



ОНЛАЙН-ОБРАЗОВАНИЕ

Обучение с подкреплением.

Он всегда был не прочь подкрепиться.
Кроме того, он было поэт.

Артур Кадулин
Преподаватель



- 1. Обучение с подкреплением**
2. Бандиты
3. Стратегии



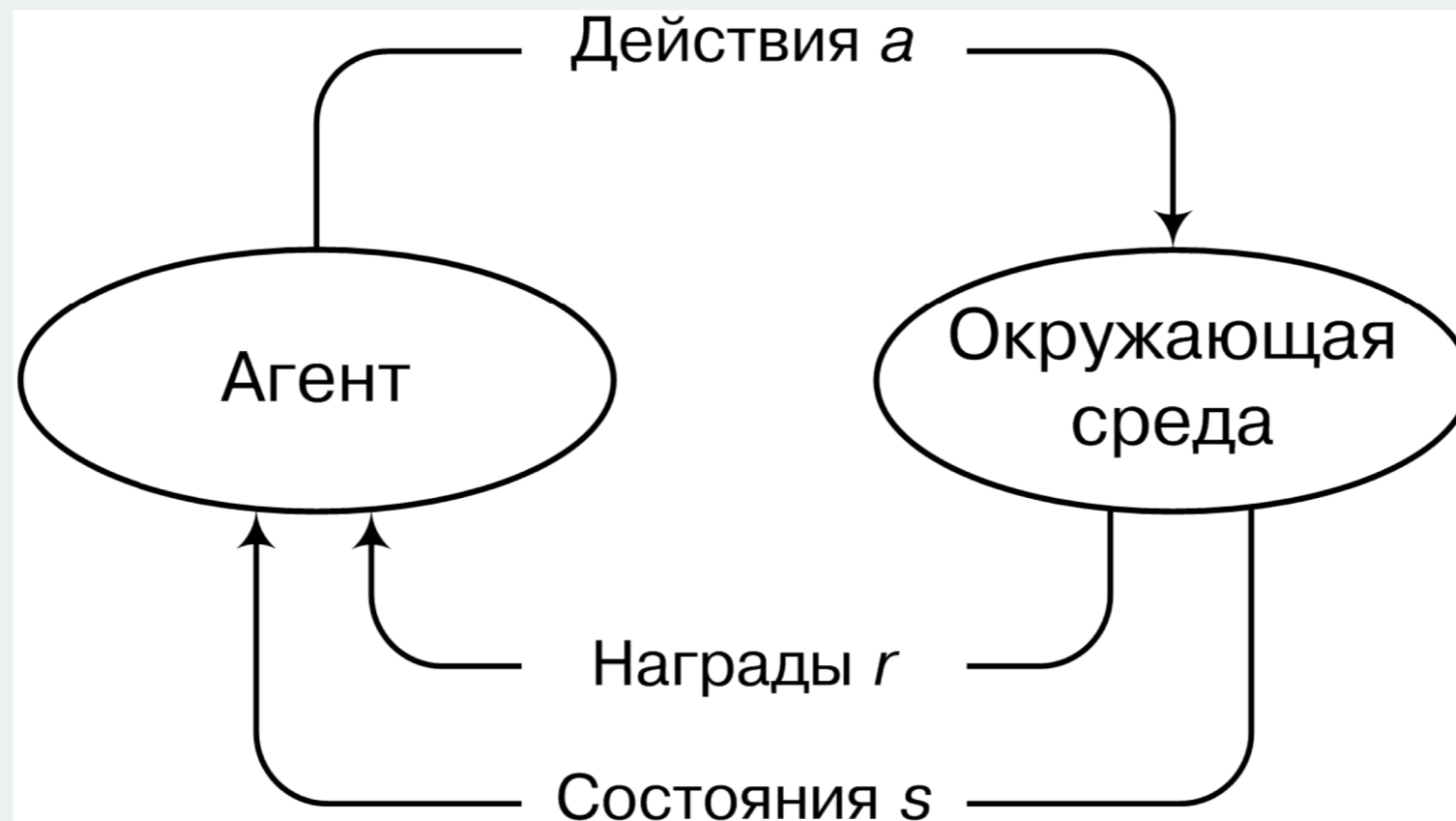
Обучение с учителем: Для заранее заданного набора объектов известны «правильные» ответы. Нужно научиться предсказывать ответы для новых объектов.

Обучение с подкреплением: Можно совершить действие и узнать ответ. Нужно научиться совершать действия приводящие к «лучшим» ответам.



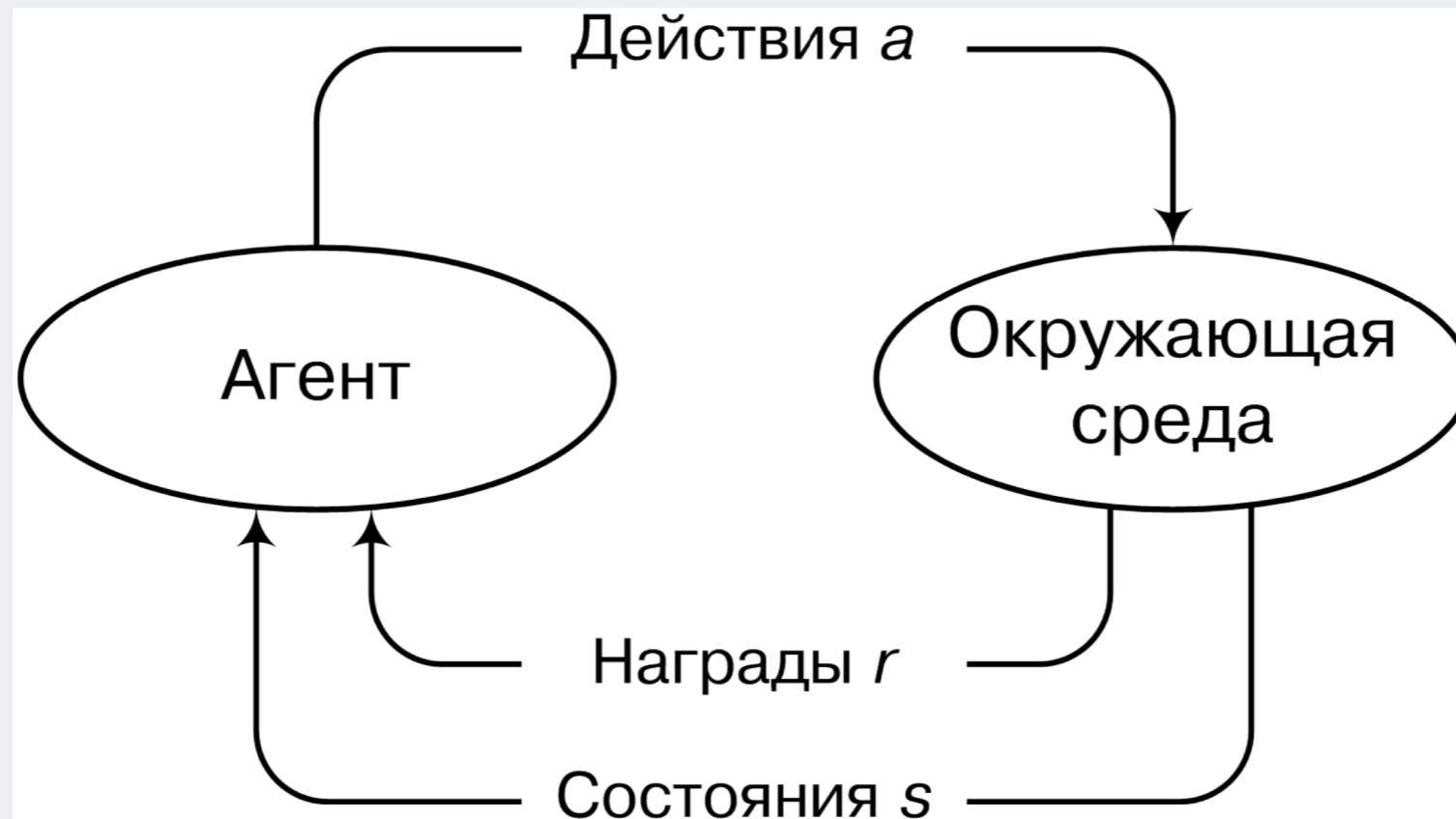
Обучение с учителем: Для заранее заданного набора объектов известны «правильные» ответы. Нужно научиться предсказывать ответы для новых объектов.

Обучение с подкреплением: Можно совершить действие и узнать ответ. Нужно научиться совершать действия приводящие к «лучшим» ответам.

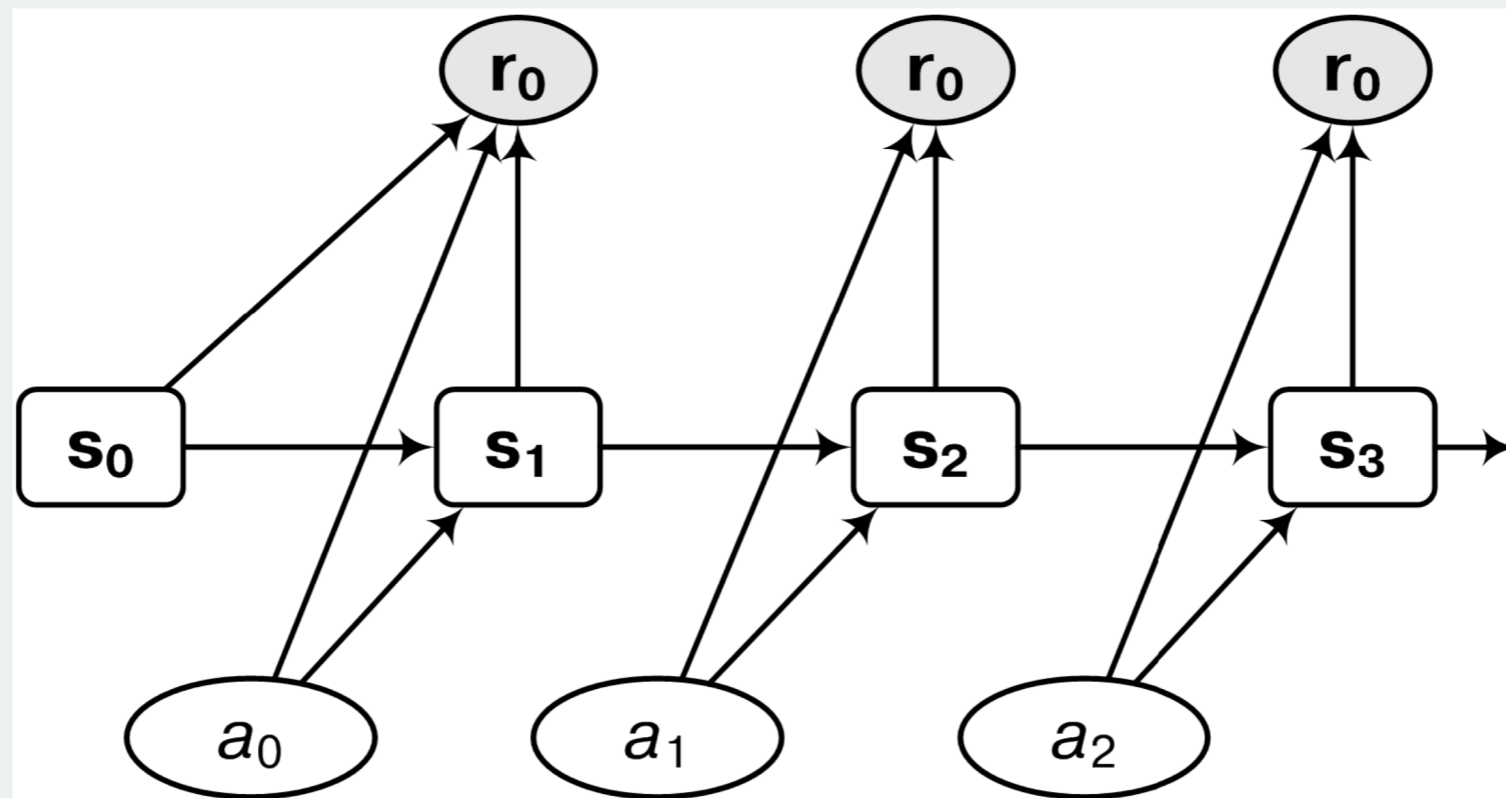


Марковский процесс принятия решений состоит из:

- Множества состояний S
- Множества действий A
- Функции вознаграждения $R: S \times A \rightarrow \mathbb{R}$, которая задает $R_{SS'}^a$,
- Функции перехода между состояниями $p_S^a: S \times A \rightarrow \Pi(S)$, которая задает $P_{SS'}^a$,

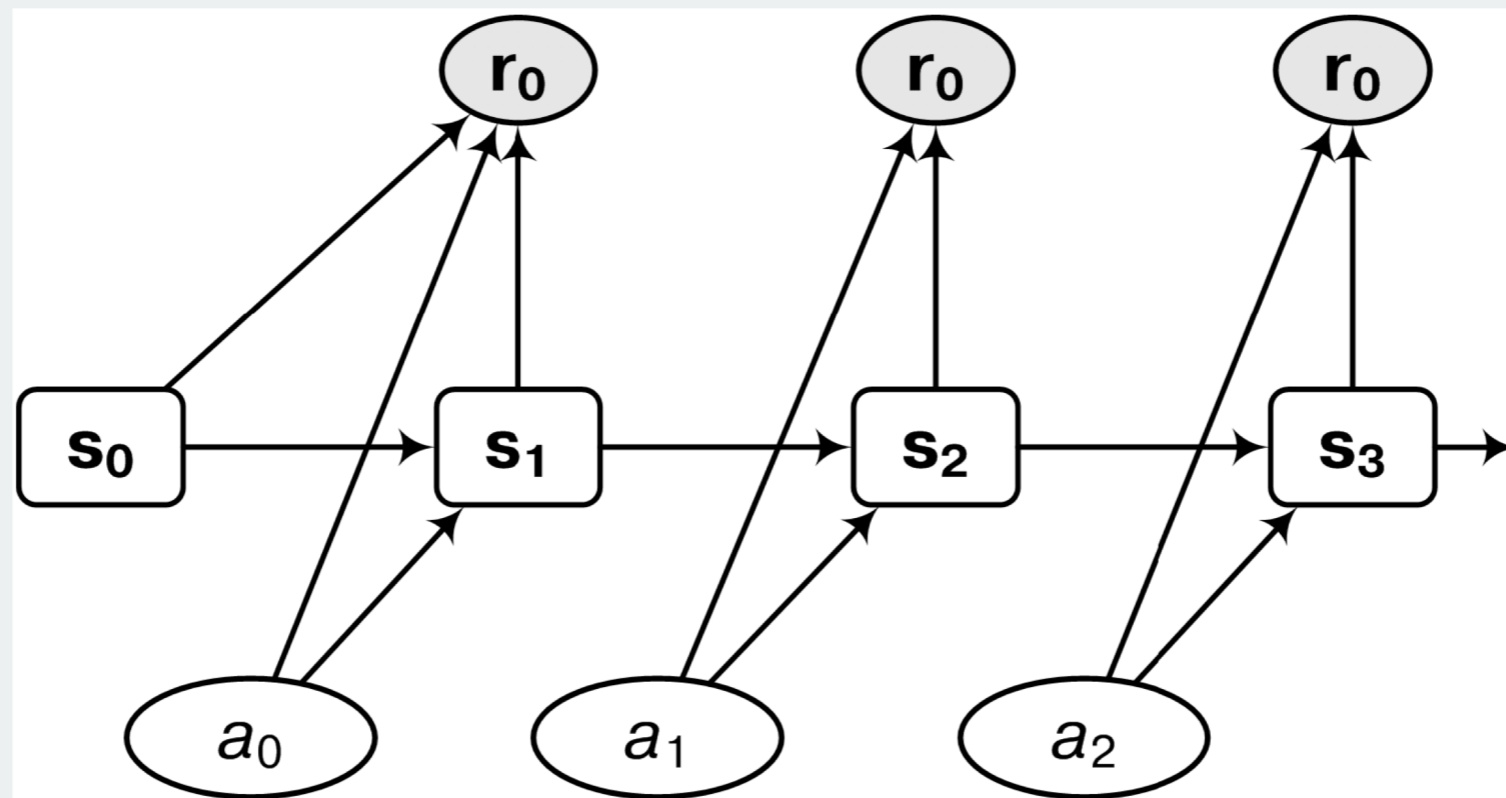


Марковское свойство заключается в том, что вероятности переходов между состояниями при выбранных действиях не зависят от истории предыдущих переходов.



Марковское свойство заключается в том, что вероятности переходов между состояниями при выбранных действиях не зависят от истории предыдущих переходов.

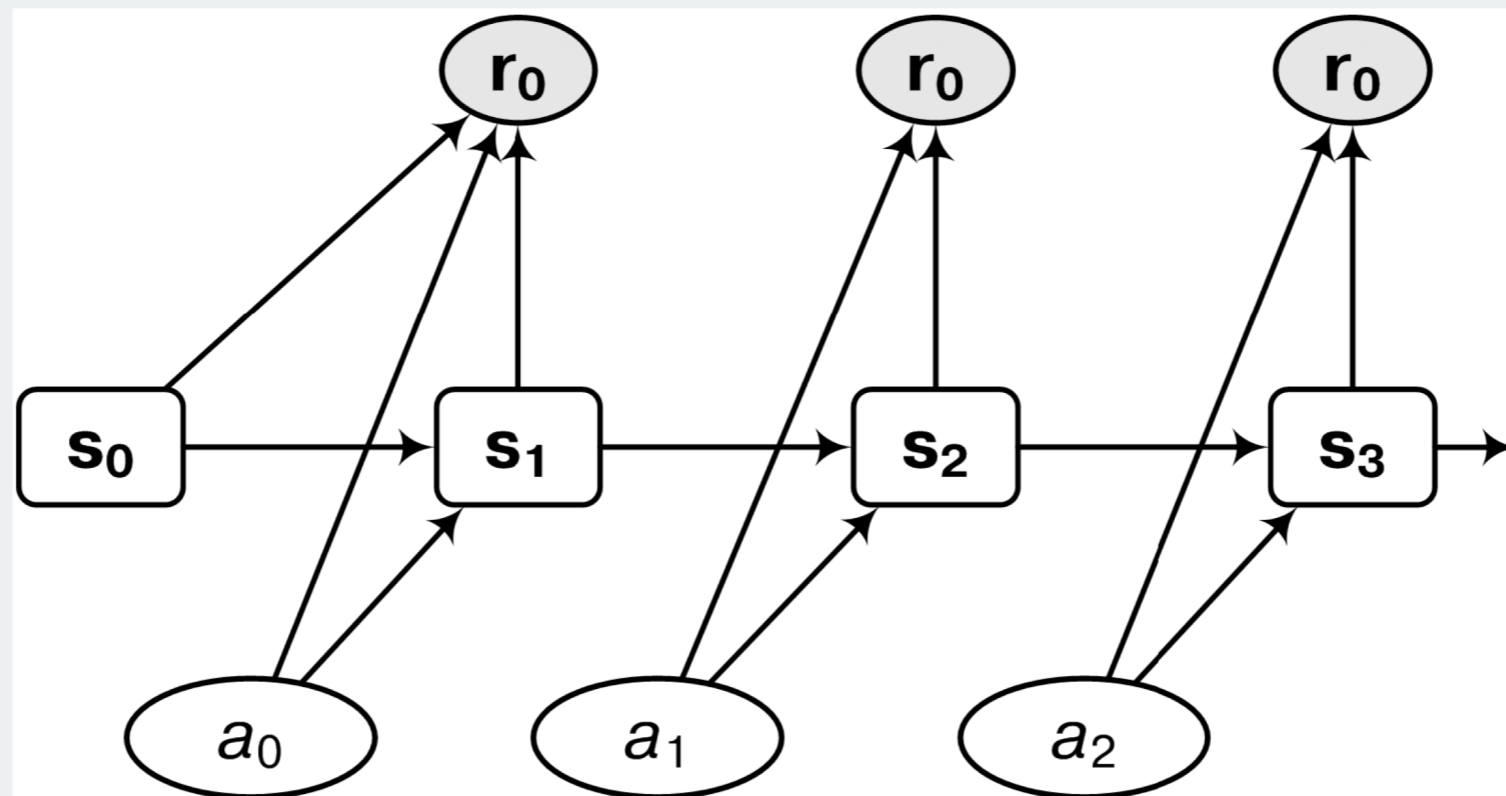
Выполняется ли это для шахмат?



Марковское свойство заключается в том, что вероятности переходов между состояниями при выбранных действиях не зависят от истории предыдущих переходов.

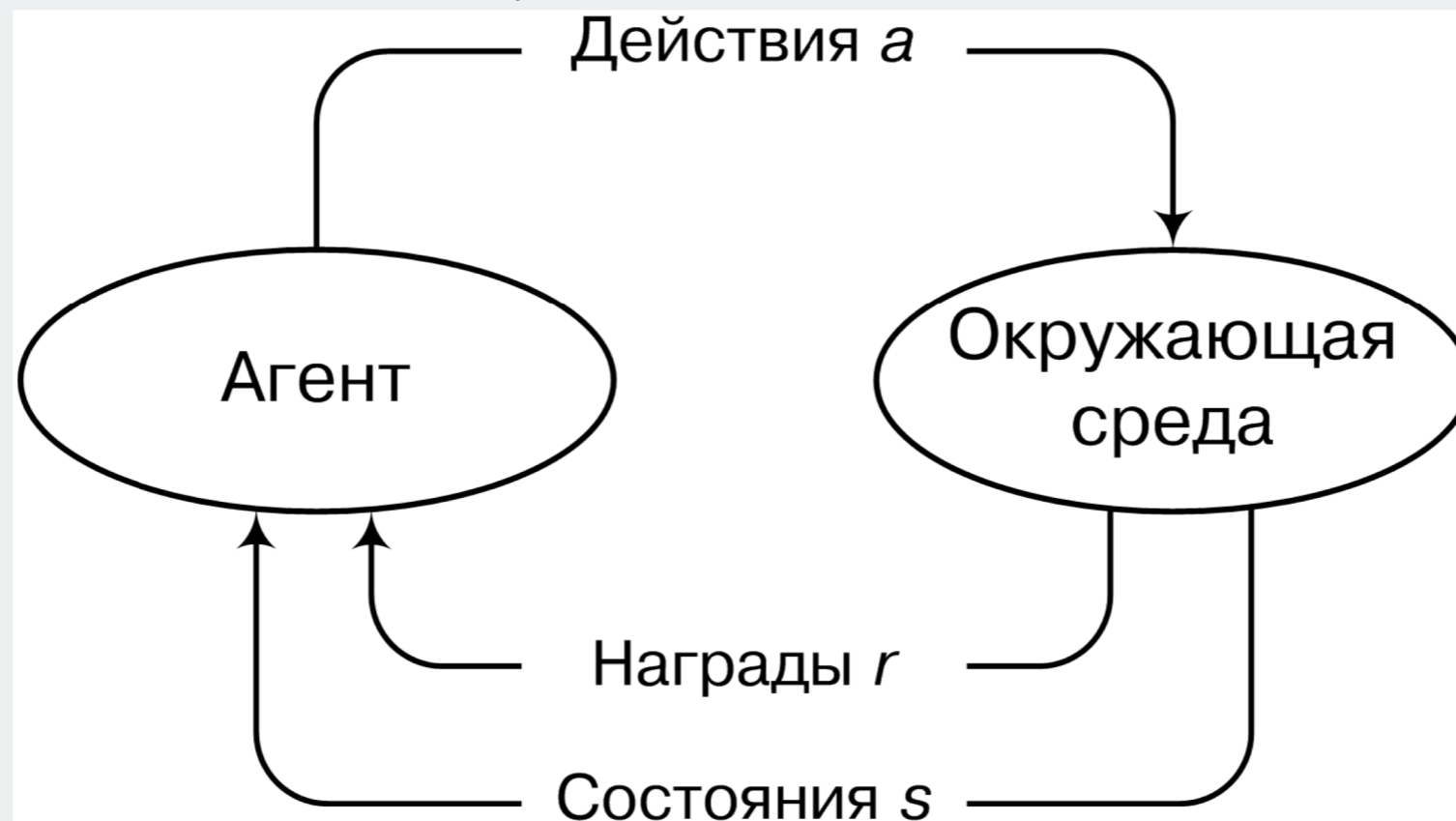
Выполняется ли это для шахмат?

А для покера?



В реальных задачах вознаграждение совсем не обязательно следует сразу же за совершенным действием, поэтому награда обычно оценивается:

- Либо за отведенное время h : $R = \mathbb{E}[\sum_{t=0}^h r_t]$
- Либо за бесконечное время с дисконтом: $R = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t]$
- Либо в среднем за 1 шаг: $R = \lim_{h \rightarrow \infty} \mathbb{E}[\frac{1}{h} \sum_{t=0}^h r_t]$

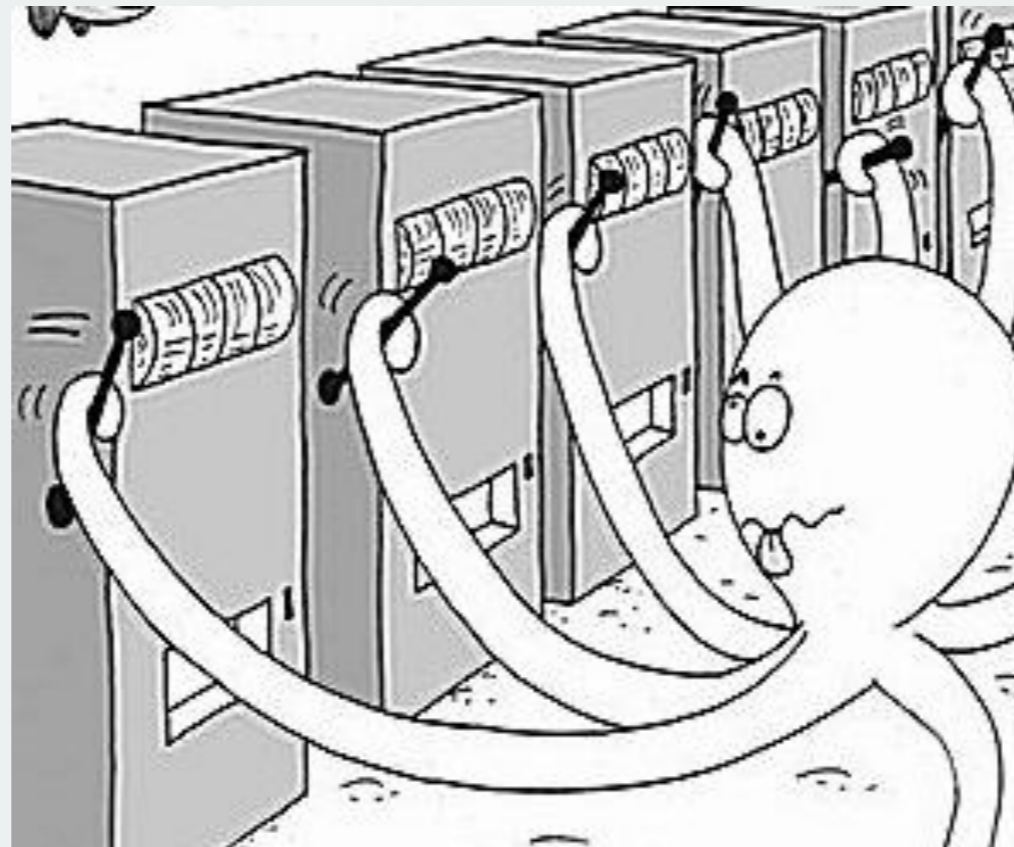


1. Обучение с подкреплением
- 2. Бандиты**
3. Стратегии



Допустим вы попали в комнату с кучей одноруких бандитов, каждый из которых выдает выигрыш со своей собственной вероятностью, которую вы не знаете. Ваша задача выиграть как можно больше (проиграть как можно меньше) денег за ночь.

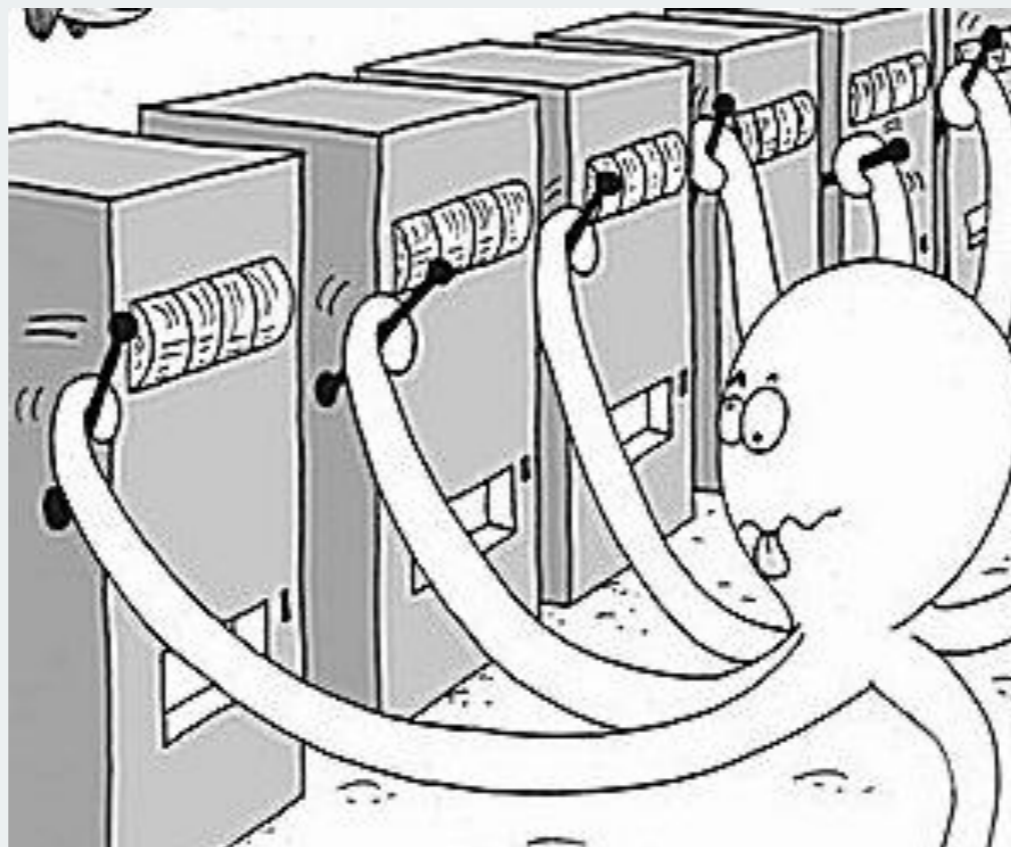
Как вы будете это делать и почему это обучение с подкреплением?



Жадный алгоритм:

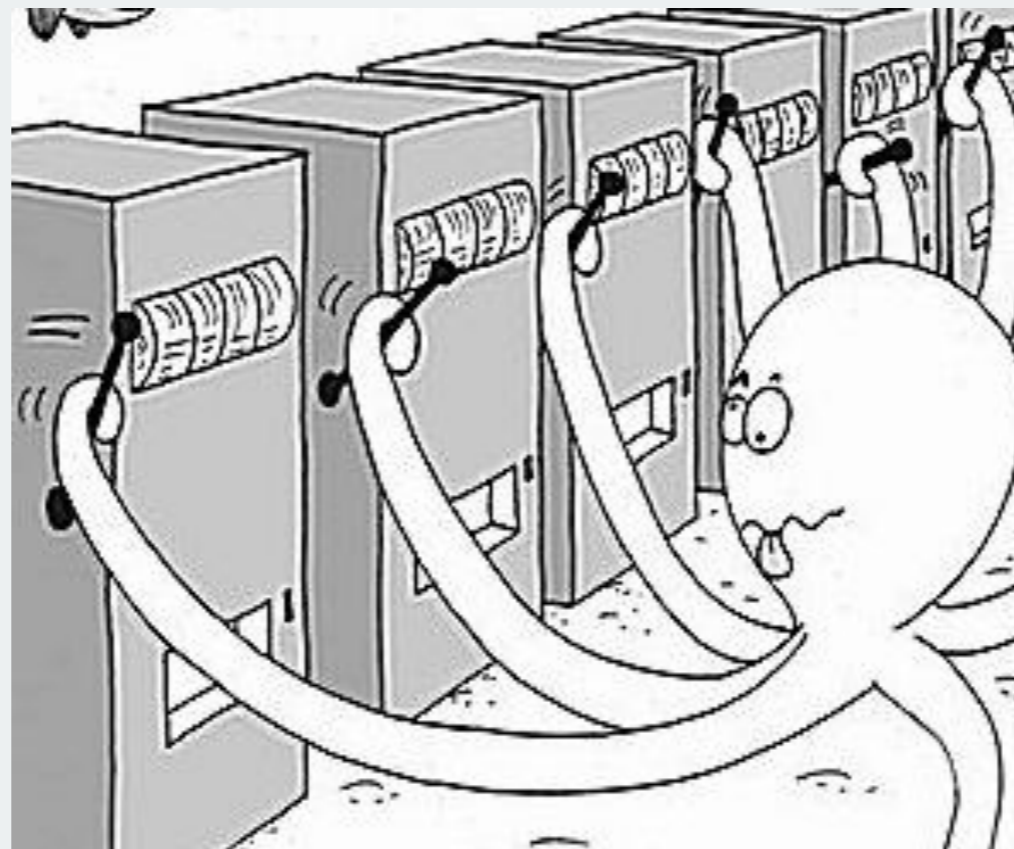
Дернем каждую ручку n раз, а потом будем дергать только ту, у которой средний выигрыш после очередного шага является наибольшим.

Что с этой стратегией не так?



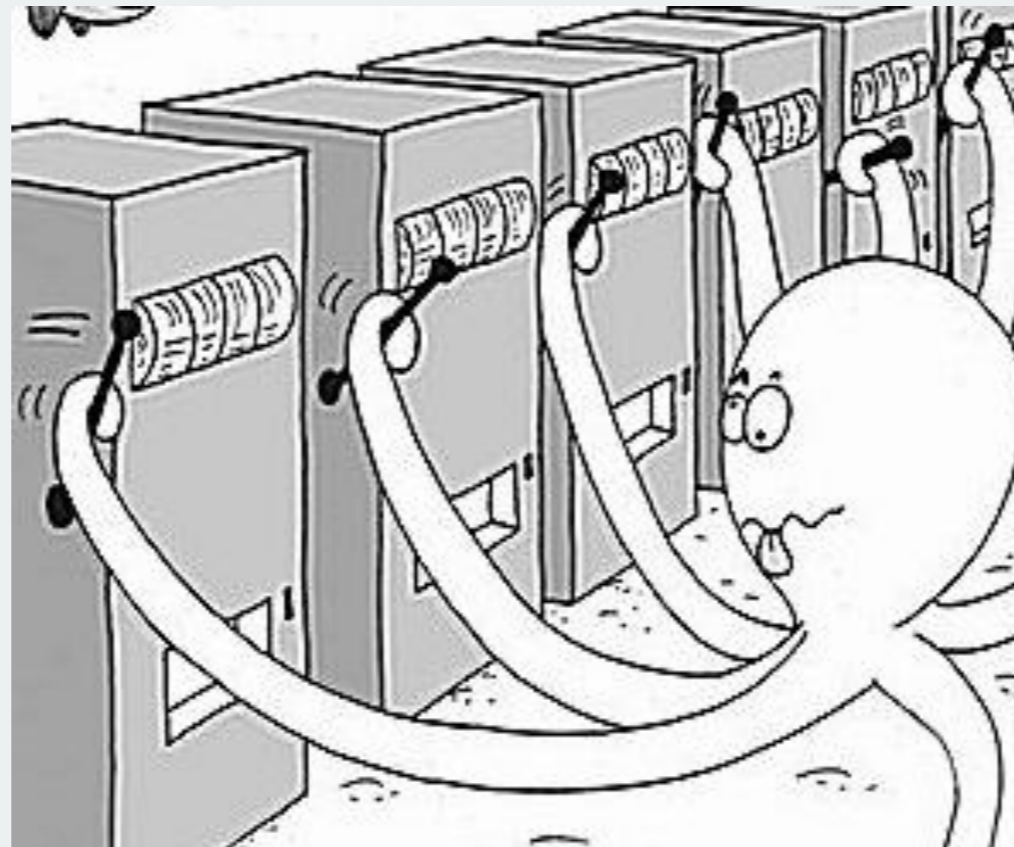
Эпсилон-жадный алгоритм:

Давайте теперь закладывать определенный уровень «неуверенности» в нашу стратегию и, скажем, с вероятностью ε будем совершать случайное действие. Мы гарантируем некоторую вероятность дернуть каждую ручку, и, значит, даже промахнувшись вначале, в конце концов найдем лучшую



Оптимистично-жадный алгоритм:

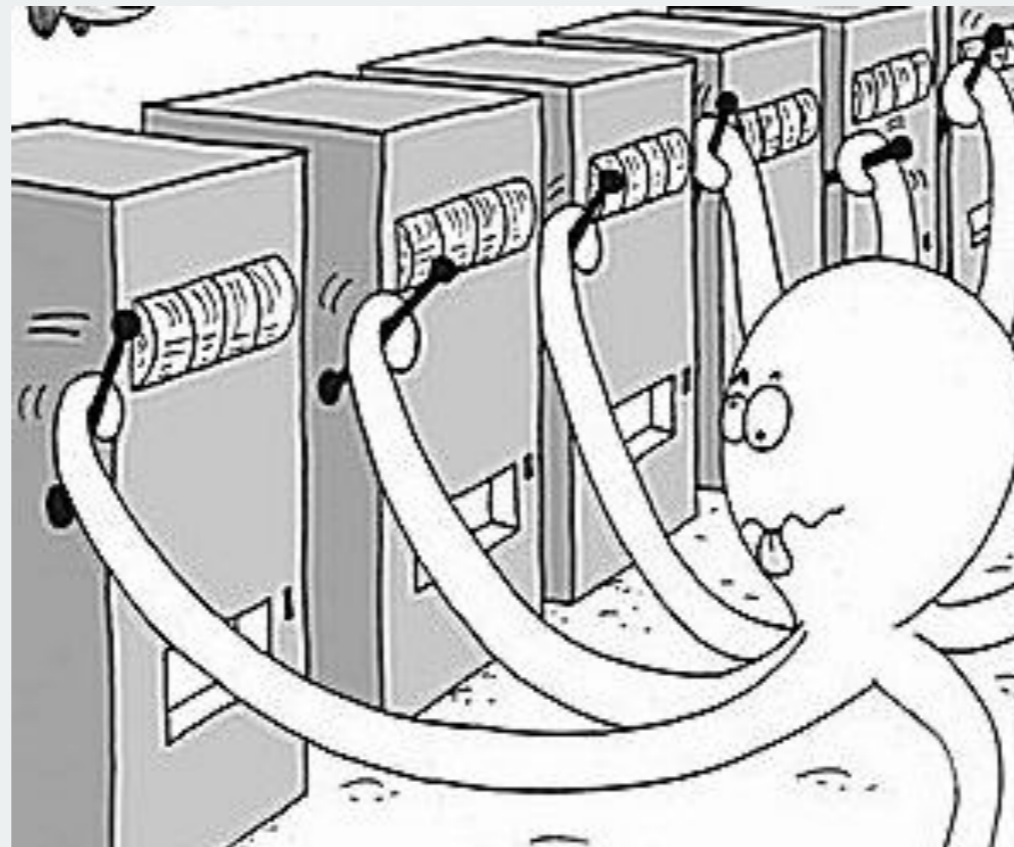
Давайте оценивать не только среднее вознаграждение, но и доверительный интервал для каждой ручки, а выбирать ручки в соответствии с верхней границей доверительного интервала. Тогда, изначально мы оптимистично оцениваем все ручки и постепенно улучшаем оценку.



Мы не будем вдаваться в подробности алгоритмов, но в большинстве своем разные алгоритмы сводятся к разным оптимистичным способам оценки ручек.

Можно посмотреть:

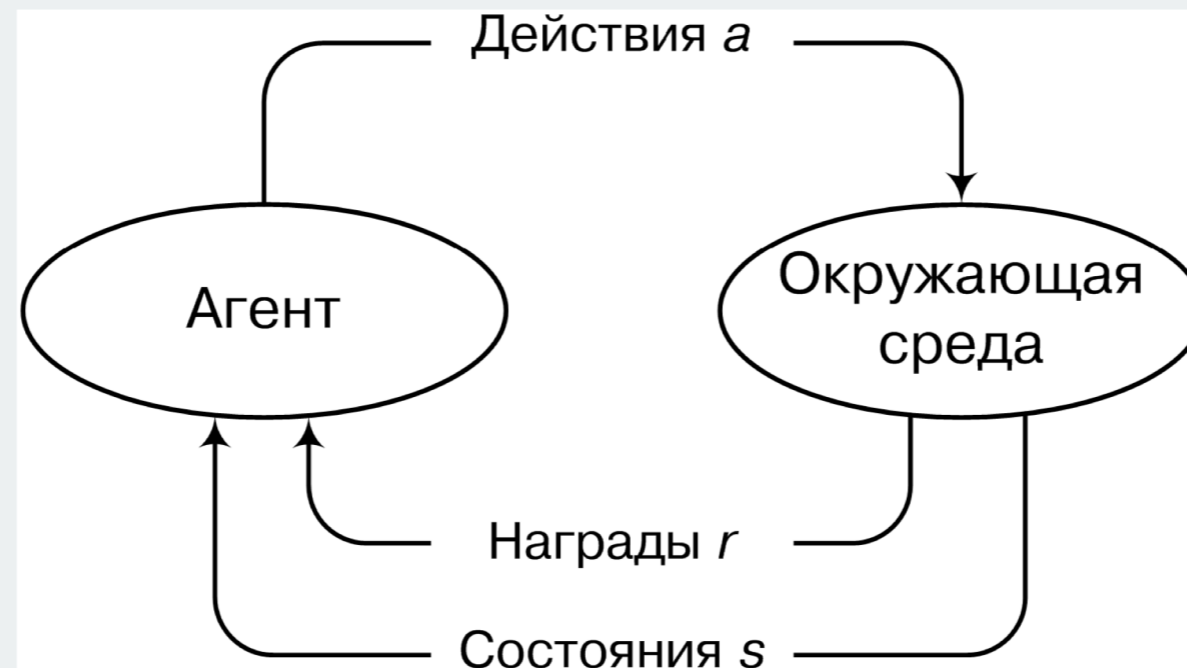
- UCB1 и аналоги
- Thompson sampling



1. Обучение с подкреплением
2. Бандиты
- 3. Стратегии**



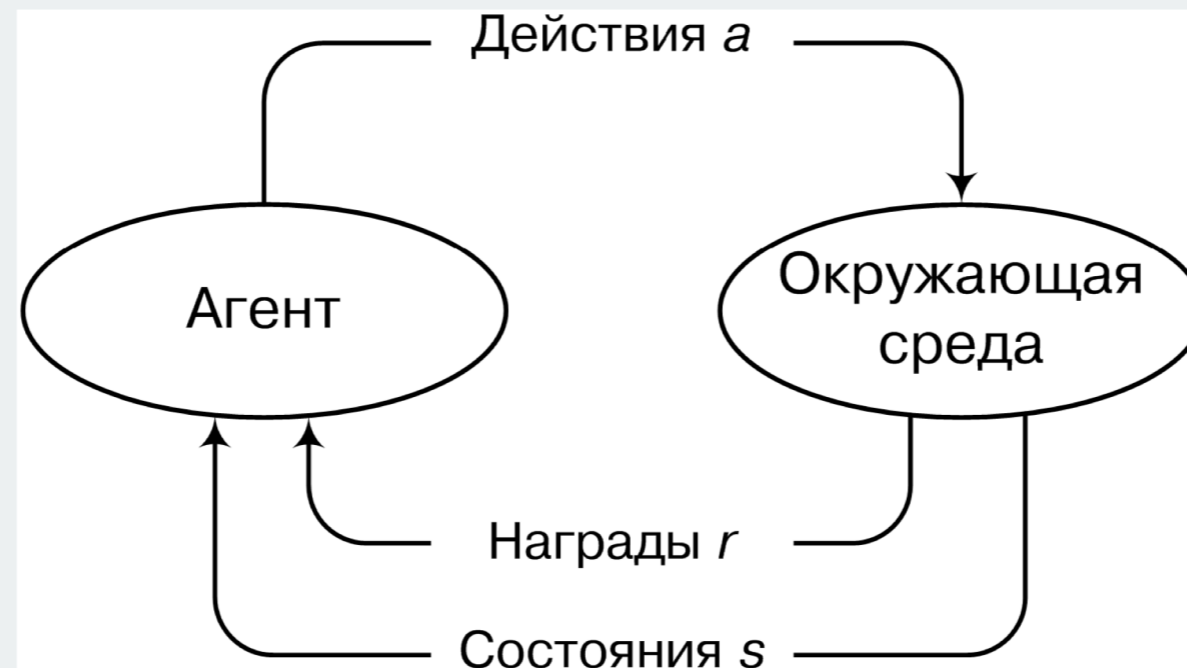
В более общем случае мы хотим уметь оценивать состояния, действия в зависимости от состояния или и то и другое одновременно.
Как можно оценить текущее состояние?



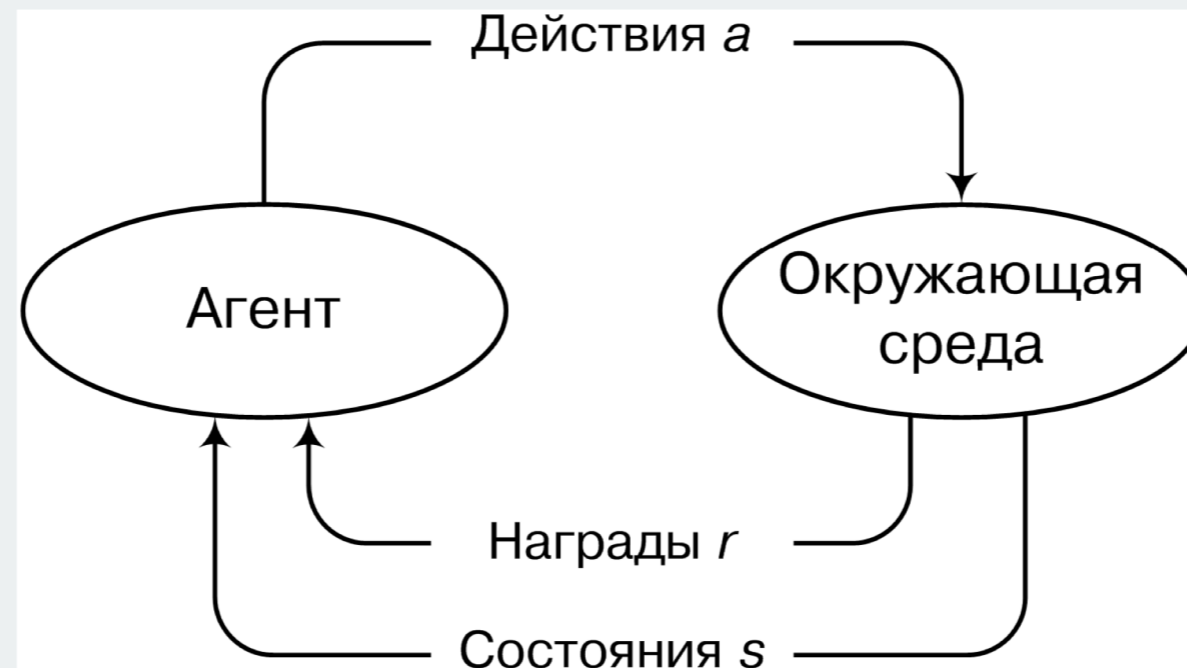
В более общем случае мы хотим уметь оценивать состояния, действия в зависимости от состояния или и то и другое одновременно.

Как можно оценить текущее состояние?

$$V^\pi(s) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s \right]$$

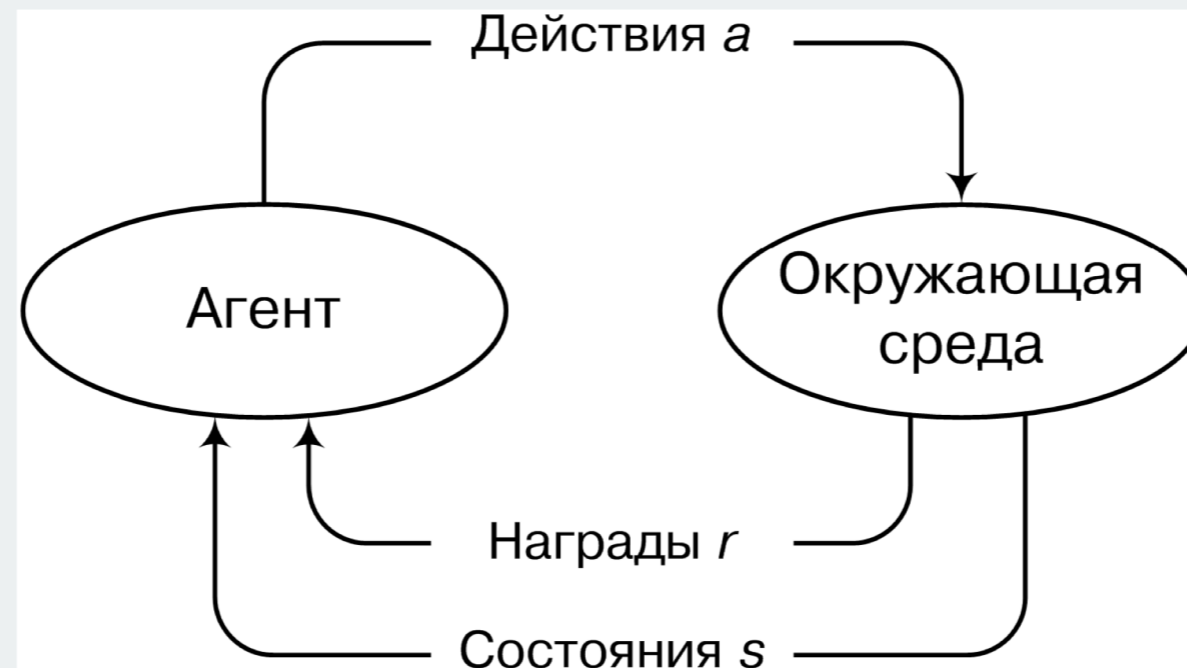


В более общем случае мы хотим уметь оценивать состояния, действия в зависимости от состояния или и то и другое одновременно.
Как можно оценить действие в заданном состоянии?



В более общем случае мы хотим уметь оценивать состояния, действия в зависимости от состояния или и то и другое одновременно.
Как можно оценить действие в заданном состоянии?

$$Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a \right]$$



Вообще, стратегия π — это функция которая задает распределение вероятностей на множестве действий в зависимости от состояния. Все три функции — функцию оценки состояния V , функцию оценки действия Q и стратегию можно связать простым уравнением:

$$V^\pi(s) = \sum_{a \in A} \pi(s, a) * Q^\pi(s, a)$$



Вообще, стратегия π — это функция которая задает распределение вероятностей на множестве действий в зависимости от состояния. Все три функции — функцию оценки состояния V , функцию оценки действия Q и стратегию можно связать простым уравнением:

$$V^\pi(s) = \sum_{a \in A} \pi(s, a) * Q^\pi(s, a)$$

При этом, Q мы тоже можем оценить с помощью V :

$$Q^\pi(s, a) = \sum_{s' \in S} P_{ss'}^a \left(R_{ss'}^a + \gamma V^\pi(s') \right)$$



Вообще, стратегия π — это функция которая задает распределение вероятностей на множестве действий в зависимости от состояния. Все три функции — функцию оценки состояния V , функцию оценки действия Q и стратегию можно связать простым уравнением:

$$V^\pi(s) = \sum_{a \in A} \pi(s, a) * Q^\pi(s, a)$$

При этом, Q мы тоже можем оценить с помощью V :

$$Q^\pi(s, a) = \sum_{s' \in S} P_{ss'}^a \left(R_{ss'}^a + \gamma V^\pi(s') \right)$$

Понятно, что подставляя первое уравнение во второе или наоборот, мы можем избавиться от V или от Q соответственно:

$$V^\pi(s) = \sum_{a \in A} \pi(s, a) * \sum_{s' \in S} P_{ss'}^a \left(R_{ss'}^a + \gamma V^\pi(s') \right)$$
$$Q^\pi(s, a) = \sum_{s' \in S} P_{ss'}^a \left(R_{ss'}^a + \gamma \sum_{a' \in A} \pi(s', a') * Q^\pi(s', a') \right)$$



$$V^\pi(s) = \sum_{a \in A} \pi(s, a) * \sum_{s' \in S} P_{ss'}^a \left(R_{ss'}^a + \gamma V^\pi(s') \right)$$

$$Q^\pi(s, a) = \sum_{s' \in S} P_{ss'}^a \left(R_{ss'}^a + \gamma \sum_{a' \in A} \pi(s', a') * Q^\pi(s', a') \right)$$

В реальности мы хотим оценивать все эти функций только с целью найти оптимальную стратегию, поэтому, нам совершенно не обязательно уметь делать оценки для совсем любой стратегии. Тогда, для оптимальной стратегии эти формулы превратятся в:

$$V^*(s) = \max_a \sum_{s' \in S} P_{ss'}^a \left(R_{ss'}^a + \gamma V^*(s') \right)$$

$$Q^*(s, a) = \sum_{s' \in S} P_{ss'}^a \left(R_{ss'}^a + \gamma \max_{a'} Q^*(s', a') \right)$$

Осталось только разобраться в том, как же научиться делать такие оценки



Оказывается, есть очень простые итеративные способы. Один из основных современных подходов основан на принципе TD-learning (temporal difference). Суть которого заключается в том, чтобы при каждом очередном переходе немножко улучшать функцию оценки для предыдущего состояния основываясь на полученном вознаграждении и новом состоянии.



Оказывается, есть очень простые итеративные способы. Один из основных современных подходов основан на принципе TD-learning (temporal difference). Суть которого заключается в том, чтобы при каждом очередном переходе немножко улучшать функцию оценки для предыдущего состояния основываясь на полученном вознаграждении и новом состоянии.

Простейший алгоритм TD-обучения можно записать так:

- Инициализировать s ;
- Для каждого шага в эпизоде:
 - выбрать a по стратегии π или по формуле для V^*
 - сделать a и получить награду r и новое состояние s'
 - обновить $V(s) := V(s) + \alpha(r + \gamma V(s') - V(s))$
 - перейти к новому состоянию s'



Оказывается, есть очень простые итеративные способы. Один из основных современных подходов основан на принципе TD-learning (temporal difference). Суть которого заключается в том, чтобы при каждом очередном переходе немножко улучшать функцию оценки для предыдущего состояния основываясь на полученном вознаграждении и новом состоянии.

Простейший алгоритм TD-обучения можно записать так:

- Инициализировать s ;
- Для каждого шага в эпизоде:
 - выбрать a по стратегии π или по формуле для V^*
 - сделать a и получить награду r и новое состояние s'
 - обновить $V(s) := V(s) + \alpha(r + \gamma V(s') - V(s))$
 - перейти к новому состоянию s'

Для функции Q мы можем сделать абсолютно то же самое. Существенное ограничение этого метода заключается в том, что для того чтобы получить математические гарантии сходимости нужно обеспечить возможность встретить любую пару (s, a) бесконечное количество раз.

Как это можно сделать?



1. Обучение с подкреплением основано на возможности и необходимости взаимодействовать с окружающей средой
2. В простейшем случае состояние окружающей среды не меняется
3. В задачах с известным поведением среды можно использовать простейшие итеративные методы.





Спасибо
за внимание!