

# Policy gradient методы

Сухарьков Александр, БПМИ 171

21.02.2020

# Содержание

- Приближение стратегии
- Policy gradient теорема
- Log-derivative trick
- REINFORCE
- REINFORCE с бейзлайном
- One-step Actor-Critic

# Приближение стратегии

Policy функция, используемая в данных методах -

$$\pi(a|s, \theta) = \Pr\{A_t = a \mid S_t = s, \theta_t = \theta\}$$

Формула градиентного подъема -

$$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$$

# Приближение стратегии

Один из примеров возможной параметризации, экспоненциальное softmax распределение -

$$\pi(a|s, \boldsymbol{\theta}) = \frac{\exp(h(s, a, \boldsymbol{\theta}))}{\sum_b \exp(h(s, b, \boldsymbol{\theta}))}$$

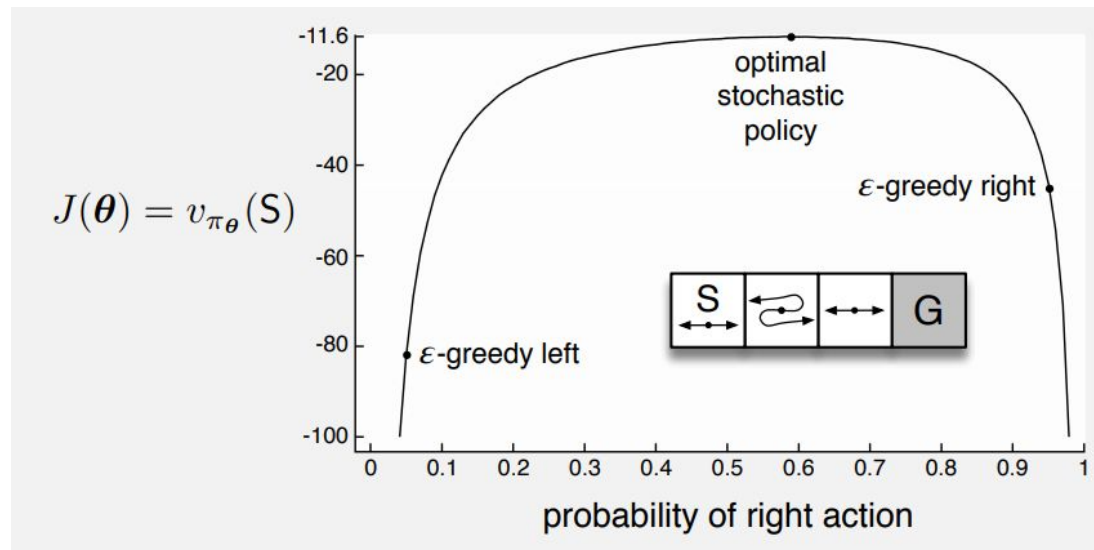
Один из примеров подсчета предпочтений -

$$h(s, a, \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{x}(s, a)$$

# Преимущества над action-value методами

- Возможность приблизиться к детерминированной стратегии
- У policy gradient методов есть возможность найти оптимальную стохастическую стратегию
- Параметризация может упростить процесс аппроксимации
- Добавляет возможные предварительные знания о стратегии

# Приближение стратегии



# Policy gradient теорема

Задаем как показатель эффективности ожидаемую награду из начального состояния с использованием стратегии с параметром -

$$J(\boldsymbol{\theta}) \doteq v_{\pi_{\boldsymbol{\theta}}}(s_0)$$

Сама теорема выглядит таким образом -

$$\nabla J(\boldsymbol{\theta}) \propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta})$$

# Доказательство теоремы

$$\begin{aligned}\nabla v_\pi(s) &= \nabla \left[ \sum_a \pi(a|s) q_\pi(s, a) \right], \quad \text{for all } s \in \mathcal{S} \\ &= \sum_a \left[ \nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \nabla q_\pi(s, a) \right] \\ &= \sum_a \left[ \nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \nabla \sum_{s', r} p(s', r | s, a) (r + v_\pi(s')) \right] \\ &= \sum_a \left[ \nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \sum_{s'} p(s' | s, a) \nabla v_\pi(s') \right] \\ &= \sum_a \left[ \nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \sum_{s'} p(s' | s, a) \right. \\ &\quad \left. \sum_{a'} [\nabla \pi(a' | s') q_\pi(s', a') + \pi(a' | s') \sum_{s''} p(s'' | s', a') \nabla v_\pi(s'')] \right] \\ &= \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} \Pr(s \rightarrow x, k, \pi) \sum_a \nabla \pi(a|x) q_\pi(x, a)\end{aligned}$$



# Доказательство теоремы

$$\begin{aligned}\nabla J(\boldsymbol{\theta}) &= \nabla v_{\pi}(s_0) \\&= \sum_s \left( \sum_{k=0}^{\infty} \Pr(s_0 \rightarrow s, k, \pi) \right) \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \\&= \sum_s \eta(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \\&= \left( \sum_s \eta(s) \right) \sum_s \frac{\eta(s)}{\sum_s \eta(s)} \sum_a \nabla \pi(a|s) q_{\pi}(s, a) \\&\propto \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a). \quad \text{Q.E.D.}\end{aligned}$$

# Log-derivative trick

$$\nabla_{\theta} \log p(\mathbf{x}; \theta) = \frac{\nabla_{\theta} p(\mathbf{x}; \theta)}{p(\mathbf{x}; \theta)}$$

Свойства score функции:

1. Главное вычисление оценки максимального правдоподобия
2. Математическое ожидание равно нулю

$$\begin{aligned}\mathbb{E}_{p(\mathbf{x}; \theta)} [\nabla_{\theta} \log p(\mathbf{x}; \theta)] &= \mathbb{E}_{p(\mathbf{x}; \theta)} \left[ \frac{\nabla_{\theta} p(\mathbf{x}; \theta)}{p(\mathbf{x}; \theta)} \right] \\ &= \int p(\mathbf{x}; \theta) \frac{\nabla_{\theta} p(\mathbf{x}; \theta)}{p(\mathbf{x}; \theta)} d\mathbf{x} = \nabla_{\theta} \int p(\mathbf{x}; \theta) d\mathbf{x} = \nabla_{\theta} 1 = 0\end{aligned}$$

3. Дисперсия является информацией Фишера

$$\mathbb{V}[\nabla_{\theta} \log p(\mathbf{x}; \theta)] = \mathcal{I}(\theta) = \mathbb{E}_{p(\mathbf{x}; \theta)} [\nabla_{\theta} \log p(\mathbf{x}; \theta) \nabla_{\theta} \log p(\mathbf{x}; \theta)^{\top}]$$

# Log-derivative trick

Хотим посчитать градиент от математического ожидания -

$$\nabla_{\theta} \mathbb{E}_{p(z; \theta)} [f(z)] = \nabla_{\theta} \int p(z; \theta) f(z) dz$$

Сведем это равенство к решаемому виду -

$$\nabla_{\theta} \mathbb{E}_{p(z; \theta)} [f(z)] = \mathbb{E}_{p(z; \theta)} [f(z) \nabla_{\theta} \log p(z; \theta)]$$

# Log-derivative trick

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p(z; \theta)}[f(z)] &= \int \nabla_{\theta} p(z; \theta) f(z) dz \\ &= \int \frac{p(z; \theta)}{p(z; \theta)} \nabla_{\theta} p(z; \theta) f(z) dz \\ &= \int p(z; \theta) \nabla_{\theta} \log p(z; \theta) f(z) dz = \mathbb{E}_{p(z; \theta)}[f(z) \nabla_{\theta} \log p(z; \theta)]\end{aligned}$$

# REINFORCE

Должны преобразовать теорему так, чтобы ее можно было бы использовать в градиентном подъеме:

$$\begin{aligned}\nabla J(\boldsymbol{\theta}) &\propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta}) \\ &= \mathbb{E}_\pi \left[ \sum_a q_\pi(S_t, a) \nabla_{\boldsymbol{\theta}} \pi(a|S_t, \boldsymbol{\theta}) \right]\end{aligned}$$

# REINFORCE

$$\begin{aligned}\nabla J(\boldsymbol{\theta}) &= \mathbb{E}_{\pi} \left[ \sum_a \pi(a|S_t, \boldsymbol{\theta}) q_{\pi}(S_t, a) \frac{\nabla_{\boldsymbol{\theta}} \pi(a|S_t, \boldsymbol{\theta})}{\pi(a|S_t, \boldsymbol{\theta})} \right] \\ &= \mathbb{E}_{\pi} \left[ q_{\pi}(S_t, A_t) \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \right] \\ &= \mathbb{E}_{\pi} \left[ G_t \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \right]\end{aligned}$$

(replacing  $a$  by the sample  $A_t \sim \pi$ )

(because  $\mathbb{E}_{\pi}[G_t|S_t, A_t] = q_{\pi}(S_t, A_t)$ )

# REINFORCE

Подставляем полученное в формулу градиентного подъема -

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha G_t \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)}$$

# REINFORCE

Input: a differentiable policy parameterization  $\pi(a|s, \boldsymbol{\theta})$

Initialize policy parameter  $\boldsymbol{\theta} \in \mathbb{R}^{d'}$

Repeat forever:

    Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \boldsymbol{\theta})$

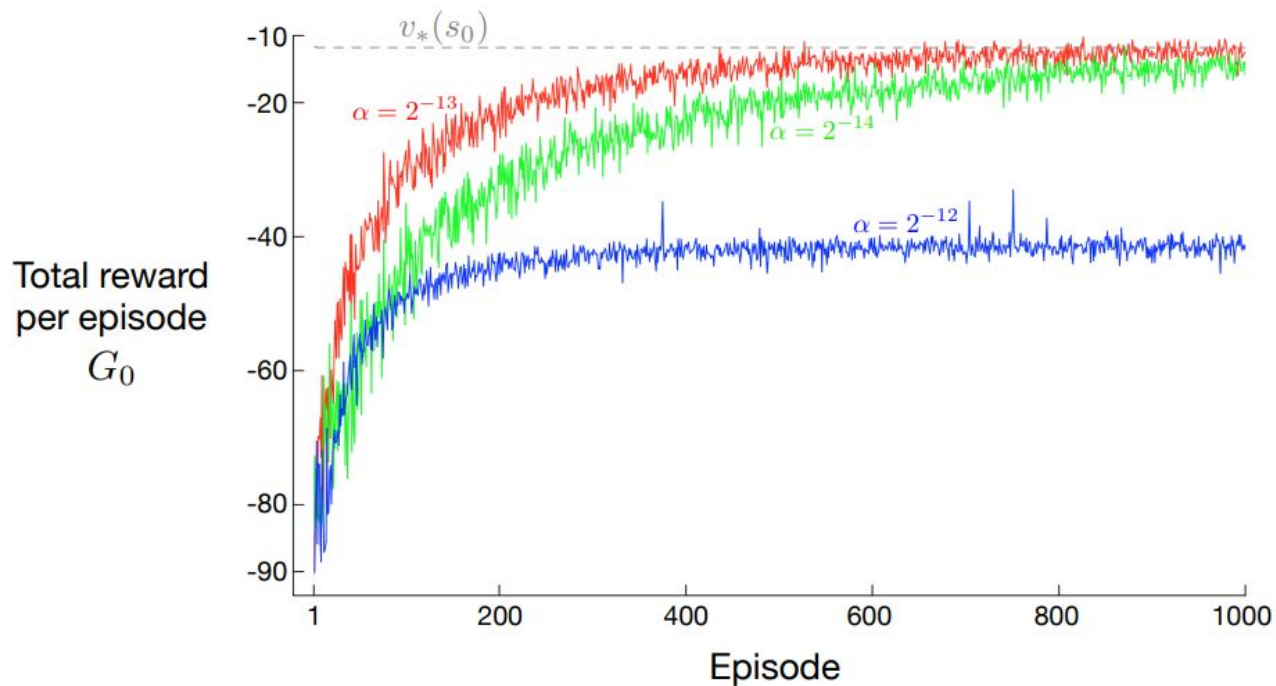
    For each step of the episode  $t = 0, \dots, T - 1$ :

$G \leftarrow$  return from step  $t$

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t G \nabla_{\boldsymbol{\theta}} \ln \pi(A_t|S_t, \boldsymbol{\theta})$



# REINFORCE



# REINFORCE с бейзлайном

Отличается от REINFORCE добавлением бейзлайна в policy gradient теорему

$$\nabla J(\boldsymbol{\theta}) \propto \sum_s \mu(s) \sum_a \left( q_\pi(s, a) - b(s) \right) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta})$$

Можем так сделать, потому что

$$\sum_a b(s) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla_{\boldsymbol{\theta}} \sum_a \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla_{\boldsymbol{\theta}} 1 = 0$$

Получаем новую формулу градиентного подъема

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \left( G_t - b(S_t) \right) \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t|S_t, \boldsymbol{\theta}_t)}{\pi(A_t|S_t, \boldsymbol{\theta}_t)}$$

# REINFORCE с бейзлайном

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$

Input: a differentiable state-value parameterization  $\hat{v}(s, \mathbf{w})$

Parameters: step sizes  $\alpha^\theta > 0$ ,  $\alpha^\mathbf{w} > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  and state-value weights  $\mathbf{w} \in \mathbb{R}^d$

Repeat forever:

    Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$

    For each step of the episode  $t = 0, \dots, T - 1$ :

$G_t \leftarrow$  return from step  $t$

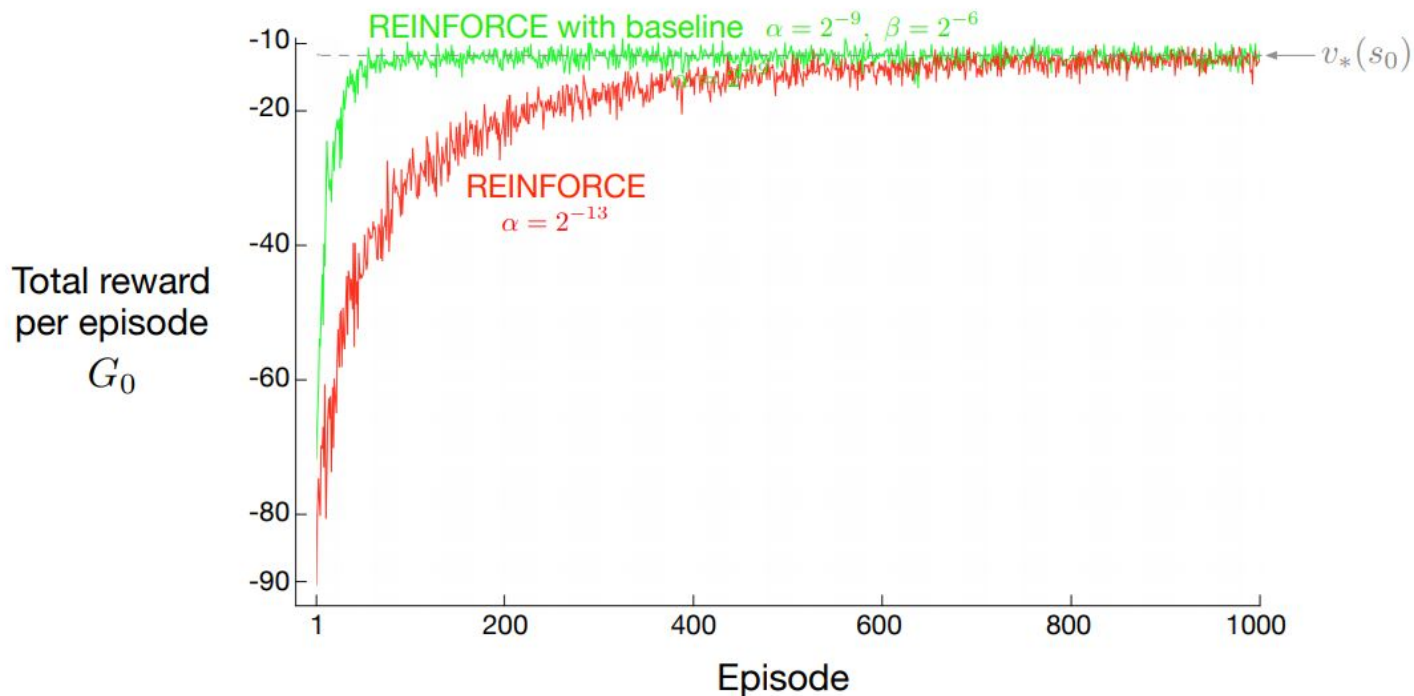
$\delta \leftarrow G_t - \hat{v}(S_t, \mathbf{w})$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha^\mathbf{w} \gamma^t \delta \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$

$\theta \leftarrow \theta + \alpha^\theta \gamma^t \delta \nabla_{\theta} \ln \pi(A_t|S_t, \theta)$

Возможный вариант задания длины шага  $\alpha^\mathbf{w}$ :  $\alpha^\mathbf{w} = 0.1 / \mathbb{E}[\|\nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})\|_\mu^2]$

# Сравнение двух вариантов REINFORCE



# One-step Actor-Critic

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &\doteq \boldsymbol{\theta}_t + \alpha \left( G_{t:t+1} - \hat{v}(S_t, \mathbf{w}) \right) \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)} \\ &= \boldsymbol{\theta}_t + \alpha \left( R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w}) \right) \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)} \\ &= \boldsymbol{\theta}_t + \alpha \delta_t \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)}.\end{aligned}$$

# One-step Actor-Critic

Input: a differentiable policy parameterization  $\pi(a|s, \boldsymbol{\theta})$

Input: a differentiable state-value parameterization  $\hat{v}(s, \mathbf{w})$

Parameters: step sizes  $\alpha^{\boldsymbol{\theta}} > 0$ ,  $\alpha^{\mathbf{w}} > 0$

Initialize policy parameter  $\boldsymbol{\theta} \in \mathbb{R}^{d'}$  and state-value weights  $\mathbf{w} \in \mathbb{R}^d$

Repeat forever:

    Initialize  $S$  (first state of episode)

$I \leftarrow 1$

    While  $S$  is not terminal:

$A \sim \pi(\cdot|S, \boldsymbol{\theta})$

        Take action  $A$ , observe  $S', R$

$\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$       (if  $S'$  is terminal, then  $\hat{v}(S', \mathbf{w}) \doteq 0$ )

$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} I \delta \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})$

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha^{\boldsymbol{\theta}} I \delta \nabla_{\boldsymbol{\theta}} \ln \pi(A|S, \boldsymbol{\theta})$

$I \leftarrow \gamma I$

$S \leftarrow S'$

# Вопросы

1. Какие свойства у score функции? Распишите одно из свойств с помощью log-derivative трюка.
2. Опишите алгоритм работы REINFORCE с бейзлайном.
3. Опишите алгоритм работы one-step actor-critic метода.

# Список литературы

1. Richard S. Sutton and Andrew G. Barto, Reinforcement Learning: An Introduction, глава 13
2. <http://blog.shakirm.com/2015/11/machine-learning-trick-of-the-day-5-log-derivative-trick/> (Log-derivative trick)
3. <https://medium.com/@aminamollaysa/policy-gradients-and-log-derivative-trick-4aad962e43e0>
4. <https://lilianweng.github.io/lil-log/2018/04/08/policy-gradient-algorithms.html>