

Q-learning

Биршерт Алексей

Февраль 2020

- Марковский процесс принятия решений
- Равенство Беллмана
- on-policy и off-policy
- Q-learning
- Deep Q-learning
- Experience replay
- вопросы

Марковский процесс принятия решений

Взаимодействие агента и среды

Общее представление

- Агент совершает действие $A_t \in A(s)$
- Агент получает награду $R_{t+1} \in R \subset \mathbb{R}$, новое состояние $S_{t+1} \in S$
- $S_0, A_0, R_1, S_1, \underbrace{A_1, R_2, S_2}_{\text{элемент цикла}}, A_2, R_3, S_3, A_3 \dots$

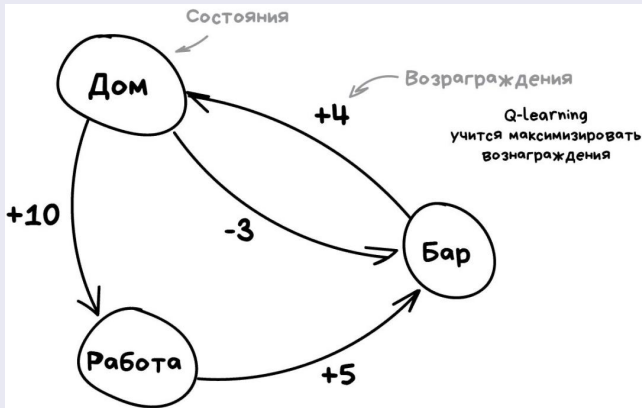
Вероятности

- $p(s', r | s, a) \doteq \Pr(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$
- $\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1$ для всех $s \in S, a \in A(s)$

Марковский процесс принятия решений

Взаимодействие агента и среды

Мем



Рутинный Марковский Процесс

Марковский процесс принятия решений

Ожидаемая награда агента

- $G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, \quad 0 \leq \gamma \leq 1$
- $G_t = R_{t+1} + \gamma G_{t+1}$
- γ определяет насколько важна будущая награда сейчас
- Если $\gamma = 0$, то важна только награда на текущем шагу, если $\gamma = 1$, то все важны равнозначно

Марковский процесс принятия решений

Policy и value функции

Policy функция

$$\pi(a|s) \doteq Pr(A_t = a | S_t = s), \quad a \in A(s), s \in S$$

Задаёт вероятностное распределение возможных действий в конкретном состоянии

Марковский процесс принятия решений

Policy и value функции

State-value функция

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

Задаёт ожидаемую награду, если мы стартуем из состояния s и следуем policy π

Action-value функция

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

Задаёт ожидаемую награду, если стартуем из состояния s и совершаем действие a , и потом следуем policy π

Связь состояния и последующих

$$\begin{aligned}v_{\pi}(s) &\doteq \mathbb{E}_{\pi} [G_t | S_t = s] \\&= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] \\&= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_{\pi} [G_{t+1} | S_{t+1} = s']] \\&= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]\end{aligned}$$

Равенство Беллмана

Оптимальность

Оптимальная policy функция

- $\pi \geq \pi' \Leftrightarrow v_{\pi}(s) \geq v_{\pi'}(s), \quad \forall s \in S$
- Оптимальная π - та, которая лучше или равна всем остальным
- Оптимальной policy функции соответствуют оптимальные state-value и action-value функции
- $v_*(s) \doteq \max_{\pi} v_{\pi}(s), \quad q_*(s, a) \doteq \max_{\pi} q_{\pi}(s, a)$

Равенство Беллмана

Оптимальность

Оптимальная state-value функция

$$\begin{aligned} v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\ &= \max_a \mathbb{E}_{\pi_*} [G_t | S_t = s, A_t = a] \\ &= \max_a \mathbb{E}_{\pi_*} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\ &= \max_a \mathbb{E} [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')] \end{aligned}$$

Ожидаемая награда у состояния при оптимальной policy должна быть равна ожидаемой награде лучшего действия из этого состояния.

Равенство Беллмана

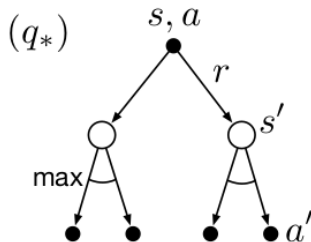
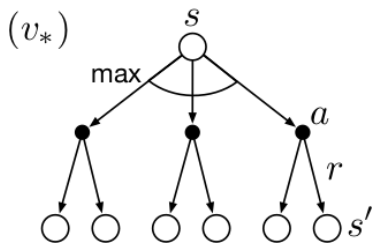
Оптимальность

Оптимальная action-value функция

$$\begin{aligned} q_*(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r | s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right] \end{aligned}$$

Равенство Беллмана

Оптимальность



on-policy и off-policy методы

on-policy

Пытаются оценить или оптимизировать policy функцию, в соответствии с которой действуют сами.

off-policy

Используют две различные policy функции - behavior и target policy, первая регулирует их исследование МППР, вторая оптимизируется в процессе исследования.

on-policy и off-policy методы

Q-learning и SARSA

SARSA, on-policy

$$Q(S_t, A_t) \leftarrow (1 - \alpha)Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})]$$

Q-learning, off-policy

$$Q(S_t, A_t) \leftarrow (1 - \alpha)Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) \right]$$

Алгоритм, $Q \approx q_*$

Инициализация $Q(\cdot, \cdot)$ произвольно

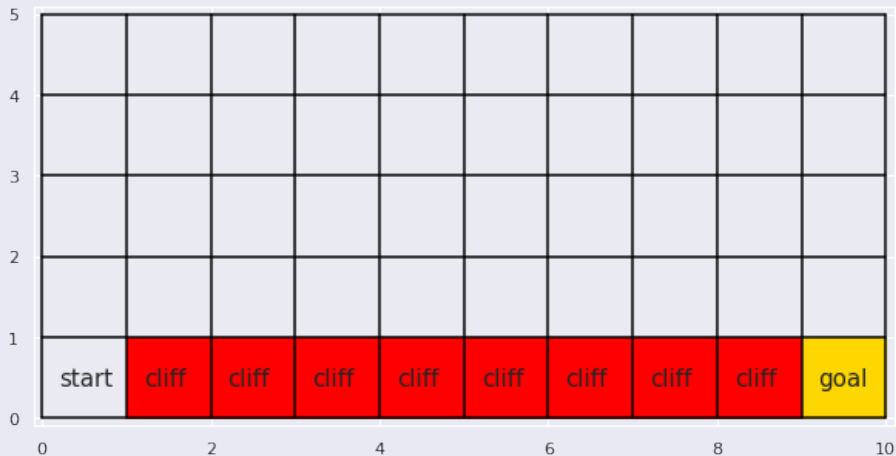
Цикл:

- Инициализировать s
- Цикл:
 - Выбрать шаг a из s с помощью behavior policy (ε – greedy, UCB)
 - Сделать шаг a , получить r и s'
 - $Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha \left[r + \max_{a'} \gamma Q(s', a') \right]$
 - $s \leftarrow s'$

Q-learning

Сравнение с SARSA на основе задачи cliff walking

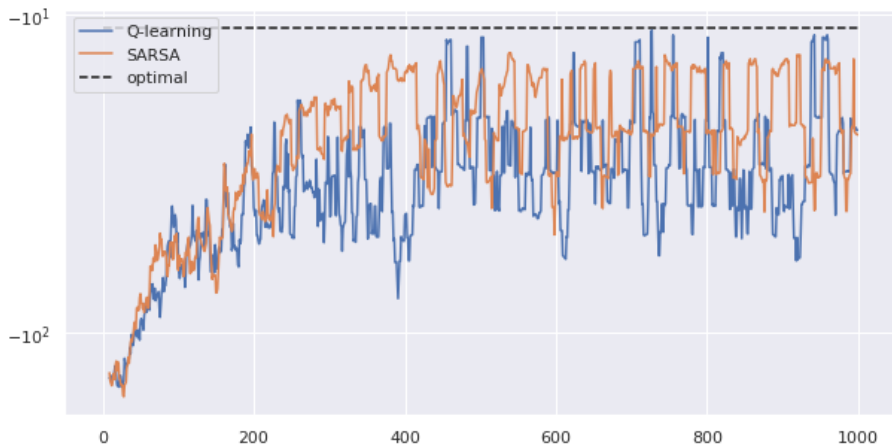
Постановка задачи



Q-learning

Сравнение с SARSA на основе задачи cliff walking

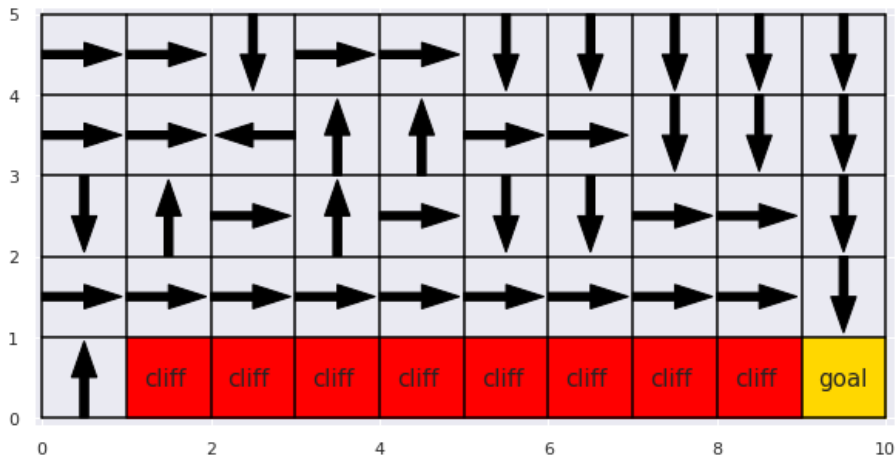
Обучение, $\varepsilon = 0.1$



Q-learning

Сравнение с SARSA на основе задачи cliff walking

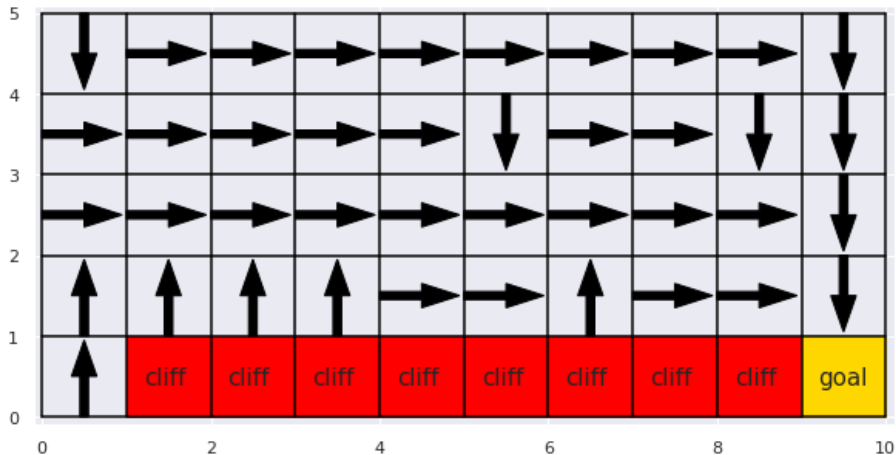
Q-learning, $\varepsilon = 0.1$



Q-learning

Сравнение с SARSA на основе задачи cliff walking

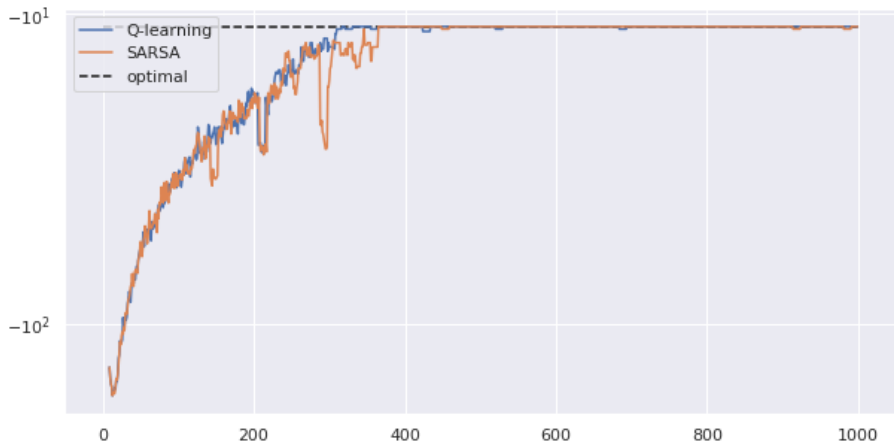
SARSA, $\epsilon = 0.1$



Q-learning

Сравнение с SARSA на основе задачи cliff walking

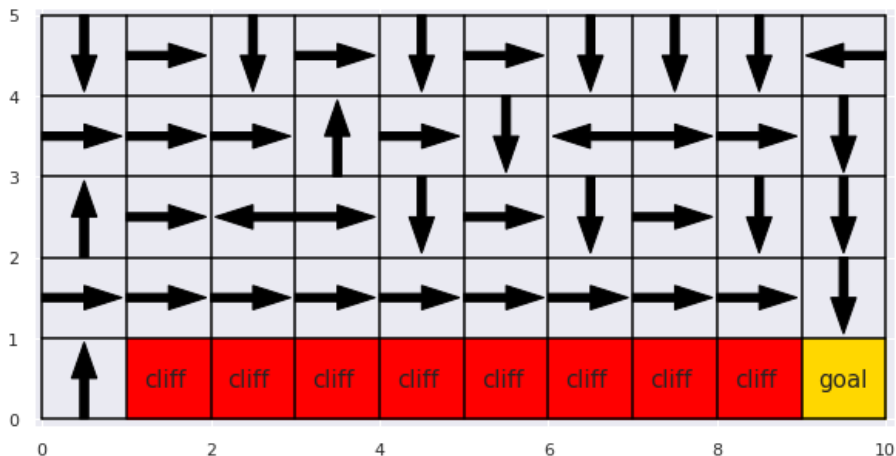
Обучение, $\varepsilon = 0.001$



Q-learning

Сравнение с SARSA на основе задачи cliff walking

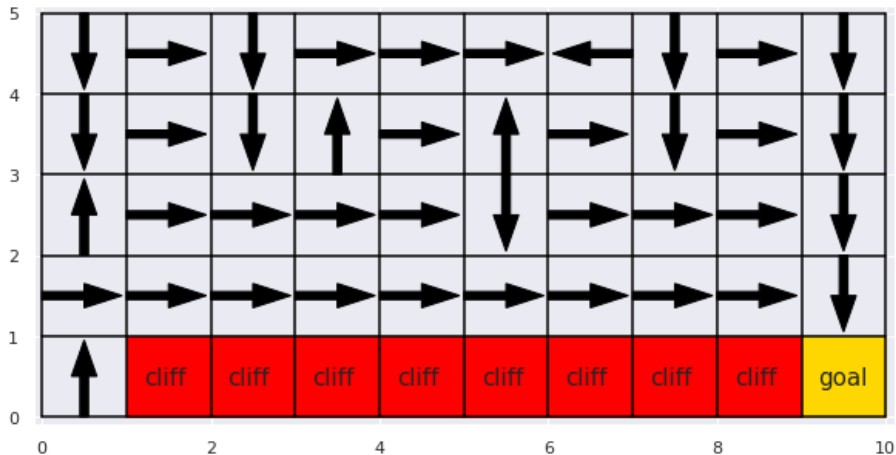
Q-learning, $\varepsilon = 0.001$



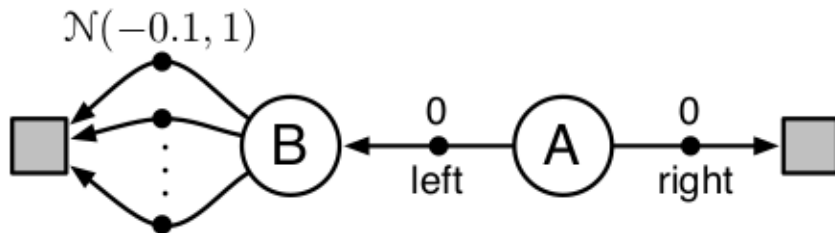
Q-learning

Сравнение с SARSA на основе задачи cliff walking

SARSA, $\epsilon = 0.001$



Максимизация смещения



Q-learning

Double Q-learning

Алгоритм, $Q_1 \approx Q_2 \approx q_*$

Инициализация $Q_1(\cdot, \cdot)$ и $Q_2(\cdot, \cdot)$ произвольно

Цикл:

- Инициализировать s
- Цикл:
 - Выбрать шаг a из s с помощью ϵ - *greedy* из $Q_1 + Q_2$
 - Сделать шаг a , получить r и s'
 - С вероятностью 0.5:

$$Q_1(s, a) \leftarrow (1 - \alpha)Q_1(s, a) + \alpha \left[r + \gamma Q_2(s', \arg \max_{a'} Q_1(s', a')) \right]$$

- Иначе:

$$Q_2(s, a) \leftarrow (1 - \alpha)Q_2(s, a) + \alpha \left[r + \gamma Q_1(s', \arg \max_{a'} Q_2(s', a')) \right]$$

- $s \leftarrow s'$

Интуиция позади использования DQN, $Q(s, a, \theta) \approx Q_*(s, a)$

Q-learning имеет проблемы со сходимостью в пространствах больших размерностей и не может обобщить предыдущий опыт для новых ситуаций

$Q(S_t, A_t, \theta) \leftarrow (1 - \alpha)Q(S_t, A_t, \theta) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a, \theta)]$, θ - веса нейронной сети

$Q(S_t, A_t, \theta)$ - предсказание

$R_{t+1} + \gamma \max_a Q(S_{t+1}, a, \theta)$ - целевая переменная

$R_{t+1} + \gamma \max_a Q(S_{t+1}, a, \theta) - Q(S_t, A_t, \theta)$ - ошибка

Интуиция позади использования DQN, $Q(s, a, \theta) \approx Q_*(s, a)$

$Q(S_t, A_t, \theta) \leftarrow (1 - \alpha)Q(S_t, A_t, \theta) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a, \theta)]$, θ -
веса нейронной сети

$$L(\theta) = \mathbb{E}_{s,a,r} \left[\left(\mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q(s', a', \theta) \mid s, a \right] - Q(s, a, \theta) \right)^2 \right]$$

$$\Delta_{\theta} L(\theta) = \mathbb{E}_{s,s',a,r} \left[\left(r + \gamma \max_{a'} Q(s', a', \theta) - Q(s, a, \theta) \right) \Delta_{\theta} Q(s, a, \theta) \right]$$

Проблемы

Однако такая нейронная сеть может не сойтись:

- Корреляция в последовательности наблюдений
 - Решение - experience replay
- Корреляция между целевой переменной и предсказанием
 - Решение - fixed Q targets

Решение проблемы с корреляцией в последовательности наблюдений

Используется буфер, в котором хранятся последние n наблюдений. Нейронная сеть обучается на мини-батчах, сэмплируемых из буфера.

Преимущества:

- Снижает влияние отдельных наблюдений на обучение
- Убирает корреляцию в обучающем батче, что позволяет нейросети сходиться лучше
- Повышает скорость обучения
- Позволяет обучать параллельно

- В чем заключается bellman optimality equation?
- В чем заключаются off-policy и on-policy methods? Почему Q-learning считается off-policy method?
- Что такое experience replay?

Список литературы

- 1 Richard S. Sutton and Andrew G. Barto. Reinforcement Learning: An Introduction. 2018
- 2 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dhharshan Kumaran, Daan Wierstra, Shane Legg & Demis Hassabis. Human-level control through deep reinforcement learning. 2015