

Использование Python для обучения с подкреплением



Даниил
Трубин

Газпром-Нефть
Региональные
Продажи

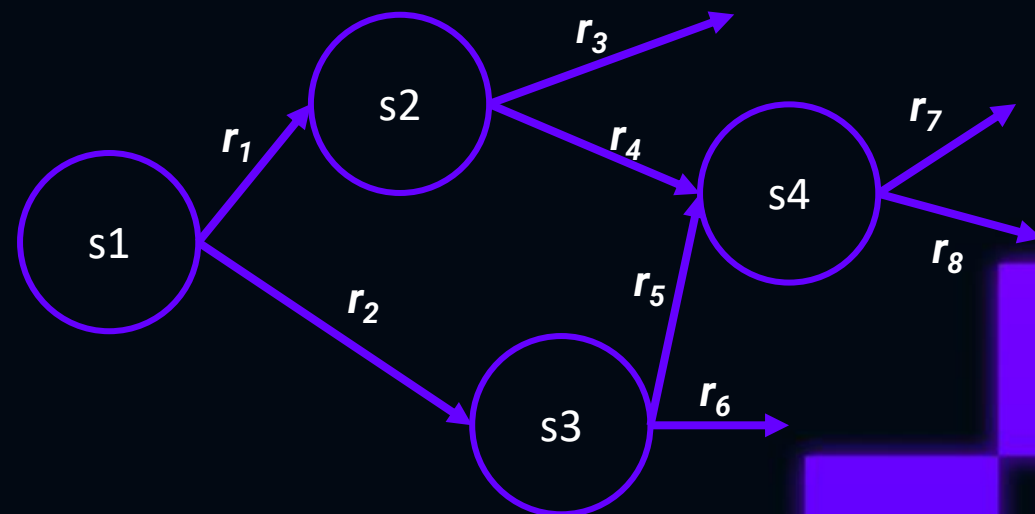
Содержание

- Что такое RL
- Как выглядит RL-проект
- Лайв-кодинг!
- Что дальше?

Уровень 1



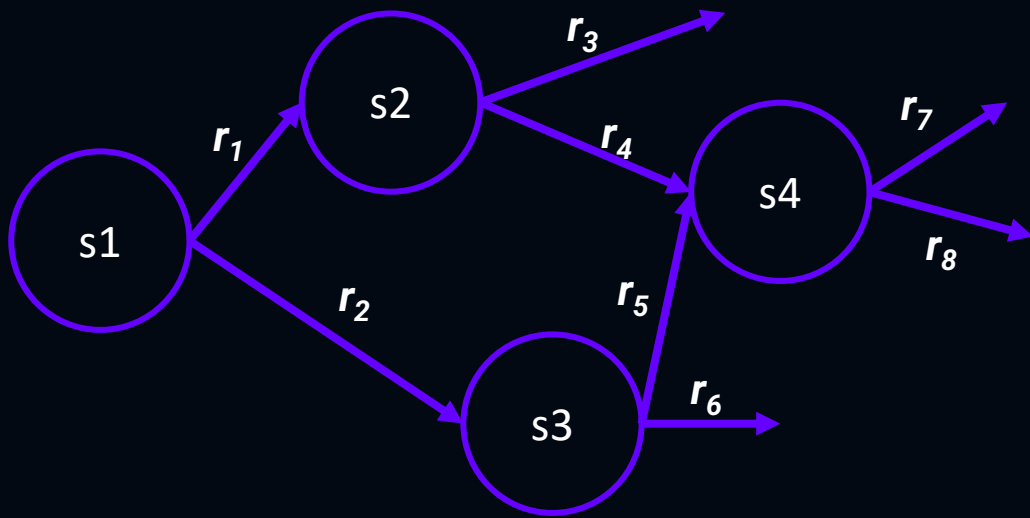
Уровень 2



Уровень 3

$$V(s) = r + V(s')$$

Марковские процесс принятия решений (MDP)



Свойства:

- Текущее состояние не зависит от предыдущих, т.е зависит только от текущего
- Состояния дискретны
- Решение находится итеративно
- Может быть бесконечным (граф цикличным)

Уравнение оптимальности

$$V(s) = \sum_{i=t}^{end} r_i = r_t + \sum_{i=t+1}^{end} r_i = r_t + V(s')$$

~~$\sum_{i=t}^{\infty} \gamma^{i-t} r_i$~~

Как выглядит RL проект

Он должен быть:

- Последовательный (стратегический характер)
- Не завесить от прошлых состояний
- Обучение методом «проб и ошибок»

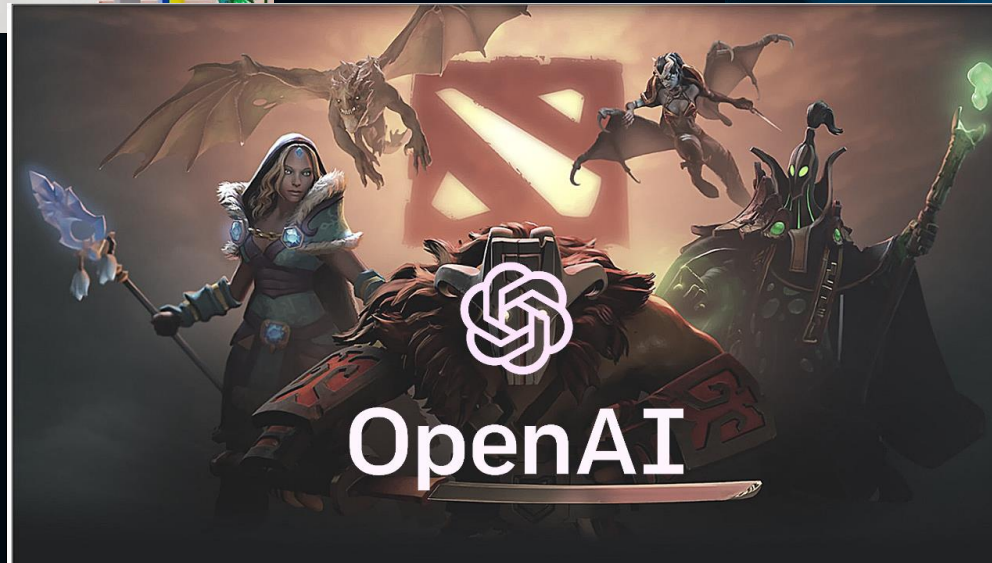
Самые модные успешные кейсы RL



*DeepMind. Охлаждение
дата центров Google*

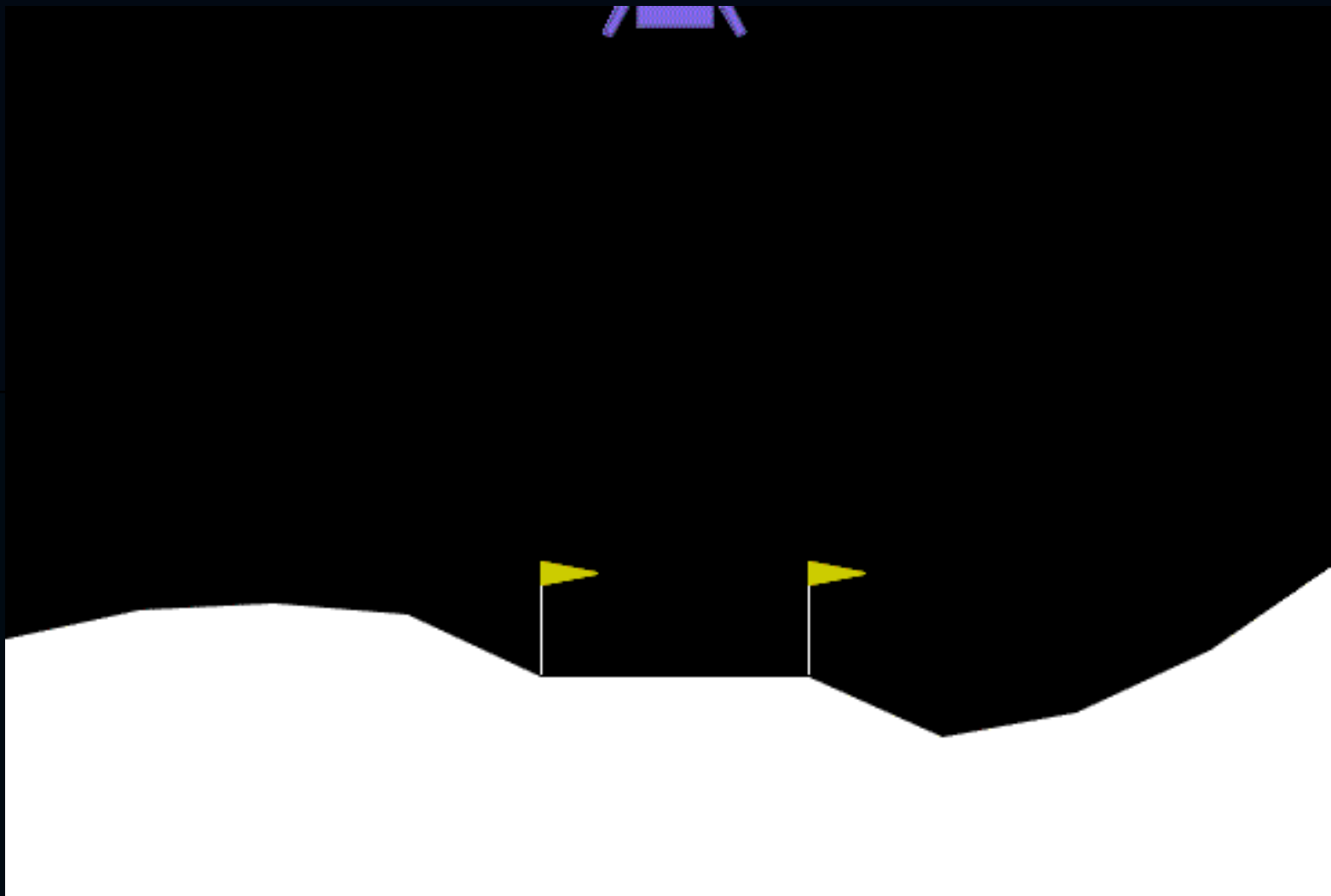


DeepMind. AlphaGo Zero



*OpenAI Five. Победа
команды людей в Дота 2*

Знакомство с проблемой: Lunar Lander

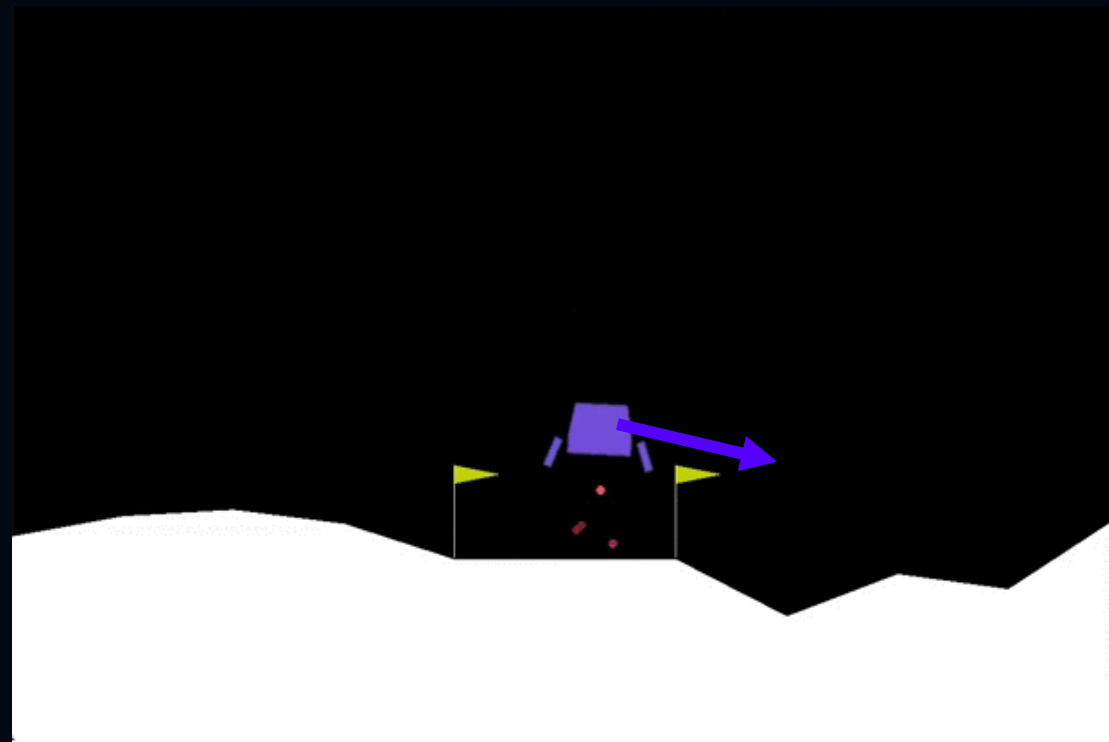
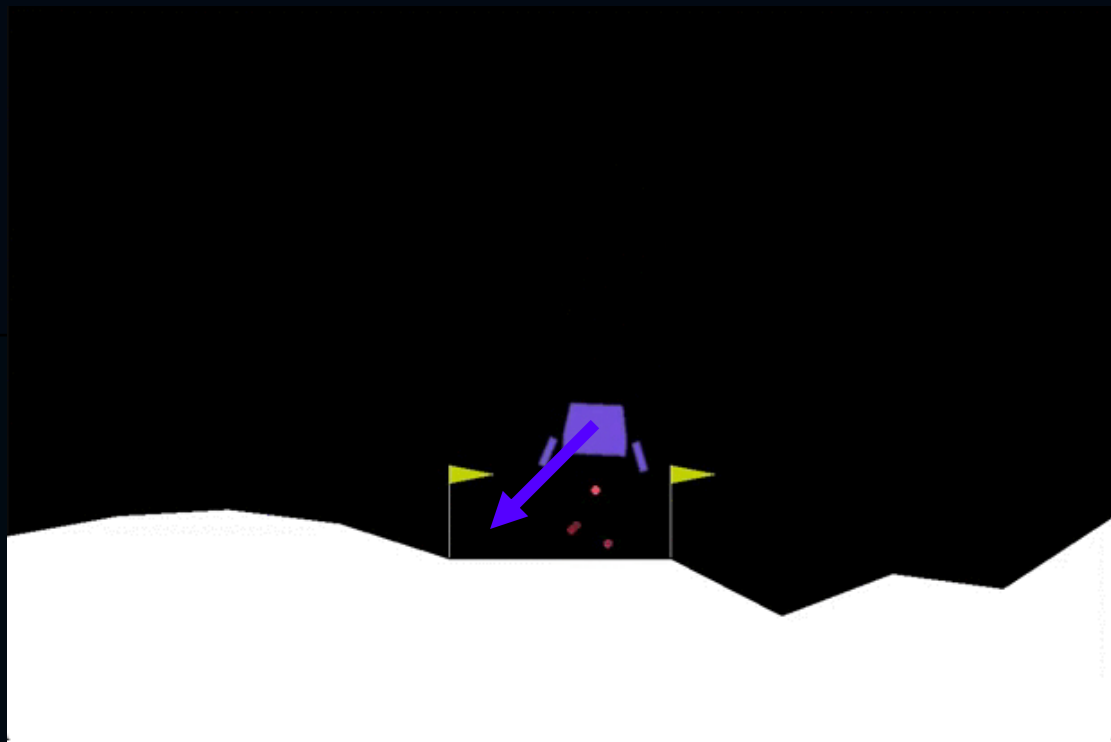


Цель: Посадить дрон на поверхность Луны между флажками

Действия:

- 0 – Ничего не делать
- 1 – Включить левый двигатель
- 2 – Включить основной двигатель
- 3 – Включить правый двигатель

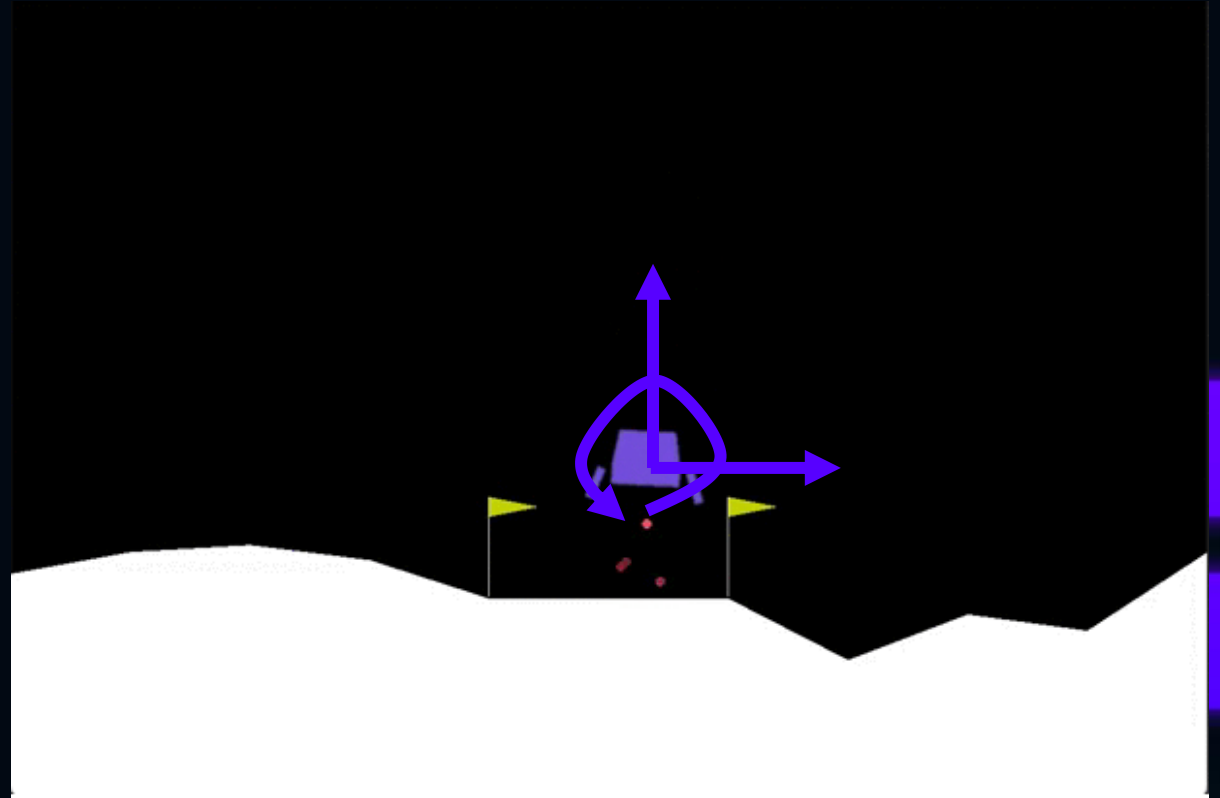
Знакомство с проблемой: Lunar Lander



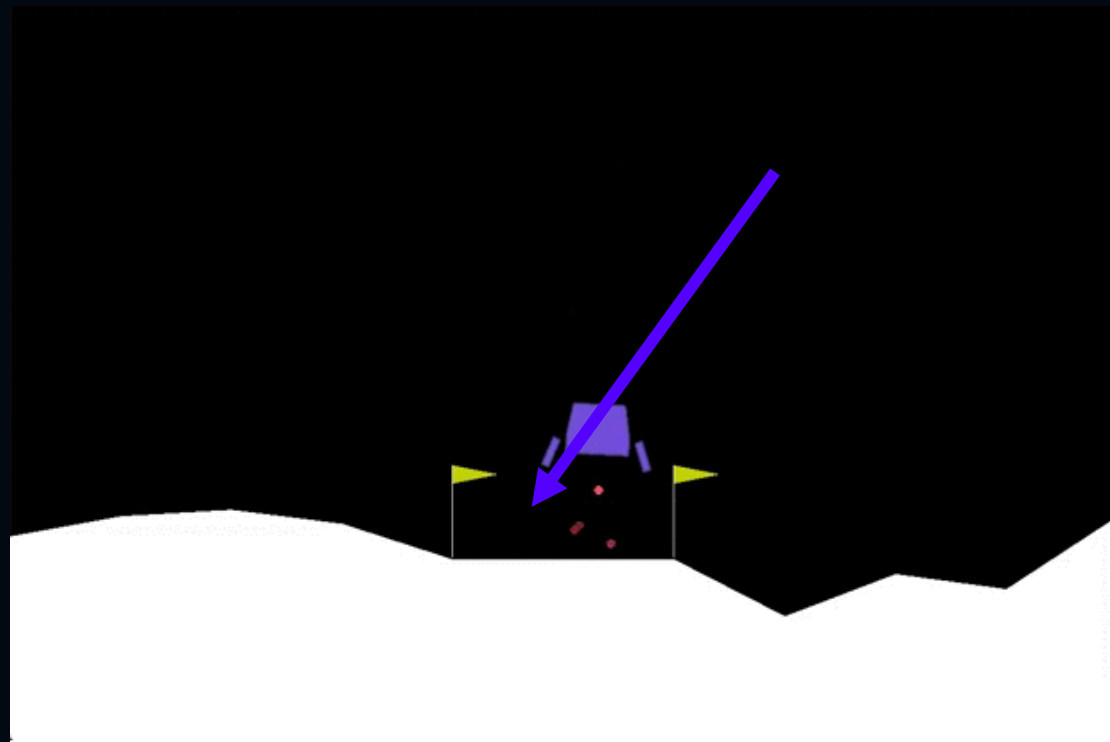
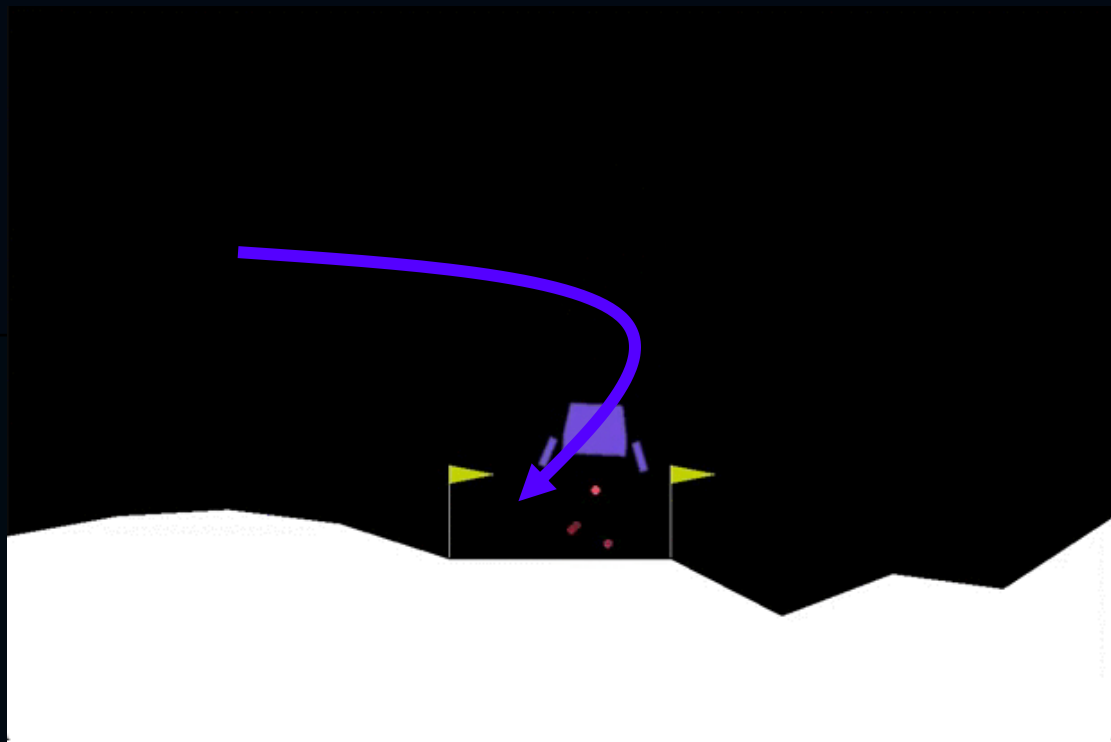
Знакомство с проблемой: Lunar Lander

Состояние описывается вектором :

1. x — Координата
2. y — Координата
3. \dot{x} — Скорость по x
4. \dot{y} — Скорость по y
5. α — Угол наклона
6. $\dot{\alpha}$ — Скорость вращения
7. Касается ли правая нога земли
8. Касается ли левая нога земли

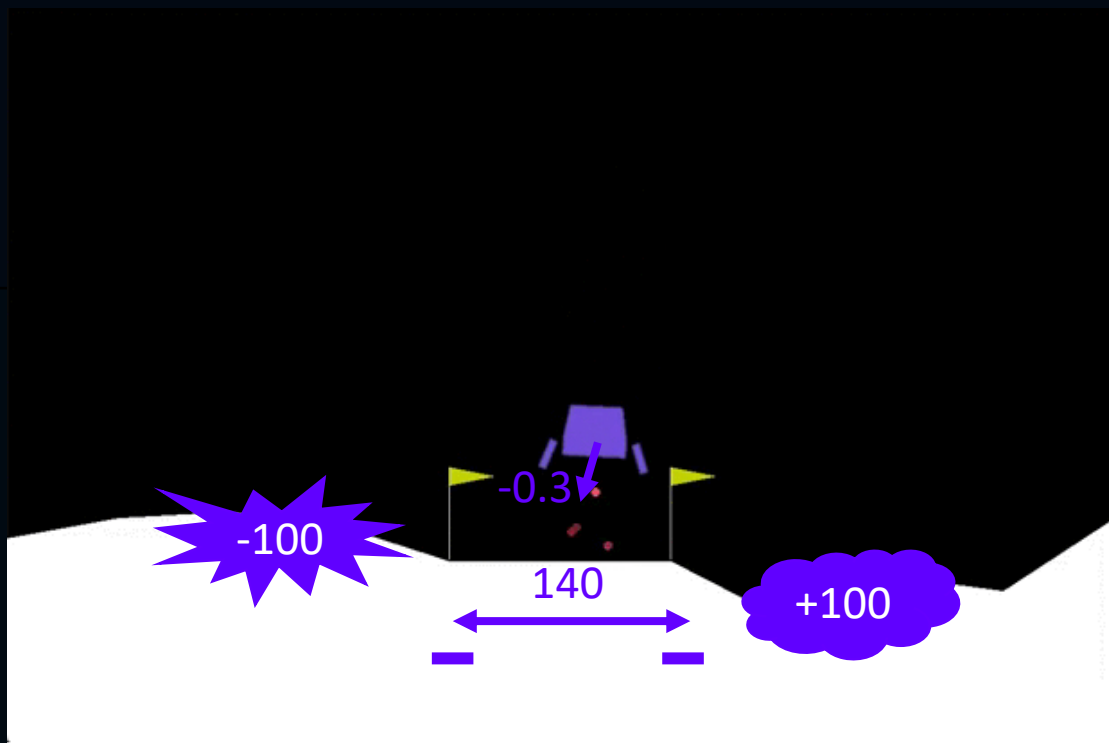


Знакомство с проблемой: Lunar Lander

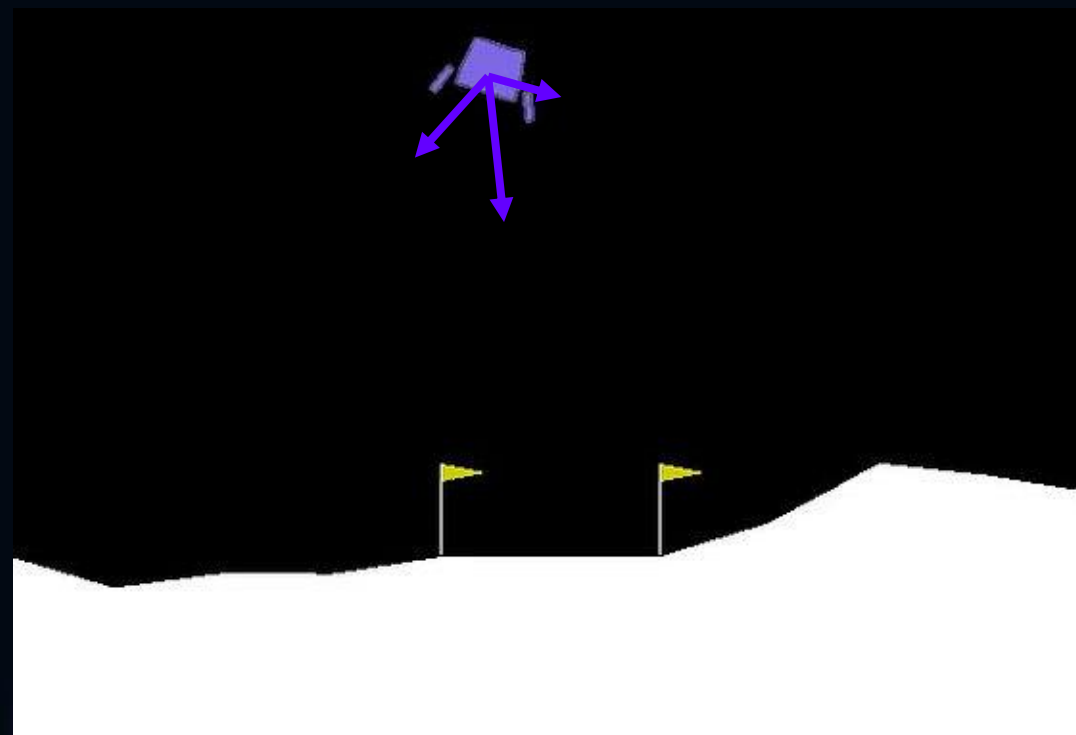


Знакомство с проблемой: Lunar Lander

Награды



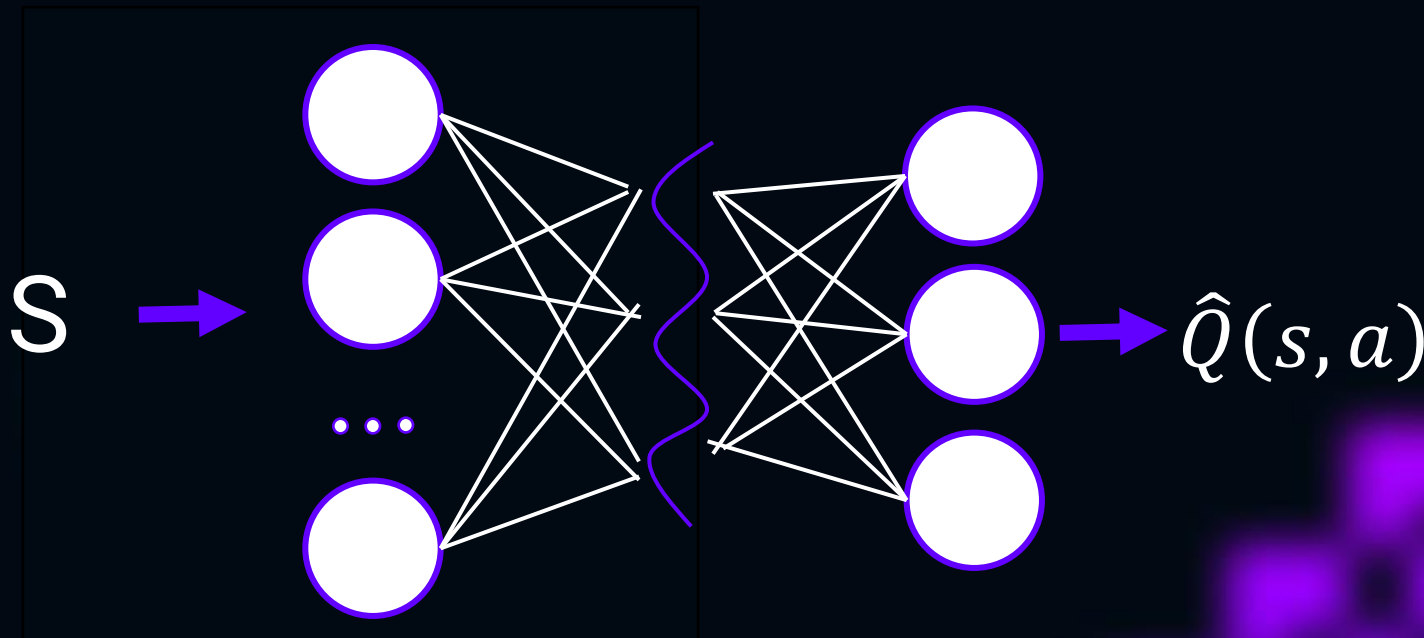
Начальные состояния



Прикинем решение: Deep Q-Network

$$\left. \begin{array}{l} Q(s, a) = r_a + V(s') \\ V(s) = \max_a Q(s, a) \end{array} \right| \Rightarrow Q(s, a) = r_a + \max_{a'} Q(s', a')$$

$$Q(s, a) := Q(s, a) + \alpha * (r_a + \max_{a'} Q(s', a') - Q(s, a))$$



Алгоритм:

1. Рассчитываем $\hat{Q}(s, a)$
2. Эпсилон-жадно выбираем действие a
3. Получаем новое состояние и награду
4. Рассчитываем $\hat{Q}(s', a)$
5. Обновляем DQN

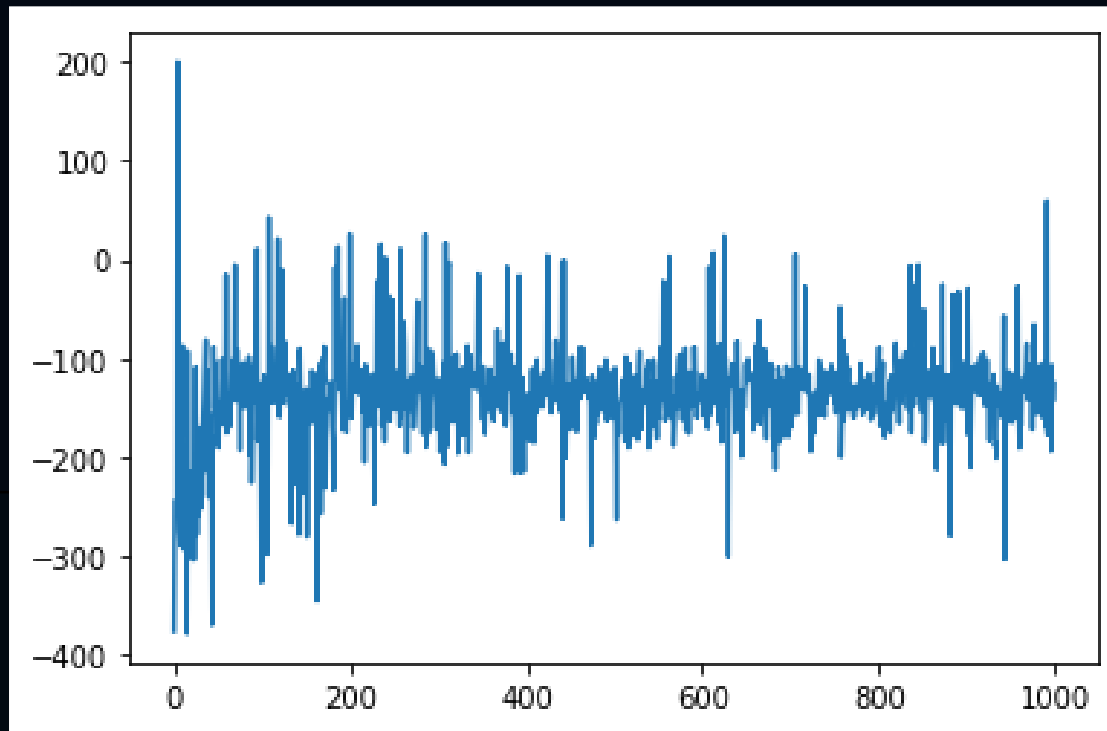
Лайв-кодинг!



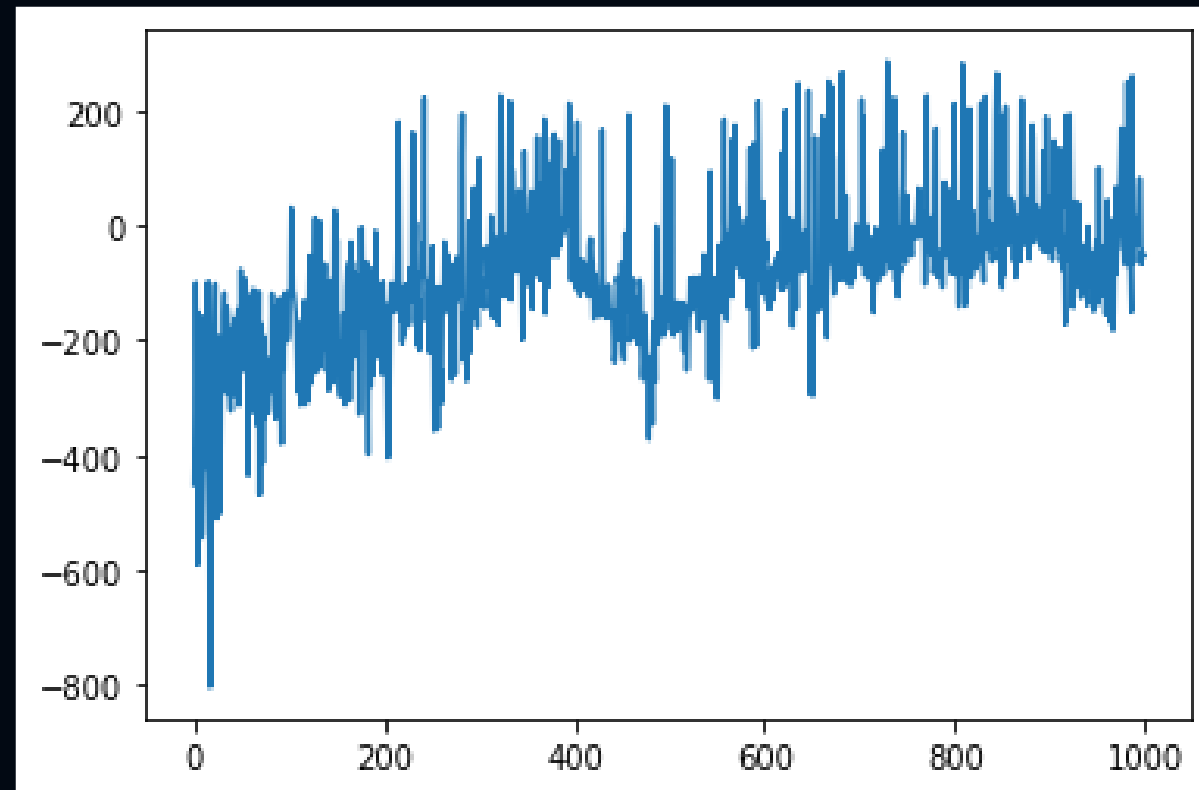
<- Ссылка на гугл-
коллаб

Но лучше использовать
jupyter-notebook

Что там с обучением?



Без Эпсилон-жадности



Эпсилон-жадная стратегия

Что можно улучшить?

- Борьба с автокорреляцией – сделать хаб памяти и обучать на случайных батчах из памяти
- Случайное исследование мира и сбор данных в начале – первичное исследование мира, сокращение ϵ относительно эпизодов
- Модификация наград

Рекомендуемая литература

- Обучение с подкреплением для реальных задач. Инженерный подход. *Фил Уиндер*
- Кто такой многорукий бандит?



- Интуитивный RL (Reinforcement Learning): введение в Advantage-Actor-Critic (A2C)



Спасибо за внимание!