

Q-learning

Владимир Морозов

Агент

Агент

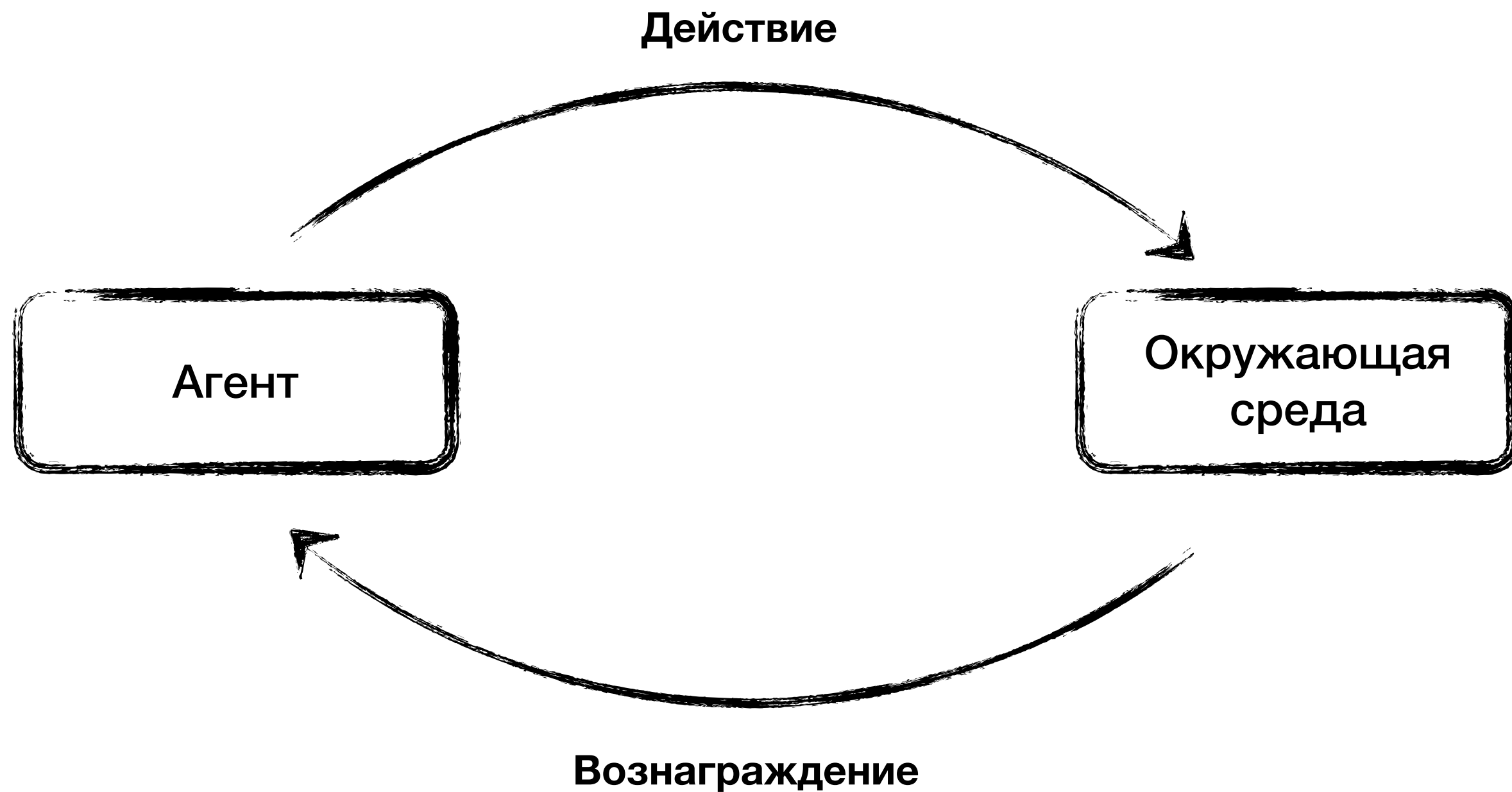
Окружающая
среда

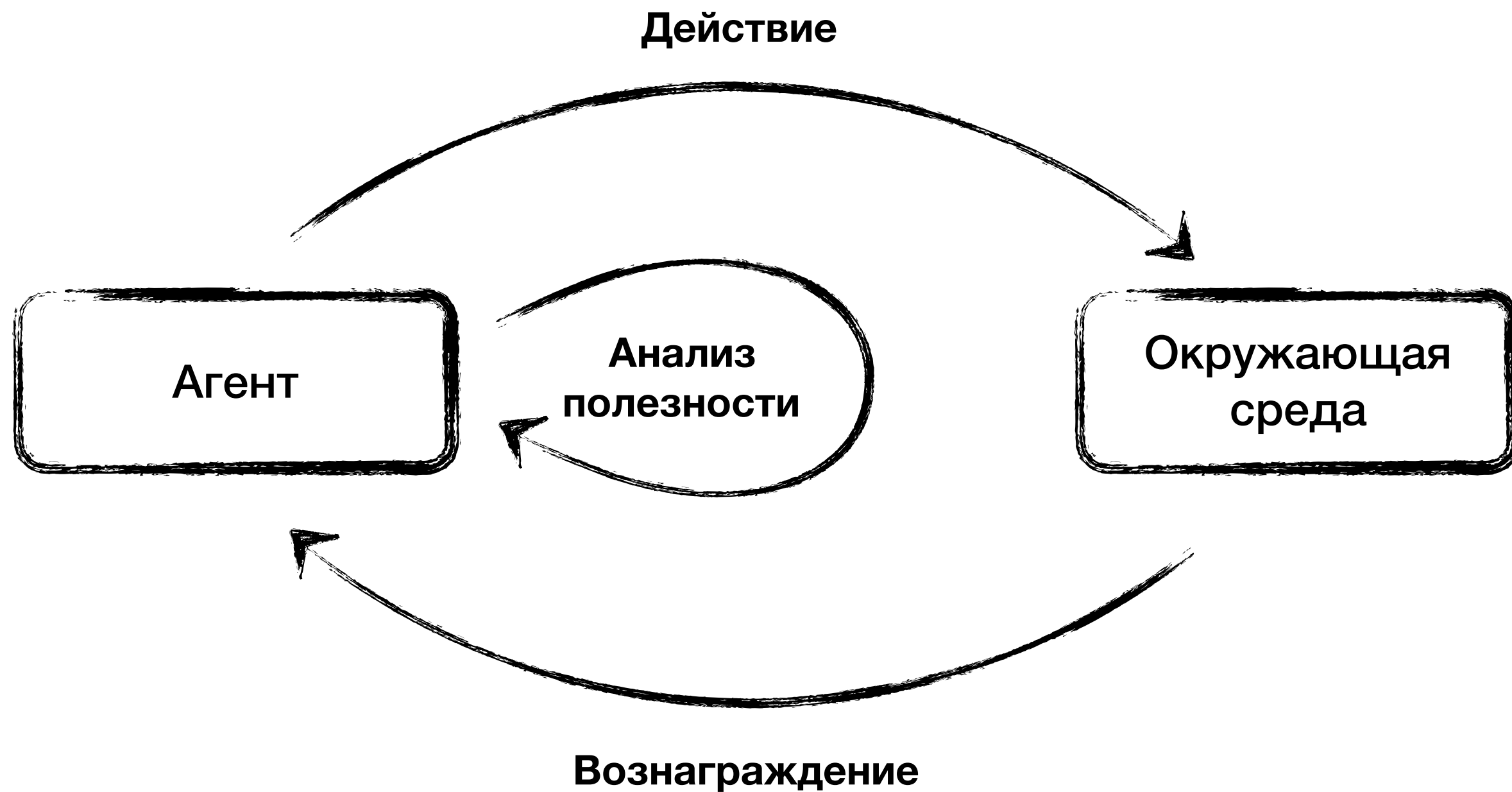
Действие

Агент

**Окружающая
среда**







Алгоритм

1.Инициализация

2.Наблюдение

3.Обновление

4.Выбор действия



Инициализация

Инициализируем функцию полезности Q
нулевыми или случайными значениями:

$$\begin{aligned}\forall i, j: Q[s_i, a_j] &= 0 \\ \forall i, j: Q[s_i, a_j] &= \text{RND}\end{aligned}$$

Пока модель ничего не знает о вознаграждениях окружающей среды.

Наблюдение

$S_{prev} = S$

Сохранить предыдущее состояние системы

$A_{prev} = A$

Сохранить предыдущее действие

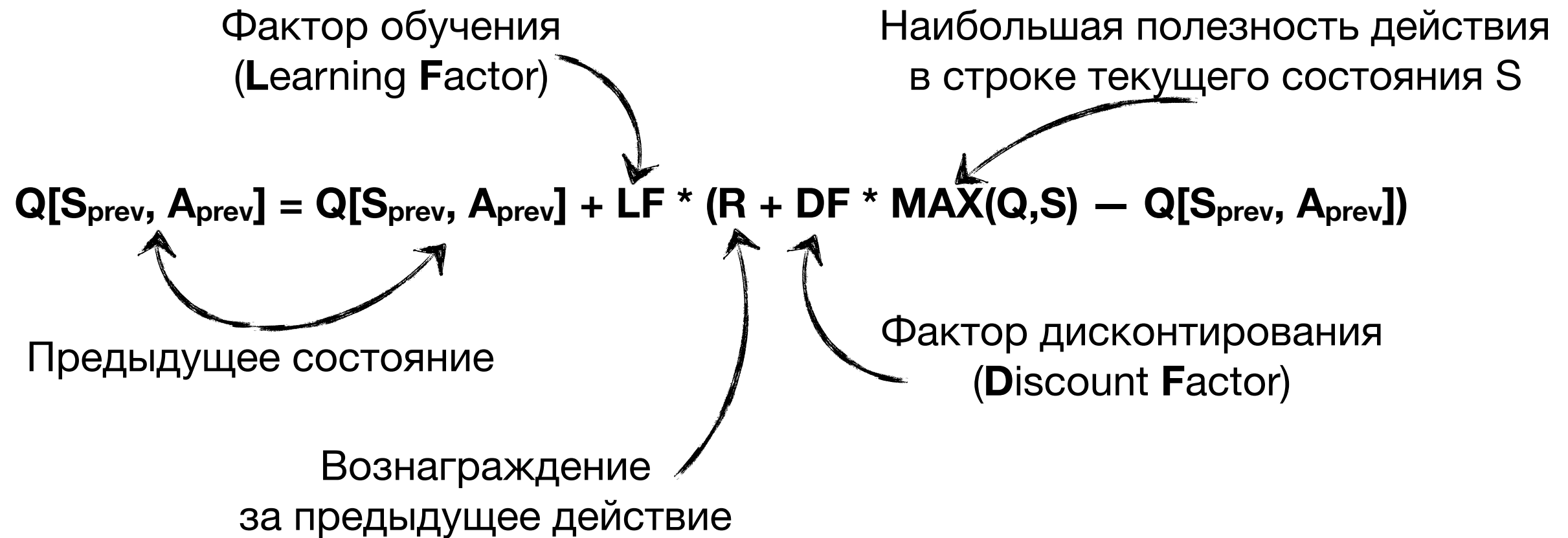
$S = \text{OBSERVE}$

Получить текущее состояние (например, с датчиков)

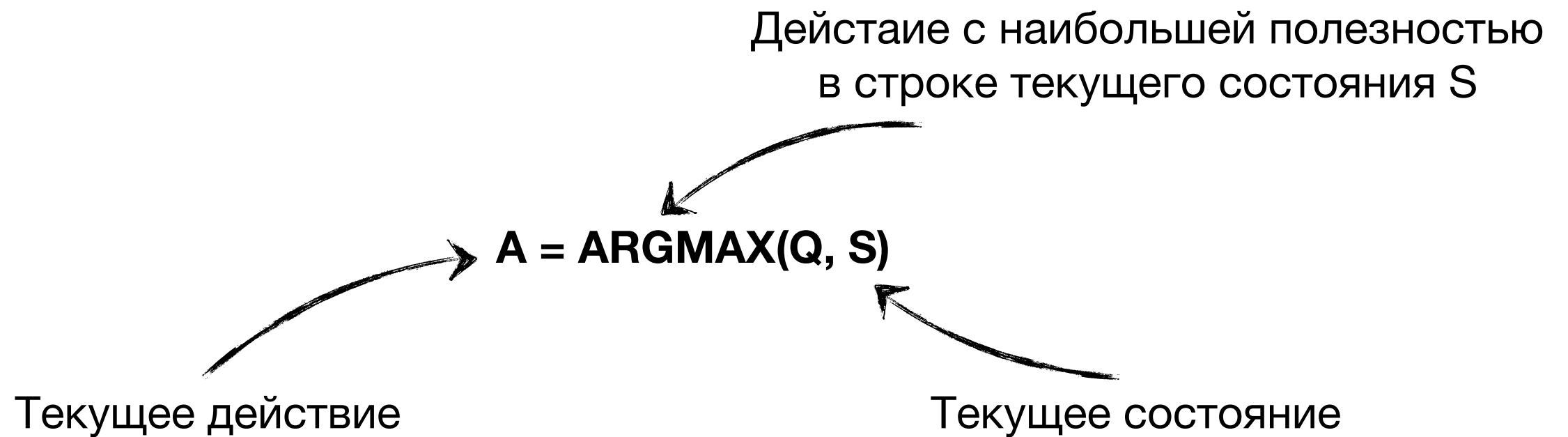
$R = \text{FROM_ENV}$

Получить награду за предыдущее действие

Обновление



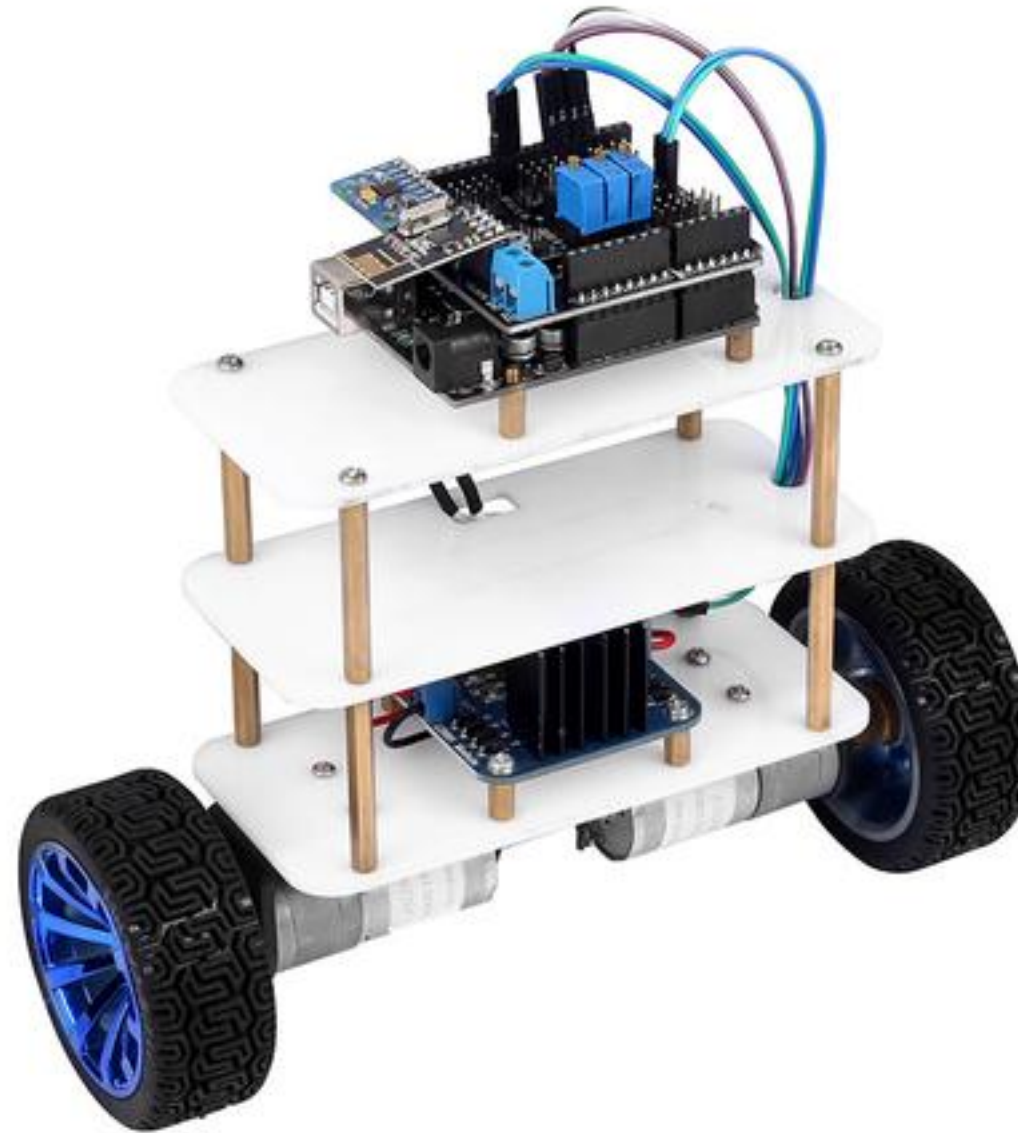
Действие



Представление Q

	Действие 1	Действие 2	...	Действие N
Состояние 1	Q_{11}	Q_{12}	...	Q_{1N}
Состояние 2	Q_{21}	Q_{22}	...	Q_{2N}
...
Состояние M	Q_{M1}	Q_{M2}	...	Q_{MN}

Пример



Самобалансирующий двухколёсный робот

Implementation of Q Learning and Deep Q Network For Controlling a Self
Balancing Robot Model

MD Muhaimin Rahman^{1,a}, SM Hasanur Rashid^{2,a} and M.M Hossain ^b

Матрица Q

	-200 рад/с	-100 рад/с	-50 рад/с	...	+200 рад/с
Континуум состояний [-10, 10]	0	0	0	...	0
	0	0	0	...	0

	0	0	0	...	0

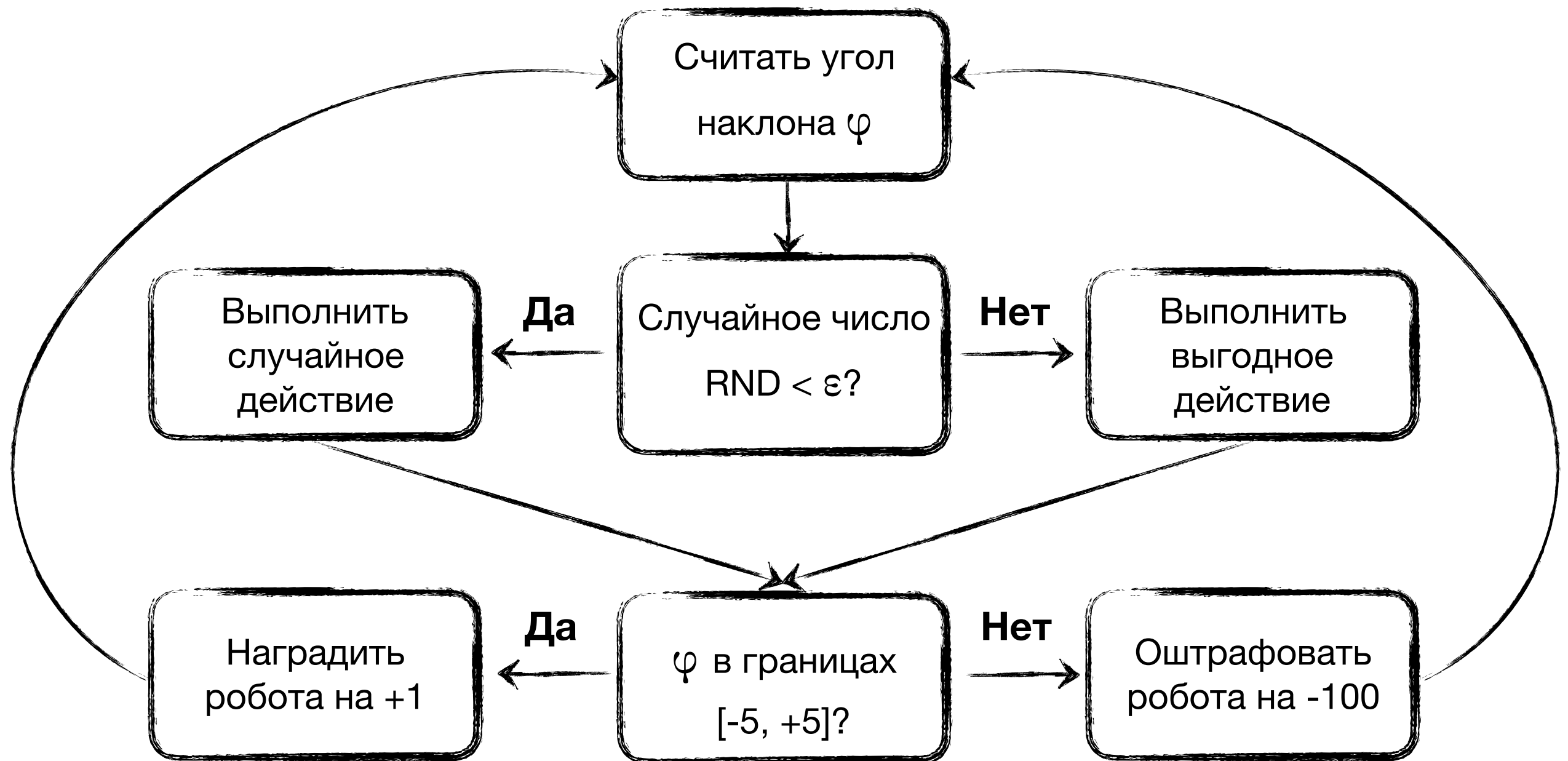
Матрица Q

	-200 рад/с	-100 рад/с	-50 рад/с	...	+200 рад/с
-10°	0	0	0	...	0
-9°	0	0	0	...	0
...
+10°	0	0	0	...	0

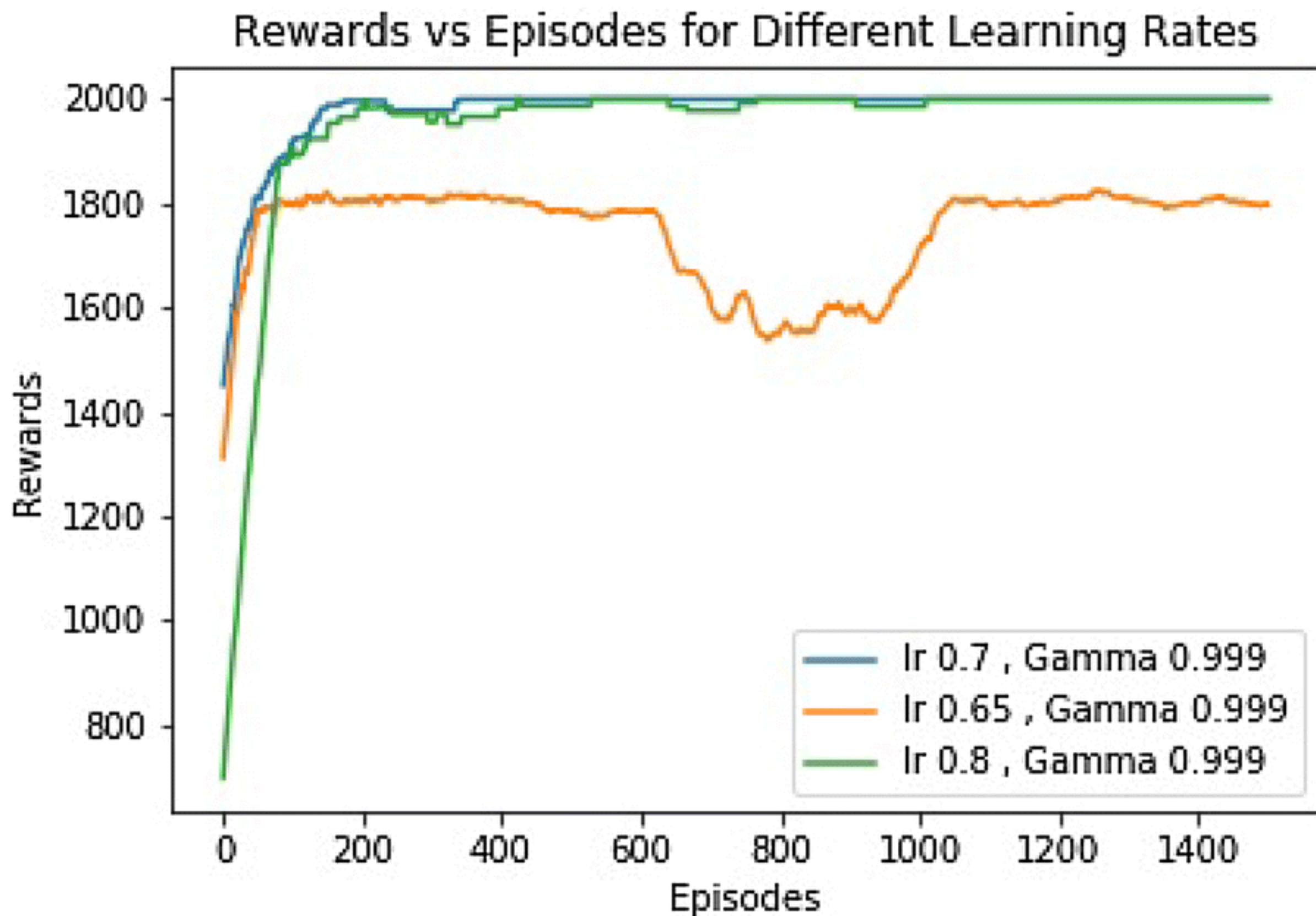
Обучение

Вознаграждения:

- Выход за границы $[-5, +5]$: -100
- Нахождение в интервале $[-5, +5]$: $+1$



Результаты



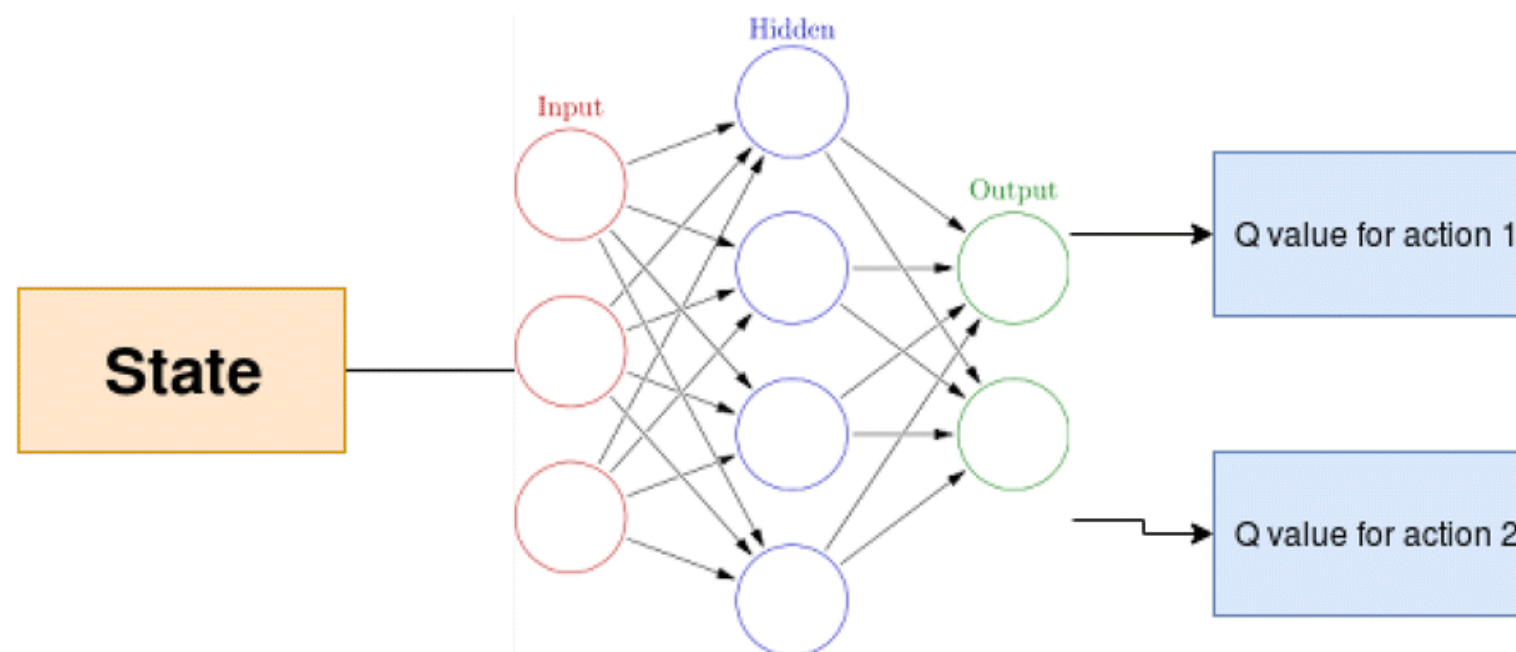
Модификации

"Experience replay"

Сохранение на каждом шаге текущих состояния, награды, действия. При новой инициализации данные не заполняются нулями, а берутся случайным образом из сохранённых

Предсказание значений Q

Вместо постоянного пересчёта значений качества их можно предсказывать по текущему состоянию.



**Спасибо за
внимание!**