

A decorative graphic on the left side of the slide consists of a grid of colored squares. The grid is 4 squares wide and 4 squares high. The colors of the squares are: Row 1: Teal, Orange, Brown, Teal; Row 2: Orange, Brown, Light Brown, Light Brown; Row 3: Orange, Teal, Light Brown, Light Brown; Row 4: Light Brown, Orange, Orange, Brown.

Введение в Reinforcement Learning

Виктор Кантор

План

1. Примеры задач

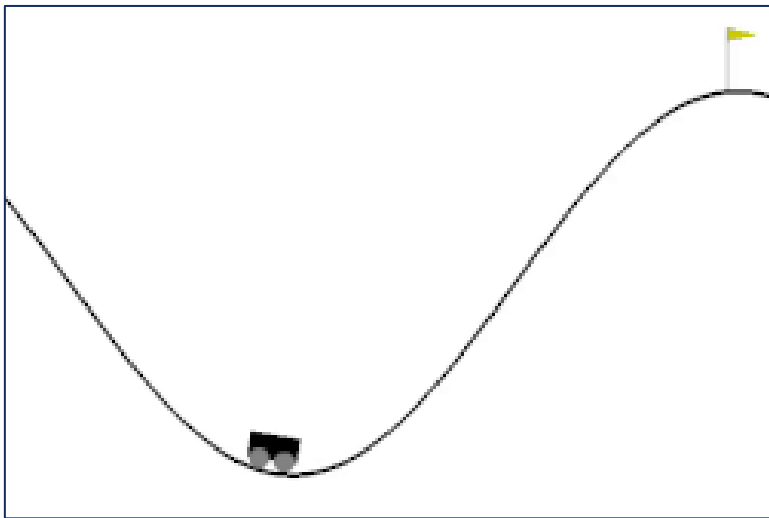
2. Базовые идеи RL

3. SARSA & Q-learning

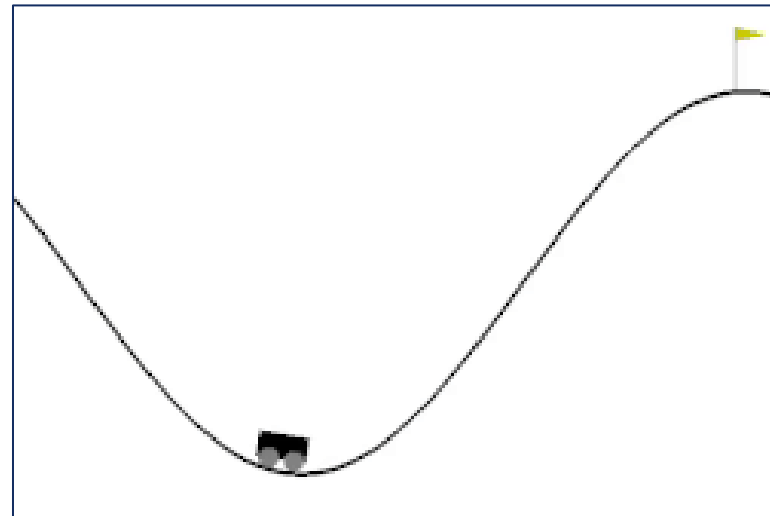
4. Policy Gradient и идеи его улучшения

1. Примеры задач

Mountain car

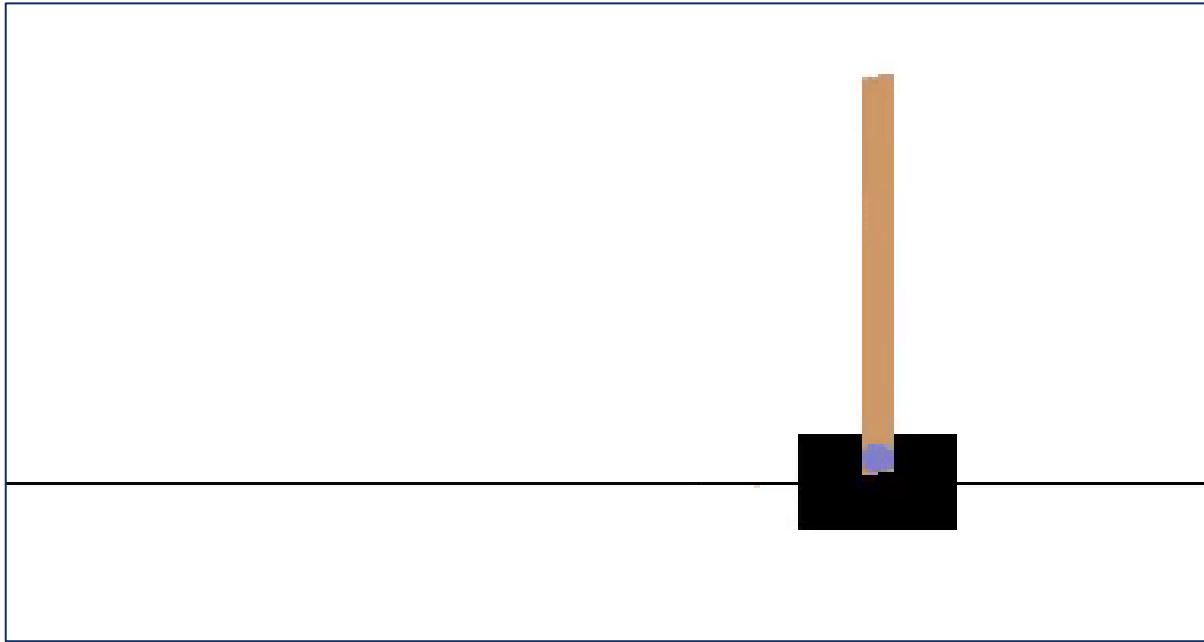


В начале обучения

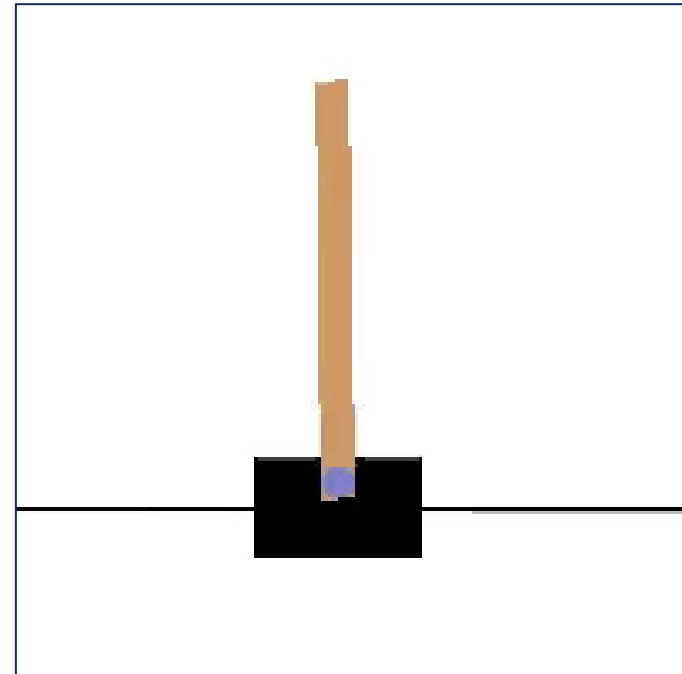


В конце обучения

Cartpole (Перевернутый маятник)



Частично решенная задача



Итоговый алгоритм


Cartpole (Перевернутый маятник)



Источник:
www.youtube.com/watch?v=XiigTGKZfks

Тестовая среда: Open AI Gym

Environments Documentation




Gym

Gym is a toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Pinball.

[View documentation >](#)
[View on GitHub >](#)

Episode 2

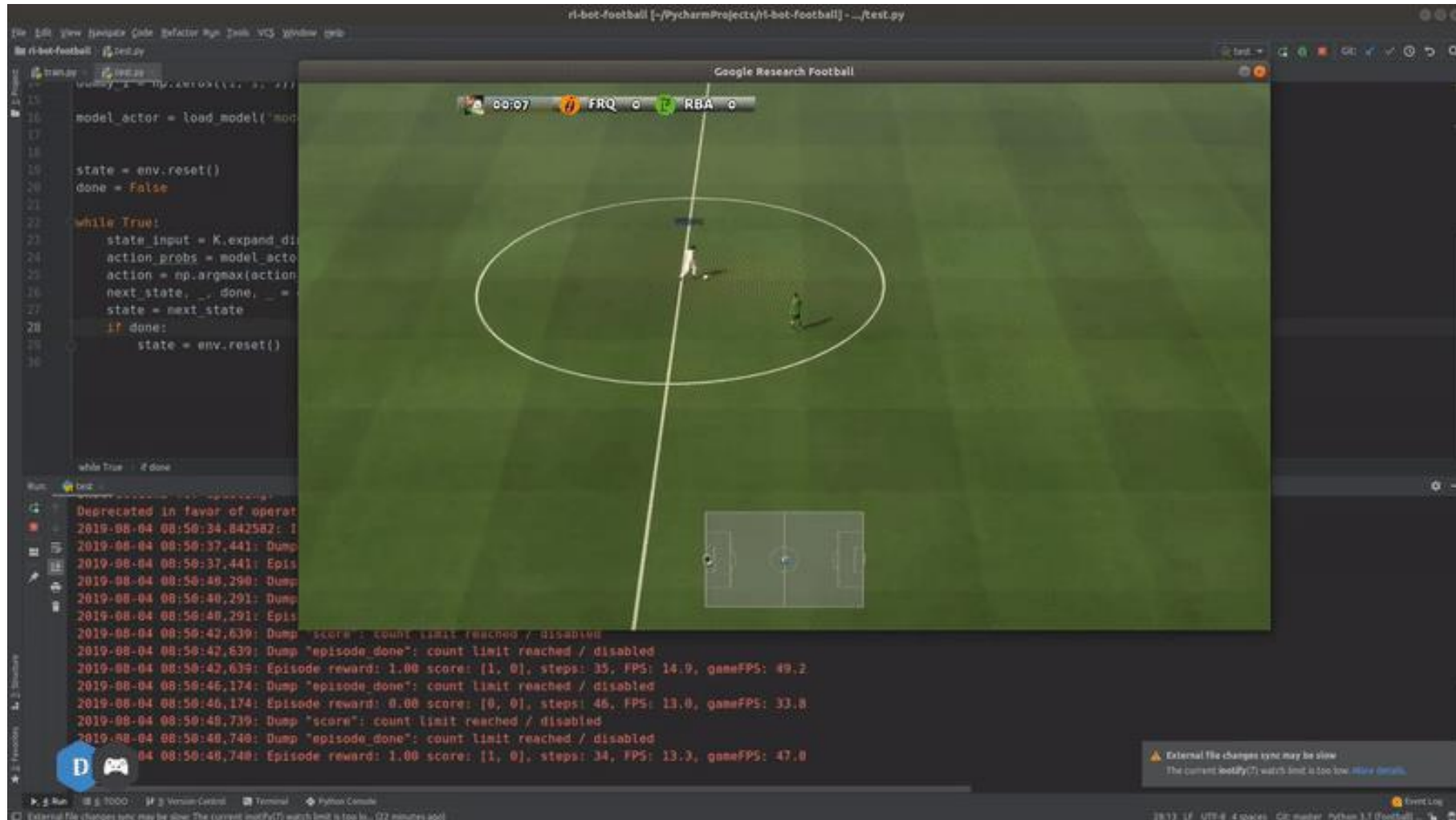
RandomAgent on Pendulum-v0



Episode 1

RandomAgent on SpaceInvaders-v0

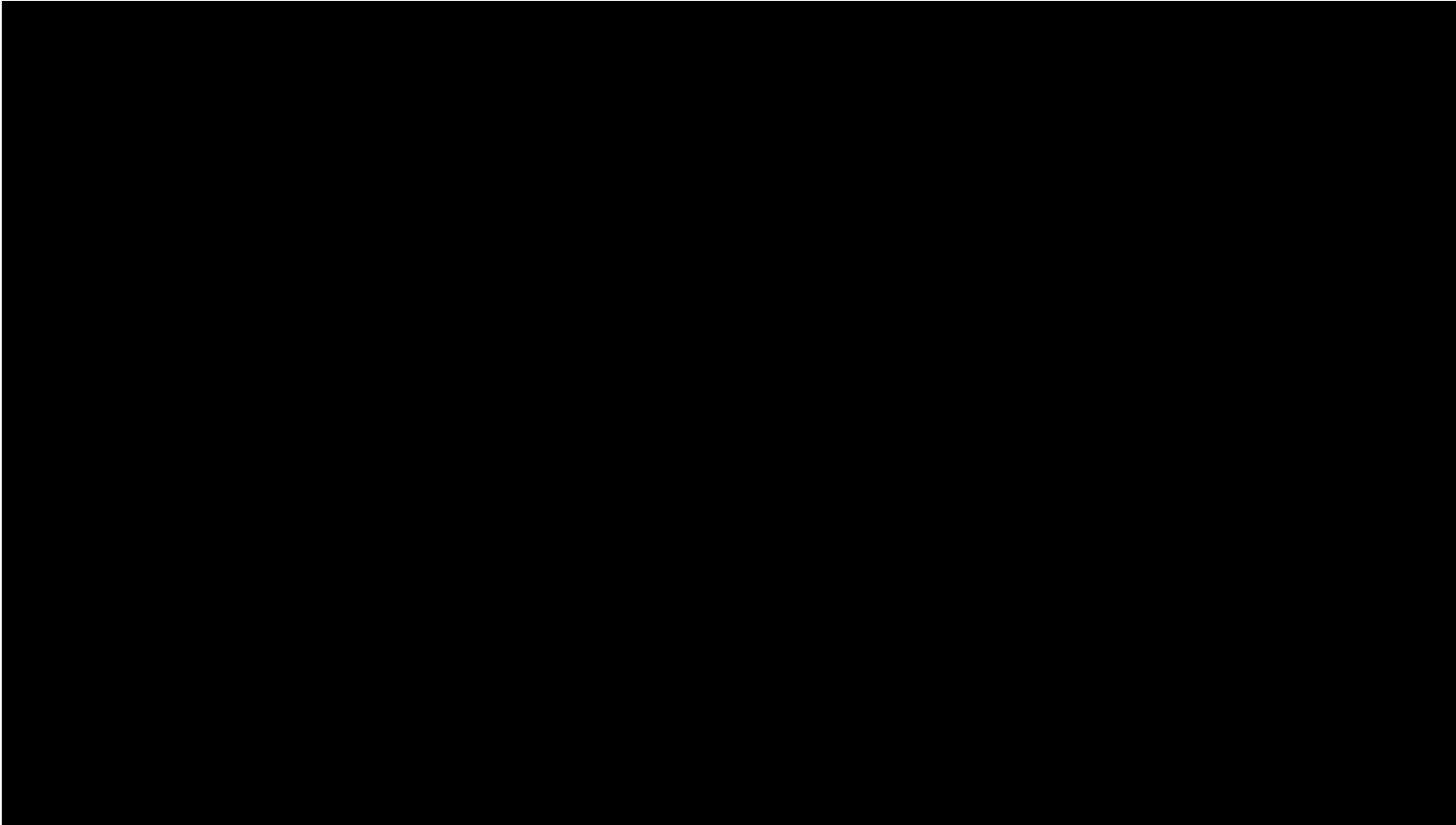
Google Research Football



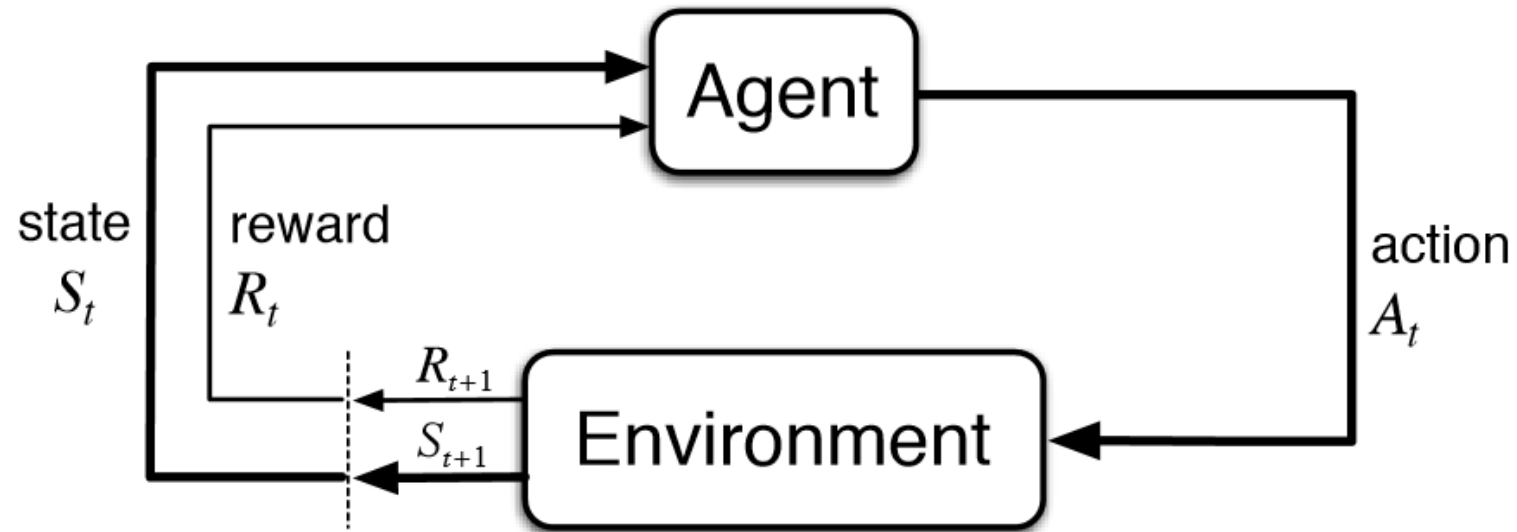
<https://github.com/google-research/football>

<https://ai.googleblog.com/2019/06/introducing-google-research-football.html>

Google Research Football



Взаимодействие агента со средой



2. Базовые идеи RL

Markov Decision Process (MDP)

$$(S, A, \{P_{sa}\}, \gamma, R)$$

- S – множество состояний
- A – множество действий
- $\{P_{sa}\}$ - распределение вероятностей перехода в новое состояние:
 $P_{sa}(s')$ – вероятность перехода из s в s' после действия a
(моделирует поведение среды)
- R – фидбек / награда от среды: $R(s)$ «баллов» получает агент попав в состояние s , в более общем виде - $R(s, a)$ (может влиять не только состояние, но и действие)
- γ – коэффициент дисконтирования награды

Стратегия (policy) и value function

$$\pi: S \rightarrow A$$

$$V^{\pi}(s) = E \left[\sum_{k=0}^{+\infty} \gamma^k R(s_k) \mid s_0 = s, \pi \right]$$

Стратегия (policy) и value function

$$\pi: S \rightarrow A$$

$$V^{\pi}(s) = E \left[\sum_{k=0}^{+\infty} \gamma^k R(s_k) \mid s_0 = s, \pi \right]$$

Вопрос к любителям строгости обозначений: что «не так» с условием в этой вероятности?

Уравнение Беллмана для V^π

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^\pi(s')$$

Уравнение Беллмана для V^π

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^\pi(s')$$

Уравнение Беллмана для V^π

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^\pi(s')$$

Справедливость уравнения очевидна, а польза? Чем оно нам может быть полезно?

Уравнение Беллмана для V^π

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^\pi(s')$$

Справедливость уравнения очевидна, а польза? Чем оно нам может быть полезно?

Имея политику $\pi(s)$ и уже оценив $\{P_{sa}\}$ можем записать $|S|$ линейных уравнений относительно $V^\pi(s)$, решив которые получим значения value function

Как оценить $\{P_{sa}\}$ и $R(s)$

Простой ответ: у вас есть среда, поиграйте с ней 😊

1. Инициализируем π
2. Повторяем:
 - а) Выполнить π в MDP некоторое количество раз
 - б) Обновить $R(s)$ и $P_{sa}(s') = \frac{\#(a,s) \rightarrow s' + 1}{\#(a,s) + |S|}$
 - с) Обновить V и π

Optimal value function

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

$$V^*(s) = R(s) + \max_{a \in A} \gamma \sum_{s' \in S} P_{sa}(s') V^*(s')$$

Optimal policy

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

$$\pi^*(s) = \operatorname{argmax}_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

Value iteration

1. Для всех s инициализируем $V(s) := 0$
2. Повторяем до сходимости для каждого s :

$$V^\pi(s) = R(s) + \max_{a \in A} \gamma \sum_{s' \in S} P_{sa}(s') V^*(s') \quad \leftarrow \text{2}$$

Value iteration

1. Для всех s инициализируем $V(s) := 0$
2. Повторяем до сходимости для каждого s :

$$V^\pi(s) = R(s) + \max_{a \in A} \gamma \sum_{s' \in S} P_{sa}(s') V^*(s') \quad \leftarrow \text{2}$$

Вопрос: как и когда получим стратегию?

Policy iteration

1. Инициализируем π

2. Повторяем:

a) $V := V^\pi$ ①

b) $\pi(s) = \operatorname{argmax}_{a \in A} \sum_{s' \in S} P_{sa}(s') V(s')$ ③

Policy iteration

1. Инициализируем π

2. Повторяем:

a) $V := V^\pi$ ①

b) $\pi(s) = \operatorname{argmax}_{a \in A} \sum_{s' \in S} P_{sa}(s') V(s')$ ③

Вопрос: как вы думаете, когда используют policy iteration, а когда value iteration?

Непрерывное пространство состояний

- С конечным набором состояний теперь справимся
- Что делать, если у нас непрерывное пространство состояний, например координаты?

Непрерывное пространство состояний

- С конечным набором состояний теперь справимся
- Что делать, если у нас непрерывное пространство состояний, например координаты?

Вариант 1. Дискретизация пространства состояний (нарежем координаты на клетки)

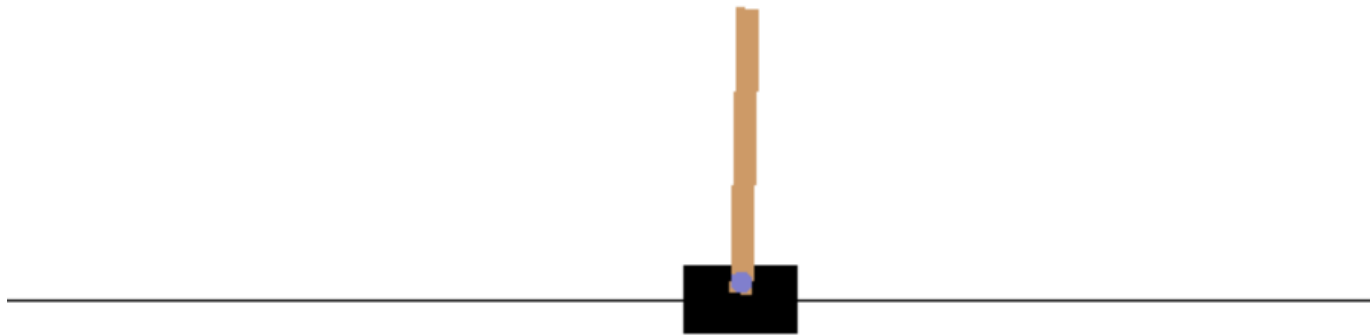
Непрерывное пространство состояний

- С конечным набором состояний теперь справимся
- Что делать, если у нас непрерывное пространство состояний, например координаты?

Вариант 1. Дискретизация пространства состояний (нарежем координаты на клетки)

Вариант 2. Прикручивать модели, которые будут прогнозировать s_{t+1} по s_t и a_t (не обязательно ML, можно физические или иные пригодные для моделирования среды). Также понадобится аппроксимировать моделью value function.

Пример задачи: перевернутый маятник



3. On-policy & off-policy: SARSA & Q-learning

Q-function (state-action value)

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E}_\pi \left(\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right) = \\ &= \mathbb{E}_\pi \left(r_{t+1} + \gamma Q^\pi(s_{t+1}, a_{t+1}) \mid s_t = s, a_t = a \right) \end{aligned}$$

SARSA

- 1: инициализация стратегии $\pi_1(a|s)$ и состояния среды s_1
- 2: **для всех** $t = 1, \dots, T, \dots$
- 3: агент выбирает действие $a_t \sim \pi_t(a|s_t)$:
 $a_t = \arg \max_a Q(s_t, a)$ — жадная стратегия
(но возможны и другие: ε -жадная, по Гиббсу, ...)
- 4: среда генерирует $r_{t+1} \sim p(r|a_t, s_t)$ и $s_{t+1} \sim p(s|a_t, s_t)$;
- 5: агент разыгрывает ещё один шаг: $a' \sim \pi_t(a|s_{t+1})$;
- 6: $Q(s_t, a_t) := Q(s_t, a_t) + \alpha_t(r_{t+1} + \gamma Q(s_{t+1}, a') - Q(s_t, a_t))$;

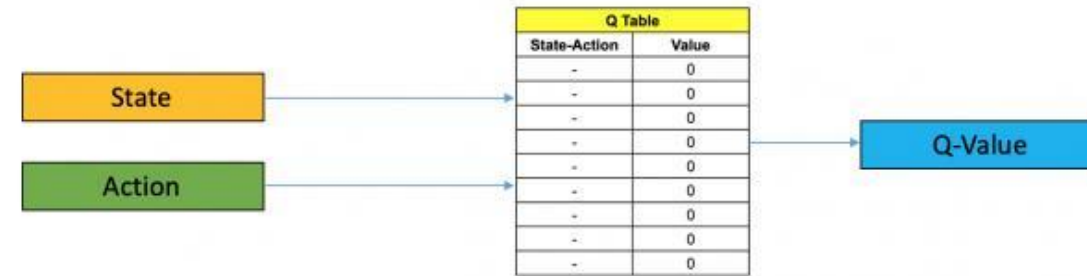
Q-learning

$$Q^*(s, a) = \mathbb{E}(\textcolor{red}{r_{t+1}} + \gamma \max_{a'} \textcolor{red}{Q^*(s_{t+1}, a')} \mid s_t = s, a_t = a)$$

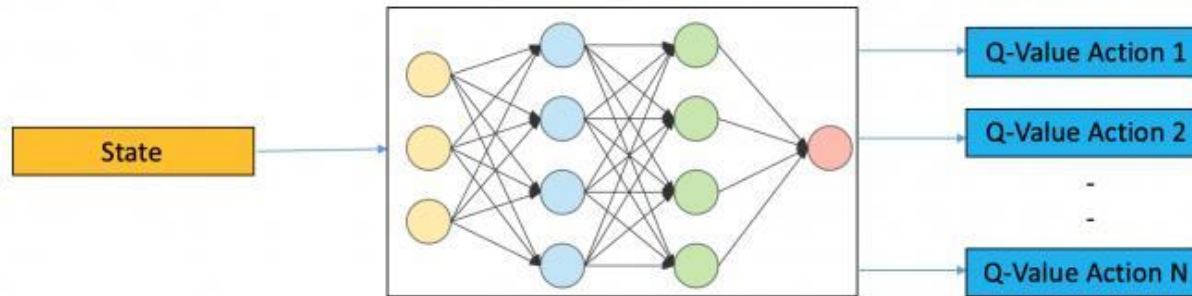
Делаем то же, что в SARSA, но:

$$Q(s_t, a_t) := Q(s_t, a_t) + \alpha_t (\textcolor{red}{r_{t+1}} + \gamma \max_{a'} \textcolor{red}{Q(s_{t+1}, a')} - Q(s_t, a_t))$$

Deep Q-learning



Q Learning



Deep Q Learning

Как добавить модель в обучение

Вместо шага с обновлением state-action value function, появляется шаг с обновлением таргета модели и обновлением параметров модели:

$$y(s') = R_{sa}(s') + \gamma \max_{a'} Q_k(s', a')$$

$$\theta_{k+1} = \theta_k - \alpha \nabla_{\theta} \mathbf{E}_{s' \sim P_{sa}(s')} [(Q_k(s', a') - y(s'))^2] |_{\theta = \theta_k}$$

4. Policy Gradient и идеи его улучшения

Policy Gradient

- Рассмотрим случай стохастической стратегии: вместо детерминированной функции $\pi(s)$ будем обучать распределение $\pi_\theta(a|s)$ в виде модели (например, нейросети) с параметрами θ
- Попробуем выписать градиент оптимизируемого функционала по параметрам стратегии и обновлять их по направлению этого градиента

Policy Gradient

- Рассмотрим случай стохастической стратегии: вместо детерминированной функции $\pi(s)$ будем обучать распределение $\pi_\theta(a|s)$ в виде модели (например, нейросети) с параметрами θ
- Попробуем выписать градиент оптимизируемого функционала по параметрам стратегии и обновлять их по направлению этого градиента

Вопрос: почему по направлению градиента, а не против?

Policy Gradient

- Инициализируем стратегию и насэмплируем T траекторий из MDP
- Изучим оптимизируемый функционал:

$$\eta(\theta) \triangleq \mathbb{E} \left[\sum_{t=0}^{T-1} \gamma^t R(s_t, a_t) \right]$$

$$f(\tau) = \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t)$$

$$\eta(\theta) = \mathbb{E}_{\tau \sim P_\theta} [f(\tau)]$$

Policy Gradient

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \nabla_{\theta} \int P_{\theta}(\tau) f(\tau) d\tau \\ &= \int \nabla_{\theta} (P_{\theta}(\tau) f(\tau)) d\tau \\ &= \int (\nabla_{\theta} P_{\theta}(\tau)) f(\tau) d\tau\end{aligned}$$

Policy Gradient

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \nabla_{\theta} \int P_{\theta}(\tau) f(\tau) d\tau \\ &= \int \nabla_{\theta} (P_{\theta}(\tau) f(\tau)) d\tau \\ &= \int (\nabla_{\theta} P_{\theta}(\tau)) f(\tau) d\tau \\ &= \int P_{\theta}(\tau) (\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau) d\tau\end{aligned}$$

Policy Gradient

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \nabla_{\theta} \int P_{\theta}(\tau) f(\tau) d\tau \\ &= \int \nabla_{\theta} (P_{\theta}(\tau) f(\tau)) d\tau \\ &= \int (\nabla_{\theta} P_{\theta}(\tau)) f(\tau) d\tau \\ &= \int P_{\theta}(\tau) (\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim P_{\theta}} [(\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau)]\end{aligned}$$

Policy Gradient

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \nabla_{\theta} \int P_{\theta}(\tau) f(\tau) d\tau \\ &= \int \nabla_{\theta} (P_{\theta}(\tau) f(\tau)) d\tau \\ &= \int (\nabla_{\theta} P_{\theta}(\tau)) f(\tau) d\tau \\ &= \int P_{\theta}(\tau) (\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim P_{\theta}} [(\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau)]\end{aligned}$$

Вопрос: то, что это красиво, итак понятно, но зачем?

Policy Gradient

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \nabla_{\theta} \int P_{\theta}(\tau) f(\tau) d\tau \\ &= \int \nabla_{\theta} (P_{\theta}(\tau) f(\tau)) d\tau \\ &= \int (\nabla_{\theta} P_{\theta}(\tau)) f(\tau) d\tau \\ &= \int P_{\theta}(\tau) (\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim P_{\theta}} [(\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau)]\end{aligned}$$

Вопрос: то, что это красиво, итак понятно, но зачем?

Ответ: чтобы не брать градиенты матожидания f , если не можем

Как это считать

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \mathbb{E}_{\tau \sim P_{\theta}} [(\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau)] \\ &\approx \frac{1}{n} \sum_{i=1}^n (\nabla_{\theta} \log P_{\theta}(\tau^{(i)})) f(\tau^{(i)})\end{aligned}$$

Как это считать

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \mathbb{E}_{\tau \sim P_{\theta}} [(\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau)] \\ &\approx \frac{1}{n} \sum_{i=1}^n (\nabla_{\theta} \log P_{\theta}(\tau^{(i)})) f(\tau^{(i)})\end{aligned}$$

$$P_{\theta}(\tau) = \mu(s_0) \pi_{\theta}(a_0 | s_0) P_{s_0 a_0}(s_1) \pi_{\theta}(a_1 | s_1) P_{s_1 a_1}(s_2) \cdots P_{s_{T-1} a_{T-1}}(s_T)$$

Как это считать

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \mathbb{E}_{\tau \sim P_{\theta}} [(\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau)] \\ &\approx \frac{1}{n} \sum_{i=1}^n (\nabla_{\theta} \log P_{\theta}(\tau^{(i)})) f(\tau^{(i)})\end{aligned}$$

$$P_{\theta}(\tau) = \mu(s_0) \pi_{\theta}(a_0 | s_0) P_{s_0 a_0}(s_1) \pi_{\theta}(a_1 | s_1) P_{s_1 a_1}(s_2) \cdots P_{s_{T-1} a_{T-1}}(s_T)$$

$$\begin{aligned}\log P_{\theta}(\tau) &= \log \mu(s_0) + \log \pi_{\theta}(a_0 | s_0) + \log P_{s_0 a_0}(s_1) + \log \pi_{\theta}(a_1 | s_1) \\ &\quad + \log P_{s_1 a_1}(s_2) + \cdots + \log P_{s_{T-1} a_{T-1}}(s_T)\end{aligned}$$

Как это считать

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \mathbb{E}_{\tau \sim P_{\theta}} [(\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau)] \\ &\approx \frac{1}{n} \sum_{i=1}^n (\nabla_{\theta} \log P_{\theta}(\tau^{(i)})) f(\tau^{(i)})\end{aligned}$$

$$P_{\theta}(\tau) = \mu(s_0) \pi_{\theta}(a_0 | s_0) P_{s_0 a_0}(s_1) \pi_{\theta}(a_1 | s_1) P_{s_1 a_1}(s_2) \cdots P_{s_{T-1} a_{T-1}}(s_T)$$

$$\begin{aligned}\log P_{\theta}(\tau) &= \log \mu(s_0) + \log \pi_{\theta}(a_0 | s_0) + \log P_{s_0 a_0}(s_1) + \log \pi_{\theta}(a_1 | s_1) \\ &\quad + \log P_{s_1 a_1}(s_2) + \cdots + \log P_{s_{T-1} a_{T-1}}(s_T)\end{aligned}$$

$$\nabla_{\theta} \log P_{\theta}(\tau) = \nabla_{\theta} \log \pi_{\theta}(a_0 | s_0) + \nabla_{\theta} \log \pi_{\theta}(a_1 | s_1) + \cdots + \nabla_{\theta} \log \pi_{\theta}(a_{T-1} | s_{T-1})$$

Итоговый вид градиента по параметрам стратегии

$$\begin{aligned}\nabla_{\theta} \eta(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)] = \mathbb{E}_{\tau \sim P_{\theta}} \left[\left(\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \cdot f(\tau) \right] \\ &= \mathbb{E}_{\tau \sim P_{\theta}} \left[\left(\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \cdot \left(\sum_{t=0}^{T-1} \gamma^t R(s_t, a_t) \right) \right]\end{aligned}$$

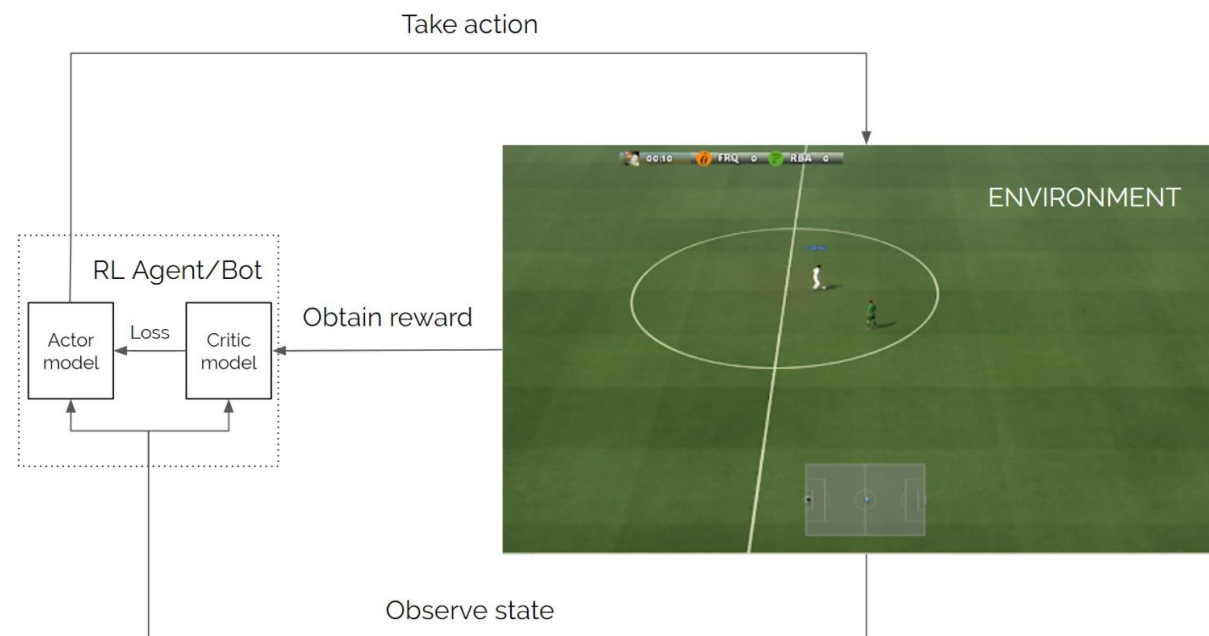
В результате можем применять метод, умея сэмплировать из $\{P_{sa}\}$ и умея получать награду по конкретному состоянию и действию $R(s, a)$ - больше нам ничего не требуется

Идеи в основе Proximal Policy Optimization

- Градиент из PG: $\hat{g} = \hat{\mathbb{E}}_t \left[\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$
- PG Loss: $L^{PG}(\theta) = \hat{\mathbb{E}}_t \left[\log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$

Подход Actor-Critic

- Одна модель (актер) оценивает $\pi_{\theta}(a|s)$ (это было в PG)
- Другая модель (критик) оценивает value function либо advantage (этого не было в PG)



Идеи в основе Proximal Policy Optimization

- Градиент из PG: $\hat{g} = \hat{\mathbb{E}}_t \left[\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$
- PG Loss: $L^{PG}(\theta) = \hat{\mathbb{E}}_t \left[\log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$
- TRPO (Trusted Region Policy Optimization):

$$\begin{aligned} & \underset{\theta}{\text{maximize}} && \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] \\ & \text{subject to} && \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]] \leq \delta. \end{aligned}$$

Идея 1: перенос ограничения на KL в штраф

$$\underset{\theta}{\text{maximize}} \mathbb{E}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t - \beta \text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)] \right]$$

- Коэффициент β перед регуляризатором можно менять динамически
- Применяют правило: если KL меньше/больше «целевого значения KL» (параметр алгоритма), то умножаем/делим β на 2

Идея 2: clipping

$$L^{CPI}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t \left[r_t(\theta) \hat{A}_t \right]$$




$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

Proximal Policy Optimization

$$L^{KL PEN}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t - \beta \text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_\theta(\cdot | s_t)] \right]$$

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

$$L_t^{CLIP+VF+S}(\theta) = \hat{\mathbb{E}}_t \left[L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) + c_2 S[\pi_\theta](s_t) \right]$$

$$(V_\theta(s_t) - V_t^{\text{targ}})^2$$


Энтропия
(для
exploration)



Proximal Policy Optimization

Algorithm 1 PPO, Actor-Critic Style

```
for iteration=1, 2, ... do
  for actor=1, 2, ...,  $N$  do
    Run policy  $\pi_{\theta_{\text{old}}}$  in environment for  $T$  timesteps
    Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
  end for
  Optimize surrogate  $L$  wrt  $\theta$ , with  $K$  epochs and minibatch size  $M \leq NT$ 
   $\theta_{\text{old}} \leftarrow \theta$ 
end for
```

План

1. Примеры задач

2. Базовые идеи RL

3. SARSA & Q-learning

4. Policy Gradient и идеи его улучшения

Резюме

1. RL рассматривает задачу обучения агента максимизировать награду при взаимодействии со средой
2. Наиболее популярные подходы – Deep Q-Learning и модификации Policy Gradients (PG)
3. DQN – value based и основан на уравнении Беллмана для Q-функции, PG – policy based и основан на оптимизации награды градиентным подъемом по стратегии
4. Развитие идеи PG привело к появлению Actor-Critic методов и добавлению регуляризации в оптимизируемом функционале (TRPO, PPO)

Что еще можно изучить

- Заметки курса ML Andrew Ng по Reinforcement learning
- Статьи по Q-learning и Policy Gradients
- Тutorials с применением различных стратегий для задач из Open AI Gym
- Статьи с более сложными методами:
 - <https://arxiv.org/pdf/1707.06347.pdf> (PPO)
 - <https://proceedings.neurips.cc/paper/2017/file/facf9f743b083008a894eee7baa16469-Paper.pdf>