

Московский государственный технический университет им. Н.Э. Баумана
Факультет «Информатика и системы управления»
Кафедра «Системы обработки информации и управления»



Рубежный контроль №2
по дисциплине
«Методы машинного обучения»

Выполнил:
Студент группы ИУ5-21М
Королёва Ю.П.

Проверил:
Гапанюк Ю.Е.

Москва, 2023

Для одного из алгоритмов временных различий, реализованных Вами в соответствующей лабораторной работе:

- SARSA
- Q-обучение
- Двойное Q-обучение

осуществите подбор гиперпараметров. Критерием оптимизации должна являться суммарная награда.

Меняя $\text{EPSILON}=0.05$, посмотрим на графики.

```
✓ [12] 3 сек.
*****
epsilon: 0.1
*****

method: sarsa

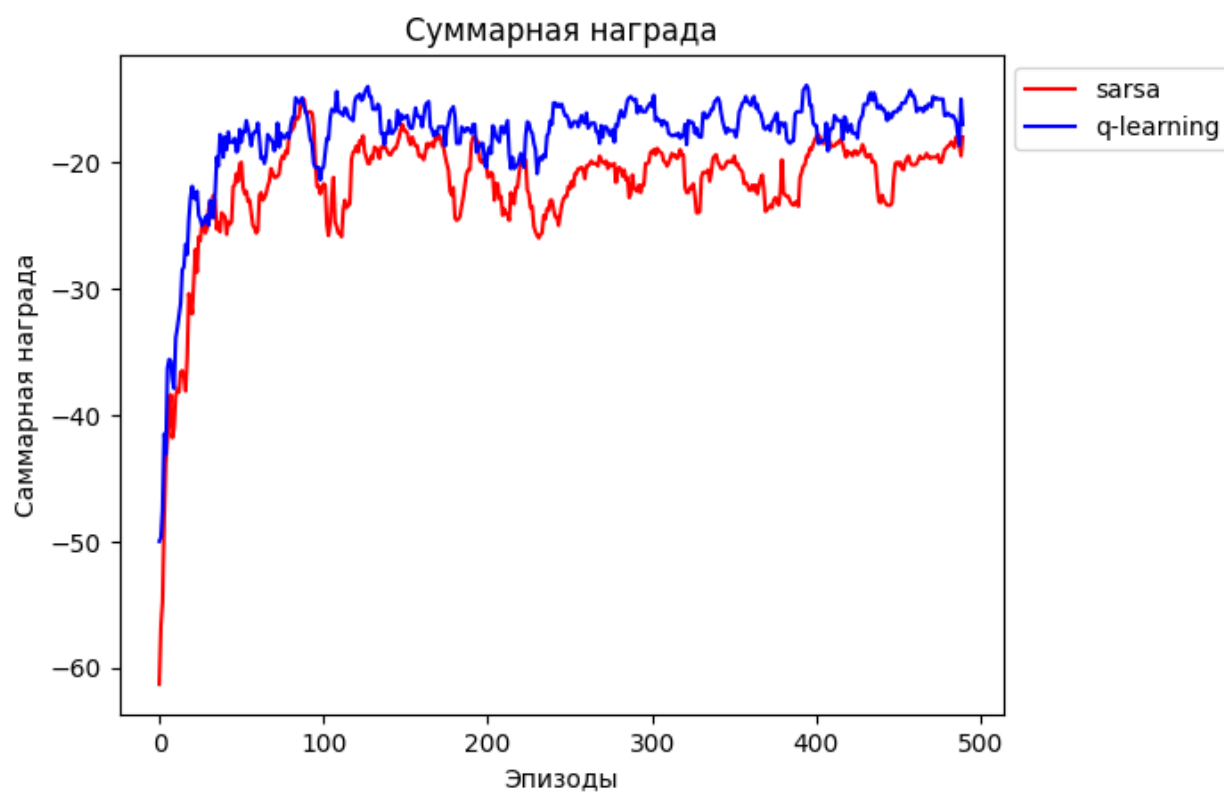
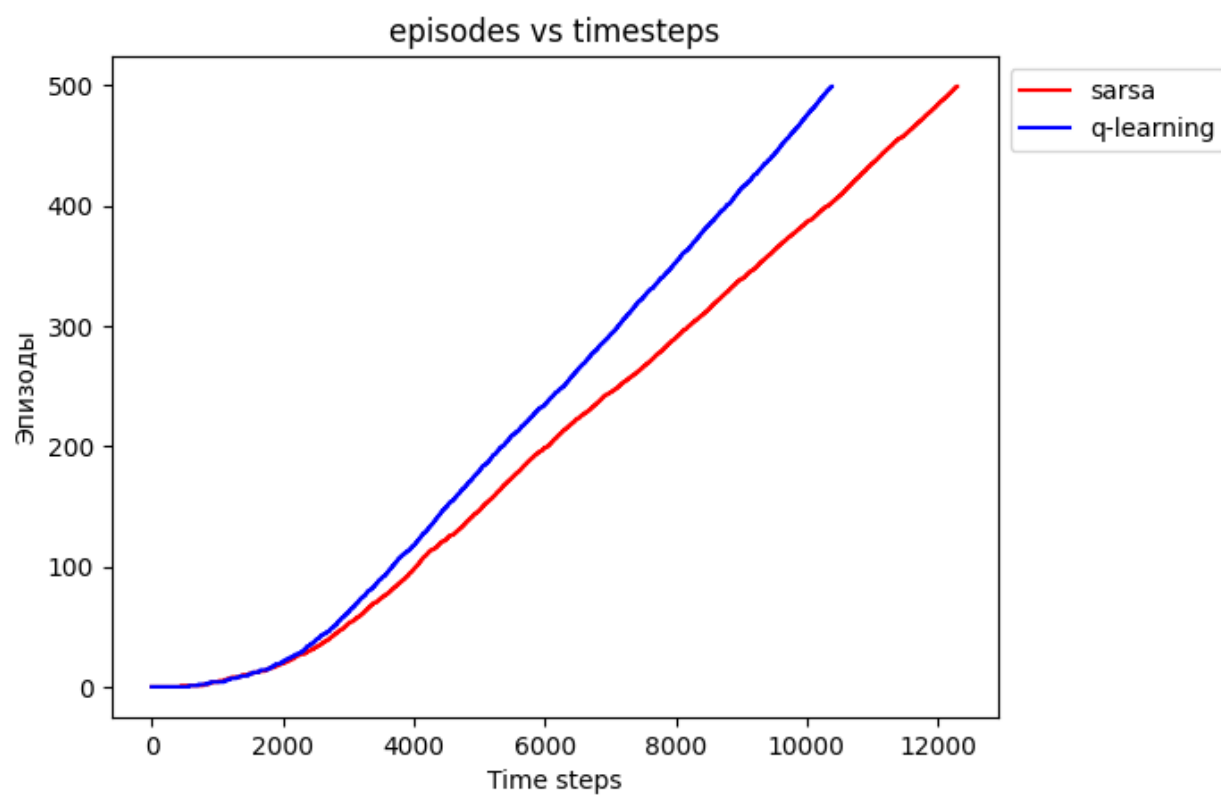
>>>> Траектория на предыдущем эпизоде
*      →      →      →      →      →      →      →      →      ↓
→      ↑      *      *      *      *      *      *      *      ↓
↑      *      *      *      *      *      *      *      *      ↓
↑      █      █      █      █      █      █      █      █      G

>>>>Оптимальная политика
→      →      →      →      →      →      →      →      →      ↓
→      ↑      →      →      →      →      →      →      →      ↓
↑      ↑      →      →      →      →      →      ↑      ↑      ↓
↑      █      █      █      █      █      █      █      █      G

method: q-learning

>>>> Траектория на предыдущем эпизоде
*      *      *      *      *      *      *      *      *      *
*      *      *      *      *      *      *      *      →      ↓
→      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      G

>>>>Оптимальная политика
↓      →      →      →      →      →      →      ↓      →      ↓
↓      ↓      ↓      →      →      →      →      ↓      ↓      ↓
→      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      G
```



epsilon: 0.2

method: sarsa

>>> Траектория на предыдущем эпизоде

```
*      *      *      *      *      *      *      *      *      *      *
*      →      →      →      →      →      →      →      →      →      ↓      *
→      ↑      *      *      *      *      *      *      *      *      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

>>>>Оптимальная политика

```
→      →      →      →      ↓      ↓      ↓      →      →      ↓      →      ↓
↑      →      →      →      →      →      →      →      →      →      ↓      ↓
→      ↑      →      →      ↑      →      →      →      →      ↑      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

method: q-learning

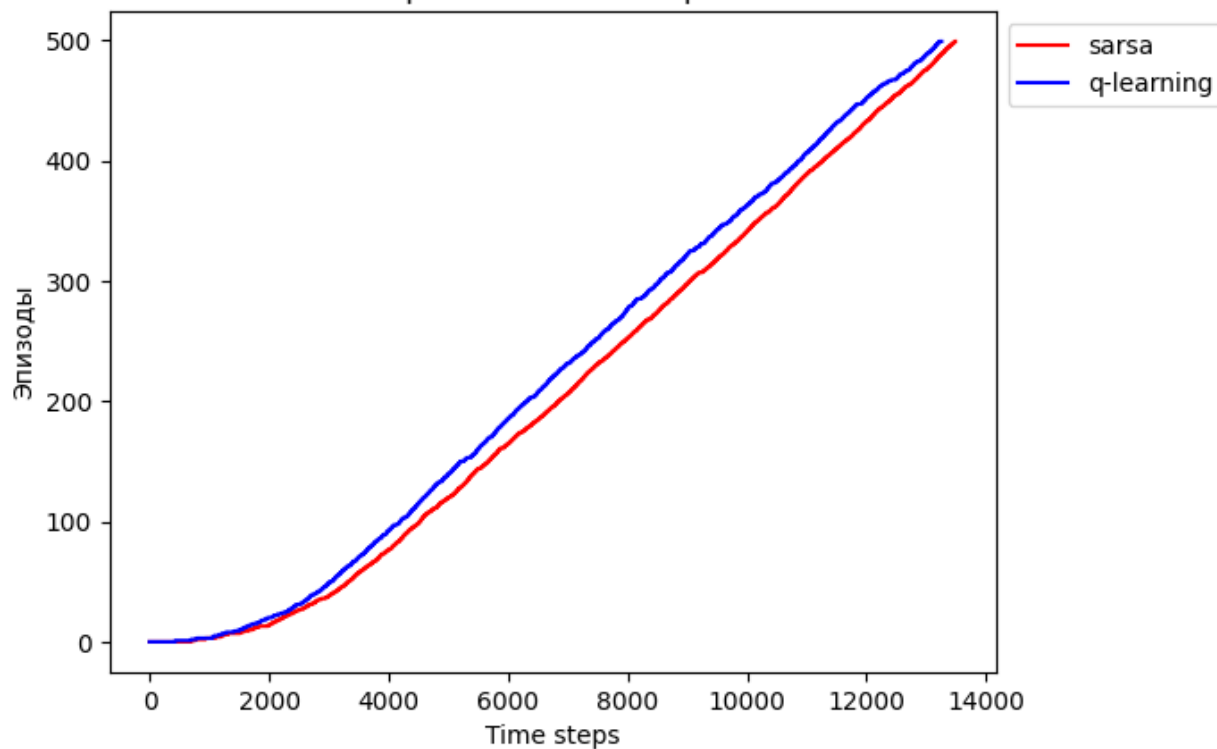
>>> Траектория на предыдущем эпизоде

```
*      *      *      *      *      *      *      *      *      *      *
*      *      ↓      *      *      *      *      *      *      *      ↓
→      →      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

>>>>Оптимальная политика

```
→      →      ↓      →      →      →      →      ↓      →      ↓      ↓      ↓
↓      →      →      →      →      →      →      →      ↓      →      →      ↓
→      →      →      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

episodes vs timesteps





```
✓ [14] epsilon: 0.3
1
сек.
```

method: sarsa

>>>> Траектория на предыдущем эпизоде

```
→      →      →      →      →      →      →      →      →      →      →      ↓
↑      *      *      *      *      *      *      *      ↑      *      ↑      ↓
↑      *      *      *      *      *      *      *      *      *      *      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

>>>>Оптимальная политика

```
→      →      →      →      →      →      →      →      →      →      ↓      ↓
↑      ↑      →      →      →      →      →      →      ↑      ↑      →      ↓
↑      ↑      ↑      ↑      →      →      ↑      →      ↑      →      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

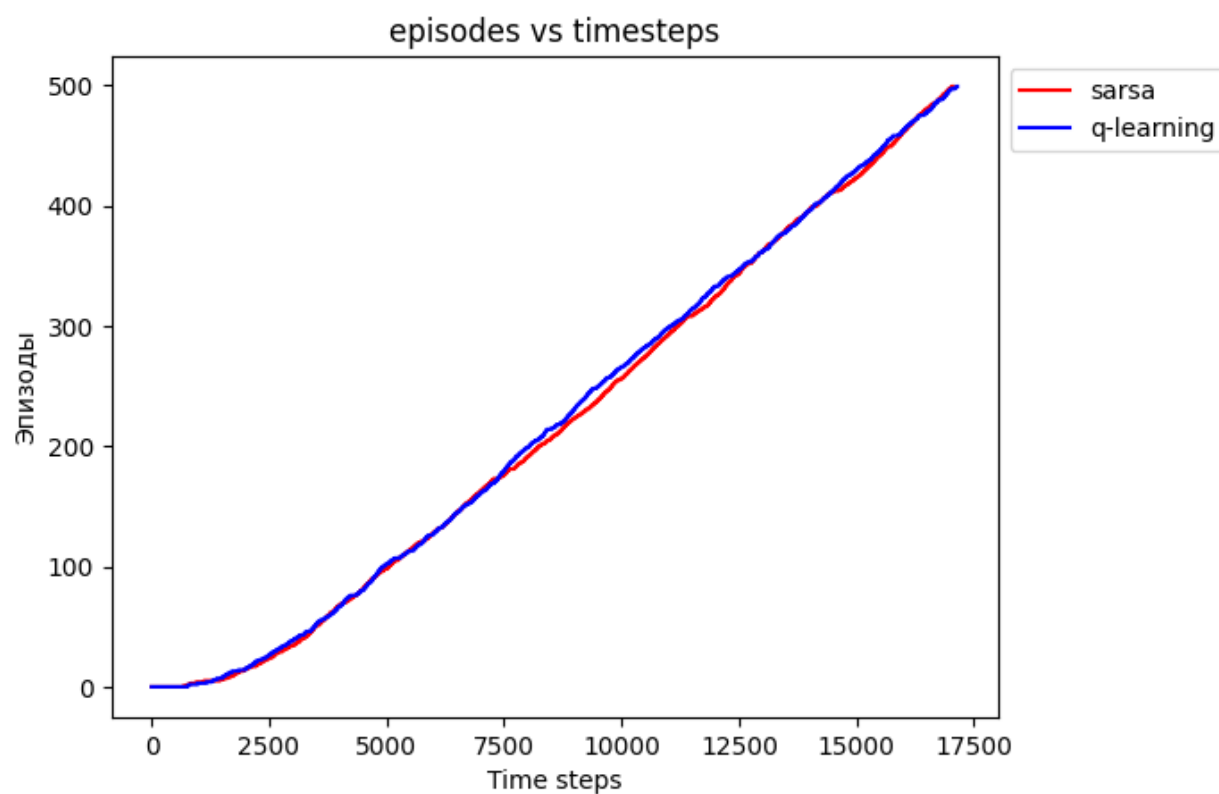
method: q-learning

>>>> Траектория на предыдущем эпизоде

```
*      *      *      *      *      *      *      *      *      *      *      *
*      *      *      ↓      *      *      *      *      *      *      *      *
→      →      →      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

>>>>Оптимальная политика

```
→      ↓      →      →      →      →      →      →      →      →      →      ↓
↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓
→      →      →      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```



✓ [15] epsilon: 0.4
2
сек.

method: sarsa

>>>> Трaектория на предыдущем эпизоде

```
→      →      →      →      →      →      →      →      →      →      →      ↓
↑      *      *      *      *      *      *      *      *      *      *      ↓
↑      *      *      *      *      *      *      *      *      *      *      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

>>>>Оптимальная политика

```
→      →      →      →      →      →      →      →      →      →      ↓      ↓
→      →      →      ↑      →      →      →      ↑      →      →      →      ↓
↑      →      ↑      ↑      →      ↑      ↑      ↑      ↑      ↑      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

method: q-learning

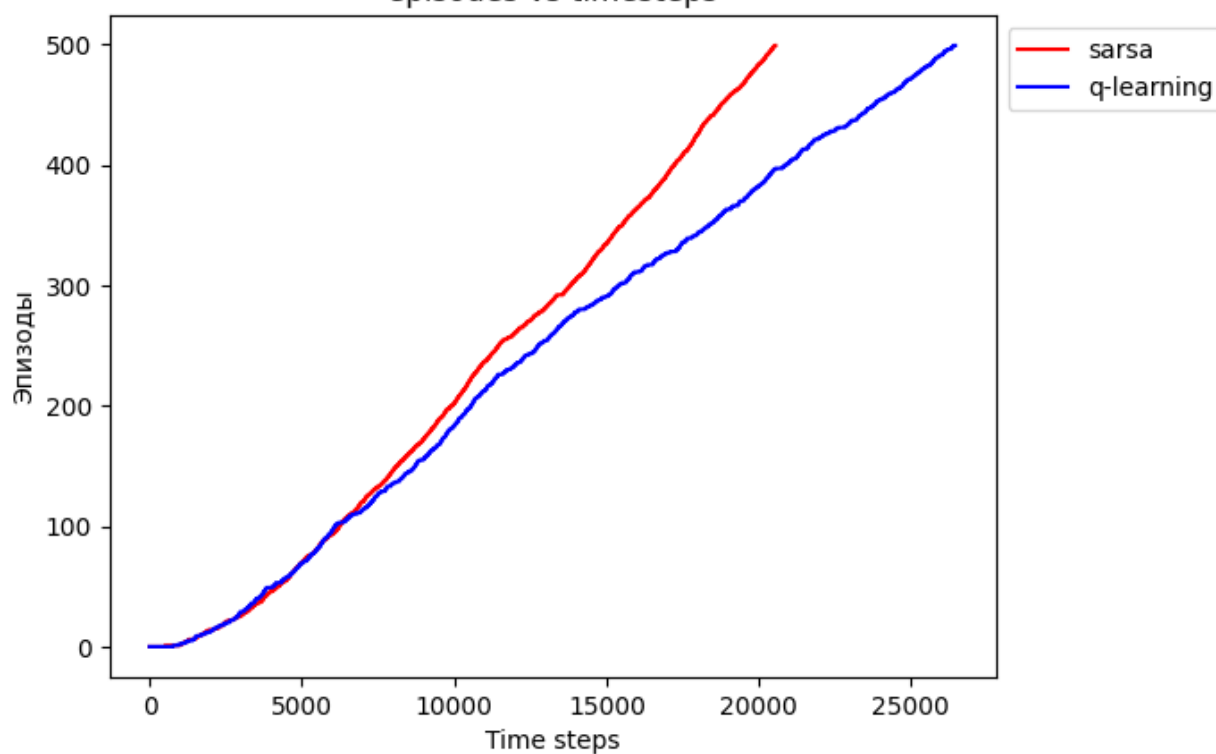
>>>> Трaектория на предыдущем эпизоде

```
→      ↓      *      →      →      →      →      →      →      →      →      ↓
↑      ↓      *      *      ↑      *      *      *      *      *      *      ↓
↑      →      →      →      ↑      ↓      *      *      *      *      *      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

>>>>Оптимальная политика

```
↓      →      →      ↓      →      →      ↓      ↓      →      ↓      ↓      ↓
↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓
→      →      →      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

episodes vs timesteps





✓ [16] *****
5 epsilon: 0.5
сек. *****

method: sarsa

>>>> Траектория на предыдущем эпизоде

```

→      →      ↓      ←      →      →      →      →      ↓      *      ↓
→      →      →      →      ↑      ←      *      ↑      ←      →      →      ↓
↑      *      *      *      ↓      *      *      *      *      *      *      ↓
↑      █      █      █      █      █      █      █      █      █      █      G

```

>>>>Оптимальная политика

```

↑      →      →      →      →      →      →      →      →      →      →      ↑
↑      →      →      →      ↑      ↑      →      ↑      →      →      →      ↓
↑      →      →      →      ↑      →      ↑      ↑      ↑      ↑      ↑      ↓
↑      █      █      █      █      █      █      █      █      █      █      G

```

method: q-learning

>>>> Траектория на предыдущем эпизоде

```

↓      →      →      ↓      ↓      *      *      ↓      *      ↓      ←      ←
↓      ↑      ↓      ↓      ↓      *      ↓      ←      ←      →      →      ↓
→      →      →      →      →      →      →      →      →      →      ↑      ↓
↑      █      █      █      █      █      █      █      █      █      █      G

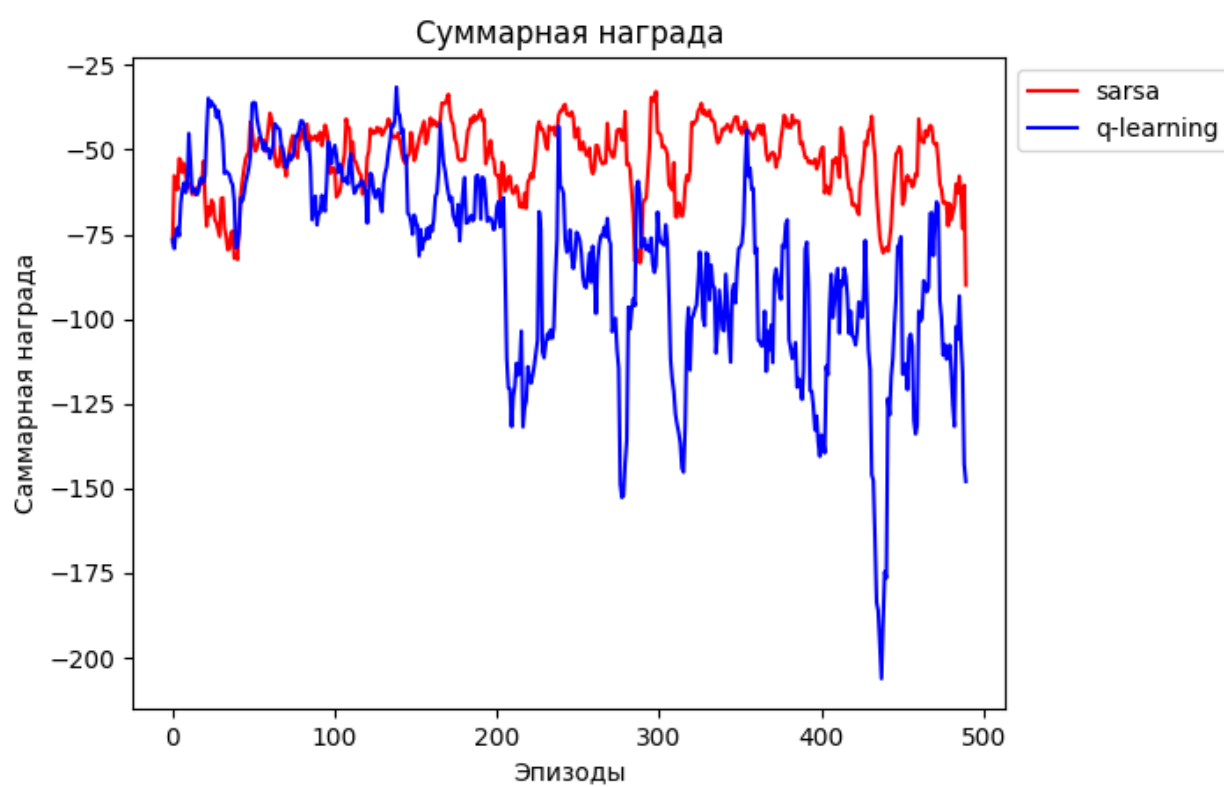
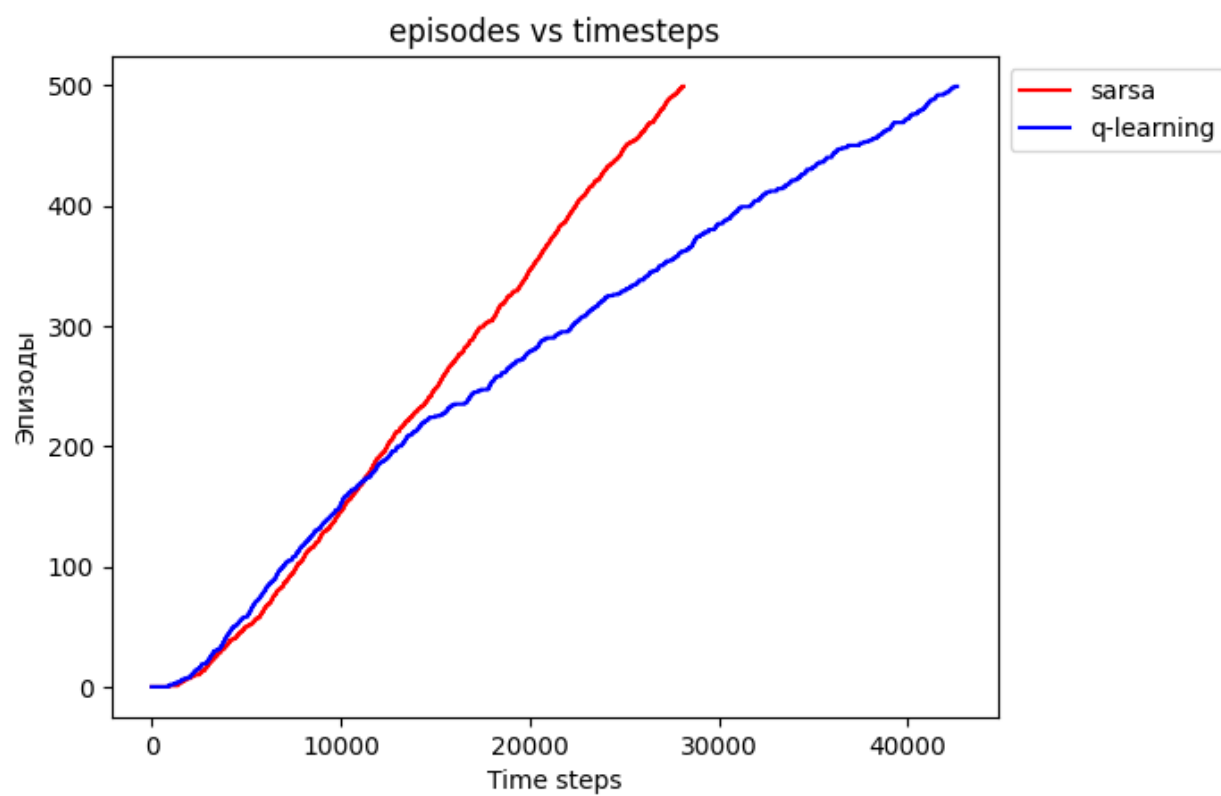
```

>>>>Оптимальная политика

```

↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓
↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓
→      →      →      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G

```

epsilon: 0.6

method: sarsa

>>>> Трактектория на предыдущем эпизоде

[illegible]

>>>>Оптимальная политика

[illegible]

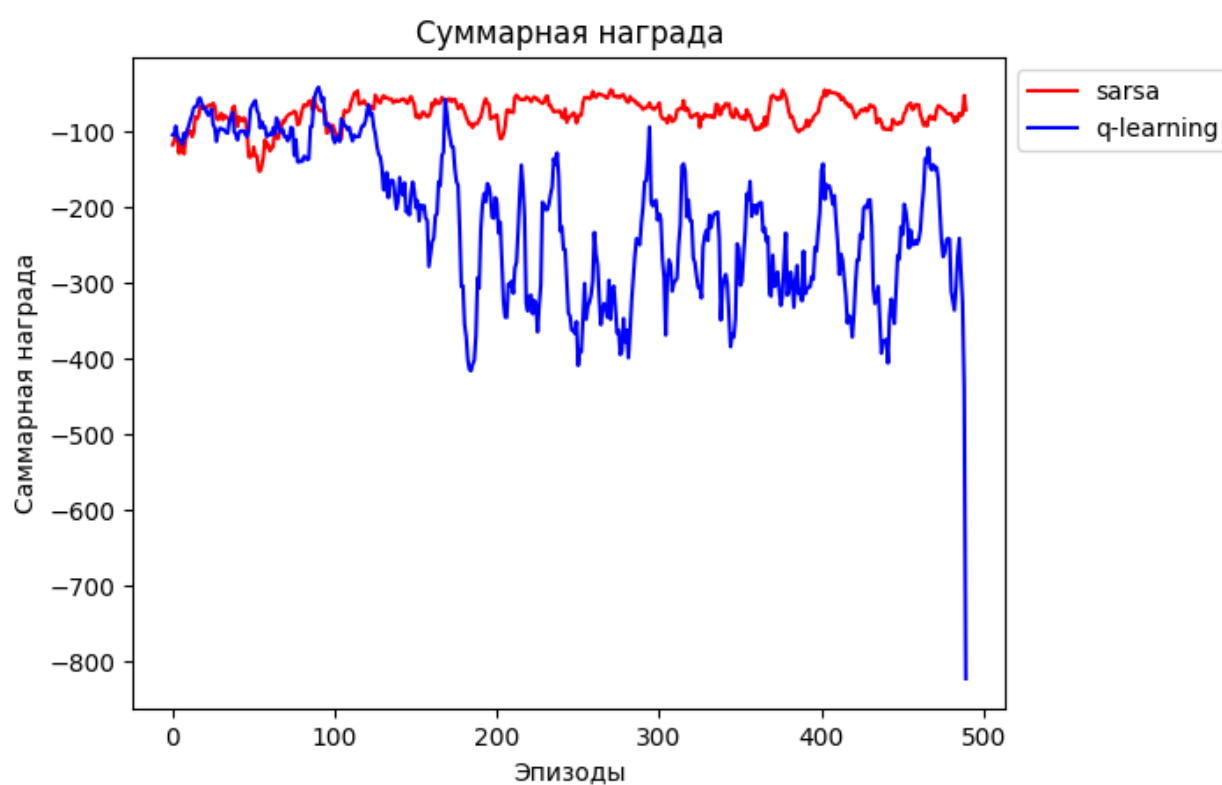
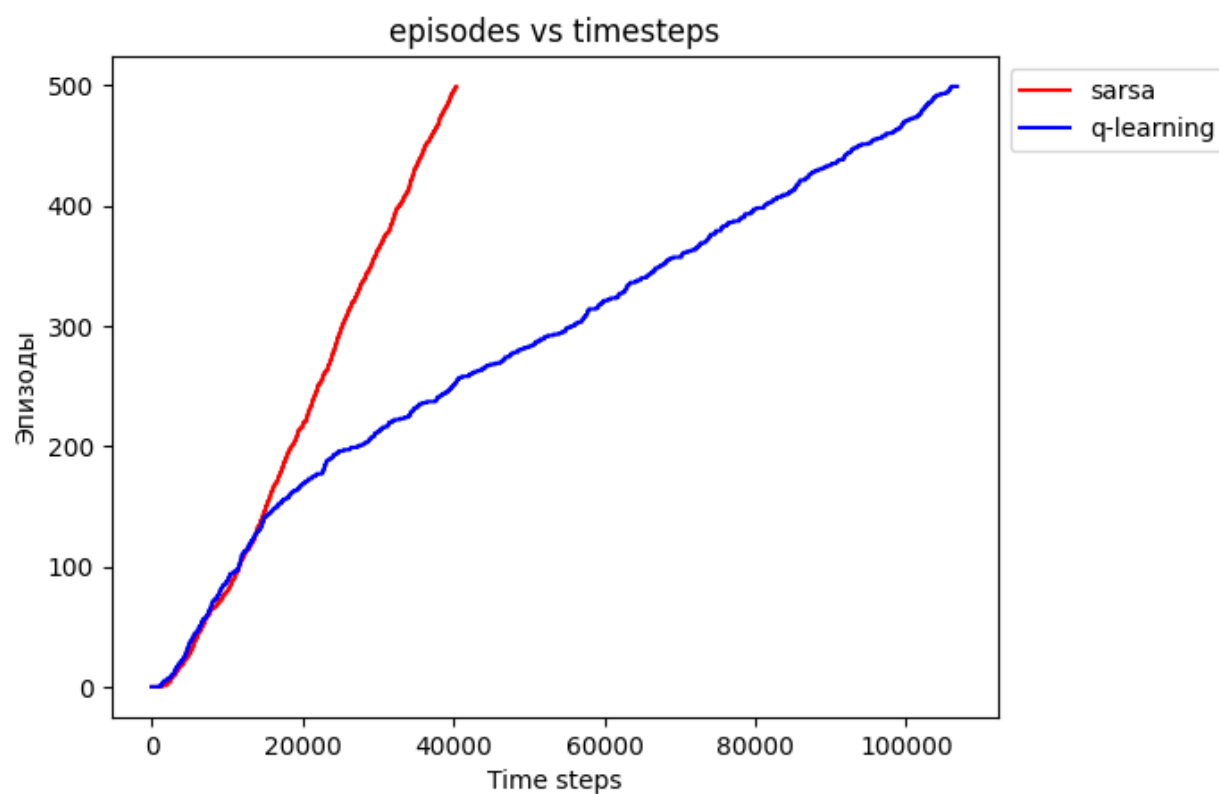
method: q-learning

>>>> Трактектория на предыдущем эпизоде

[illegible]

>>>>Оптимальная политика

Diagram illustrating a 10x10 grid structure, likely representing a 1000x1000 matrix. The grid is divided into four quadrants. The top-left quadrant (rows 1-5, columns 1-5) is white. The top-right quadrant (rows 1-5, columns 6-10) is light gray. The bottom-left quadrant (rows 6-10, columns 1-5) is dark gray. The bottom-right quadrant (rows 6-10, columns 6-10) is white. Arrows indicate the flow of data: from the top-left to the top-right, from the top-right to the bottom-right, from the bottom-right to the bottom-left, and from the bottom-left to the top-left. The bottom-right cell is labeled 'G'.



```
epsilon: 0.05
*****
```

```
method: sarsa
```

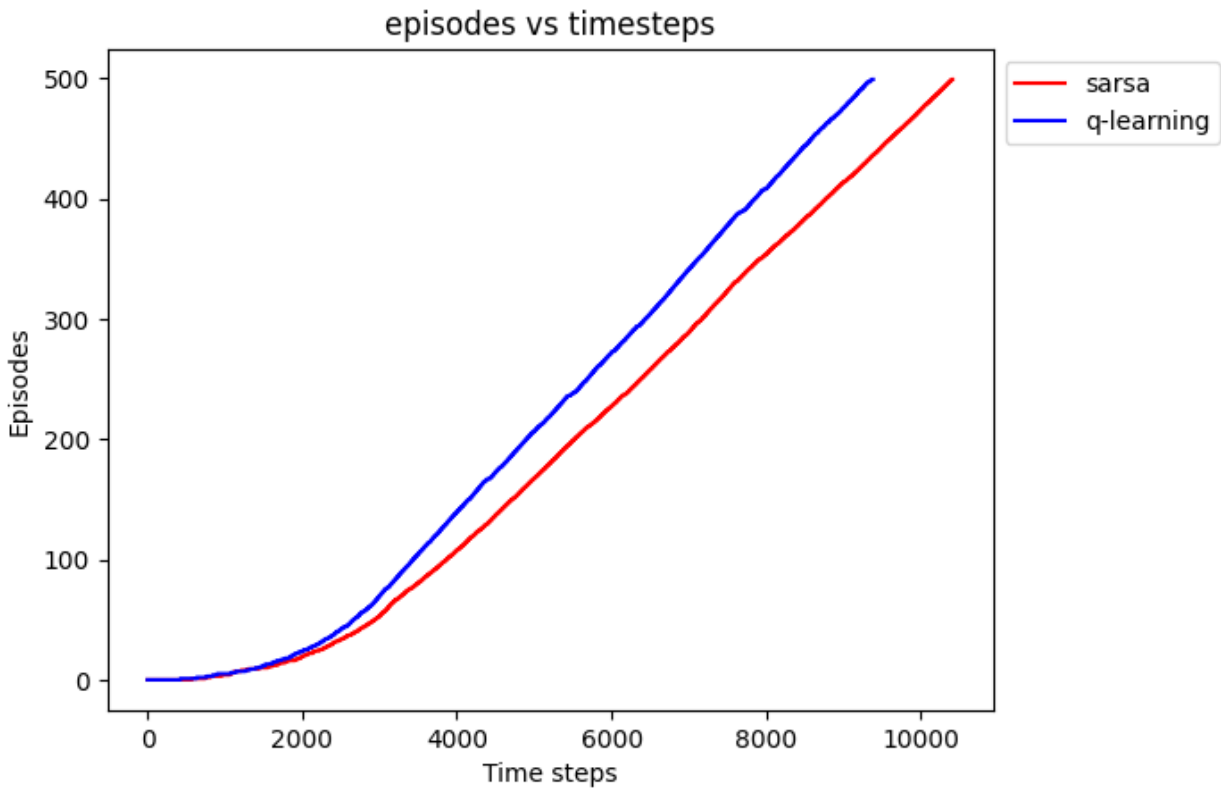
```
>>> trajectory of the latest episode
*      *      *      *      *      *      *      *      *      *      *
→      →      →      →      →      →      →      →      →      ↓      *
↑      *      *      *      *      *      *      *      *      *      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G

>>> optimal policy
→      →      →      →      →      ↓      →      →      →      ↓      →      ↓
→      →      →      →      →      →      →      →      →      →      ↓      ↓
↑      ↑      →      →      ↑      ←      ↑      ↑      →      ↑      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```

```
method: q-learning
```

```
>>> trajectory of the latest episode
*      *      *      *      *      *      *      *      *      *      *
*      *      *      *      *      *      *      *      ↓      *      *
→      →      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G

>>> optimal policy
↑      →      →      ↑      →      →      ↓      →      ↑      →      ↓      ↓
→      →      →      →      ↑      →      →      →      →      ↓      ↓      ↓
→      →      →      →      →      →      →      →      →      →      →      ↓
↑      █      █      █      █      █      █      █      █      █      █      G
```





Вывод:

Наилучший гиперпараметр $\text{EPSILON}=.05$. Дает максимальную суммарную награду. При данном параметре получается оптимальная стратегия.