из буфера;
  - 5) сделать шаг градиентного подъёма для  $Q(s, \pi(s))$  по параметрам  $\pi(s)$ ;
  - 6) посчитать таргеты  $r + \gamma Q(s', \pi(s'))$ , используя таргет-сеть;
  - 7) сделать шаг градиентного спуска для  $Q(s, a)$ ;
  - 8) если  $k : m$ , обновить веса таргет-сети.

# Модификации

DDPG (Deep deterministic policy gradient): для  $Q$  и  $\pi$  есть по две нейросети (одна меняется на каждом шаге, другая постепенно двигает веса к весам первой), действия генерируются из нормального распределения с центром в значении меняющейся  $\pi$ , а таргеты для обучения  $Q$  считаются с помощью двух стабильных сетей. Есть проблема переоценки, решается добавлением второй  $Q$ .

SAC (Soft actor-critic): в этом алгоритме случайные (исследовательские) шаги не создаются отдельно или поверх выхода  $\pi$ . Случайность  $\pi$  включается в задачу оптимизации через прибавление энтропии текущей  $\pi$  с константным весом к каждой награде  $r$ .

# Выводы

Мы рассмотрели алгоритм, для обучения которого можно использовать весь накопленный опыт (off-policy). Очевидно, есть большой простор для улучшений и оптимизаций под конкретные задачи.

Однако мы могли бы подойти к задаче RL и со стороны функции политики. Тогда бы мы попали в область on-policy алгоритмов, которые могут оказаться эффективнее, но это тема следующего семинара.

# Ссылки

Статья про Q-learning:

[https://education.yandex.ru/handbook/ml/article/obuchenie-s-podkrepleniem.](https://education.yandex.ru/handbook/ml/article/obuchenie-s-podkrepleniem)

Лекция про DDPG и SAC:

[https://dzen.ru/video/watch/6521821790b03d45f18118ce.](https://dzen.ru/video/watch/6521821790b03d45f18118ce)

Применение Deep Q-network:

[https://deepmind.com/blog/article/Agent57-Outperforming-the-human-Atari-benchmark.](https://deepmind.com/blog/article/Agent57-Outperforming-the-human-Atari-benchmark)