



Advanced Policy-Based methods. TRPO, PPO

Резяпкин Вячеслав

23.07.2019

Tinkoff.ru



1. Line search

1. Выбрать направление шага (антиградиент, например)
2. Выбрать длину шага (learning rate)

2. Trust region

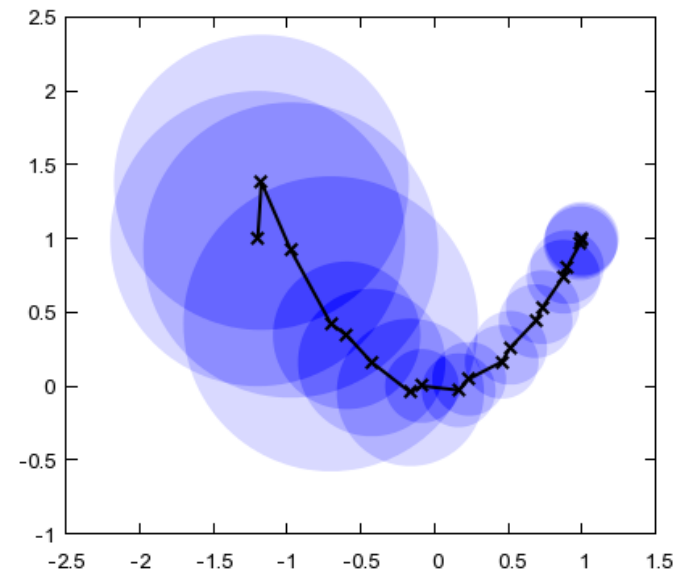
1. Выбрать окрестность для поиска
2. Выбрать точку в окрестности



Trust region methods

Идея:

- Аппроксимировать исходную сложную функцию $f(x)$ (нейронную сеть) более простой $t(x) \cong f(x)$
- Найти минимум простой приближающей функции $t(x)$ в области, где она хорошо приближает $f(x)$





Trust region methods

- В качестве $m(x)$ часто выбирается разложение функции в ряд Тейлора до второго порядка

$$f(x) \cong m(x) = f(x_k) + (x - x_k)\nabla f(x_k) + \frac{1}{2}(x - x_k)\nabla^2 f(x_k)(x - x_k)$$

Почему? Потому что задача оптимизации квадратичных функций хорошо изучена

- В качестве trust region используется Δ -окрестность точки x_k

$$||x - x_k|| < \Delta$$

Алгоритм:

Повторять

1. решить $x_{k+1} = \operatorname{argmin}_{||x-x_k||<\Delta} m(x)$
2. изменить Δ при необходимости



Если гессиан $\mathbf{H} = \nabla^2 f(\mathbf{x}_k)$ положительно определенный,
то это **хорошо**, ведь:

- Задача выпуклая
 - Есть аналитическое решение $\mathbf{x}^* = \mathbf{x}_k - \mathbf{H}^{-1} \nabla f(\mathbf{x}_k)$
- $$m(x) = f(x_k) + (x - x_k) \nabla f(x_k) + \frac{1}{2} (x - x_k) \mathbf{H} (x - x_k)$$
- $$\nabla m(x) = \nabla f(x_k) + \mathbf{H} (x - x_k) = 0$$
- $$\mathbf{x}^* = \mathbf{x}_k - \mathbf{H}^{-1} \nabla f(x_k)$$

Иначе задача **невыпуклая**, но существуют
эффективные методы, основанные на градиентном
спуске



Policy gradients

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}}[R(\tau)] = \mathbb{E}_{s_0 \sim p_0}[V(s_0)] \rightarrow \max$$

Собираем траектории и оцениваем градиент

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \sum_t r_t * \sum_t \nabla \log \pi(s_t, a_t)$$

$$= \mathbb{E}_{\tau \sim \pi_{\theta}} \sum_{t=1}^T \nabla \log \pi(s_t, a_t) * \sum_{t'=t}^T r_{t'}$$

$$= \mathbb{E}_{\tau \sim \pi_{\theta}} \sum_t Q(s_t, a_t) * \nabla \log \pi(s_t, a_t) =$$

$$\mathbb{E}_{\tau \sim \pi_{\theta}} \sum_t (Q(s_t, a_t) - V(s_t)) * \nabla \log \pi(s_t, a_t) =$$

$$\mathbb{E}_{\tau \sim \pi_{\theta}} \sum_t A(s_t, a_t) * \nabla \log \pi(s_t, a_t)$$

Обучение модели:

1. Сэмплирование траекторий
2. Обновление параметров
 $\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} J(\theta)$

Reward-to-go

Определение Q-функции

бейзлайн

Определение Advantage



Policy gradients. Baseline

$$\mathbb{E}_{\tau \sim \pi_{\theta}} \sum_t Q(s_t, a_t) * \nabla \log \pi(s_t, a_t) =$$

$$\mathbb{E}_{\tau \sim \pi_{\theta}} \sum_t (Q(s_t, a_t) - V(s_t)) * \nabla \log \pi(s_t, a_t)$$

бейзлайн

$$\mathbb{E}_{\tau \sim \pi_{\theta}} \sum_t V(s_t) * \nabla \log \pi(s_t, a_t) =$$

$$\sum_t \mathbb{E}_{s_t \sim p, \pi_{\theta}} \mathbb{E}_{a_t \sim \pi_{\theta}(s_t)} V(s_t) * \nabla \log \pi(s_t, a_t) =$$

$$\sum_t \mathbb{E}_{s_t \sim p, \pi_{\theta}} V(s_t) \mathbb{E}_{a_t \sim \pi_{\theta}(s_t)} \nabla \log \pi(s_t, a_t) =$$

$$\sum_t \mathbb{E}_{s_t \sim p, \pi_{\theta}} V(s_t) \nabla \mathbb{E}_{a_t \sim \pi_{\theta}(s_t)} \pi(s_t, a_t) =$$

$$\sum_t \mathbb{E}_{s_t \sim p, \pi_{\theta}} V(s_t) \nabla 1 = 0$$



Trusted Region Policy Optimization(TRPO)

Можно доказать равенство – связь оценок двух разных политик:

$$J(\tilde{\pi}) = J(\pi_{old}) + \mathbb{E}_{\tau \sim \tilde{\pi}} \left[\sum_t A_{\pi_{old}}(s_t, a_t) \right]$$

Ожидаемая
награда с
новой
политикой

Ожидаемая
награда со
старой
политикой

Траектории
сэмплируются
новой
политикой

Оценка
Advantage
функции по
старой
политике

Доказывать не будем, но убедимся, что оно
выполняется при равенстве политик

$$\begin{aligned} A(s, a) &= Q(s, a) - V(s) \\ \mathbb{E}_{a \sim \pi_{\theta}(s)} A(s, a) &= \mathbb{E}_{a \sim \pi_{\theta}(s)} [Q(s, a) - V(s)] = \\ &= \mathbb{E}_{a \sim \pi_{\theta}(s)} [Q(s, a)] - V(s) = V(s) - V(s) = 0 \end{aligned}$$



Trusted Region Policy Optimization(TRPO)

Можно доказать равенство – связь оценок двух разных политик:

$$J(\tilde{\pi}) = J(\pi_{old}) + \mathbb{E}_{\tau \sim \tilde{\pi}} [\sum_t A_{\pi_{old}}(s_t, a_t)]$$

Ожидаемая
награда с
новой
политикой

Ожидаемая
награда со
старой
политикой

Траектории
сэмплируются
новой
политикой

Оценка
Advantage
функции по
старой
политике

Как это максимизировать?

Мат.ожидание берется по траекториям новой политики, которых у нас еще нет.

Зато есть старые траектории. Можно ли как-то оценить эту функцию на старых траекториях?



Trusted Region Policy Optimization(TRPO)

Предыдущее неравенство можно переписать как:

$$J(\tilde{\pi}) = J(\pi_{old}) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$



Discounted visitation frequency

$$\rho_{\pi}(s) = P(s_0 = s) + \gamma P(s_1 = s) + \gamma^2 P(s_2 = s) + \dots$$



Trusted Region Policy Optimization(TRPO)

$$J(\tilde{\pi}) = J(\pi_{old}) + \overset{1}{\sum_s \rho_{\tilde{\pi}}(s)} \overset{2}{\sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)}$$

Discounted visitation frequency
 $\rho_{\pi}(s) = P(s_0 = s) + \gamma P(s_1 = s) + \gamma^2 P(s_2 = s) + \dots$

2

importance sampling

$$\begin{aligned} E_{x \sim p(x)}[f(x)] &= \int p(x) f(x) dx \\ &= \int \frac{q(x)}{q(x)} p(x) f(x) dx \\ &= \int q(x) \frac{p(x)}{q(x)} f(x) dx \\ &= E_{x \sim q(x)} \left[\frac{p(x)}{q(x)} f(x) \right] \end{aligned}$$

2

$$\mathbb{E}_{a \sim \pi_{new}(s)} A_{\pi_{old}}(s, a) = \mathbb{E}_{a \sim \pi_{old}(s)} \frac{\pi_{new}(s)}{\pi_{old}(s)} A_{\pi_{old}}(s, a)$$



$$J(\tilde{\pi}) = J(\pi_{old}) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$



$$L(\tilde{\pi}) = J(\pi_{old}) + \sum_s \rho_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

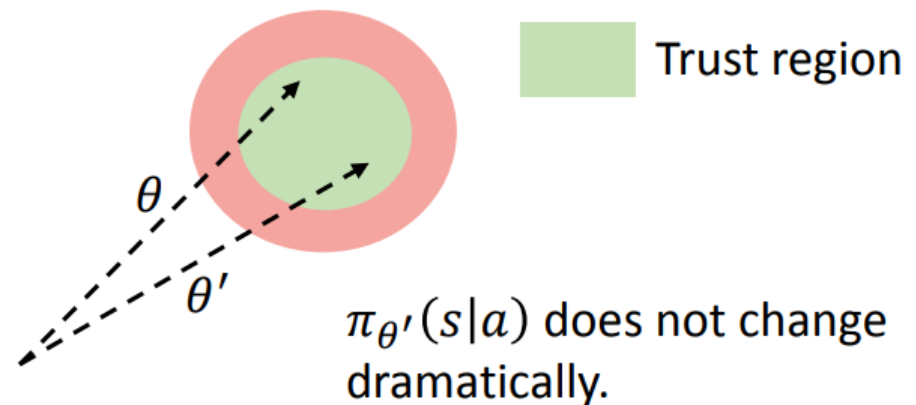
L – локальная аппроксимация J



Trusted Region Policy Optimization(TRPO)

$$L(\tilde{\pi}) = J(\pi_{old}) + \sum_s \rho_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

L – локальная аппроксимация J
в доверительной области параметров





Trusted Region Policy Optimization(TRPO)

$$J(\tilde{\pi}) = J(\pi_{old}) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

Trust Region trick:

If $E_s [KL(\pi || \tilde{\pi})]$ is small,

$$J(\tilde{\pi}) \approx J(\pi) + \sum_s \rho_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

Дивергенция Кульбака-Лейблера:

$$D_{KL}(p||q) = \int_X p(x) \log \frac{p(x)}{q(x)} dx$$

$$\begin{aligned} J(\pi_{old}) &= L(\pi_{old}) \\ \nabla J \Big|_{\pi=\pi_{old}} &= \nabla L \Big|_{\pi=\pi_{old}} \end{aligned}$$

ДЗ



Trusted Region Policy Optimization(TRPO)

$$L(\tilde{\pi}) = J(\pi_{old}) + \sum_s \rho_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

Можно показать, что выполняется неравенство

$$J(\tilde{\pi}) \geq L(\tilde{\pi}) - C D_{KL}^{\max}(\pi, \tilde{\pi})$$

Дивергенция Кульбака-Лейблера:

$$D_{KL}(p||q) = \int_X p(x) \log \frac{p(x)}{q(x)} dx$$

$$C = \frac{4\gamma}{(1-\gamma)^2} \max_{s,a} |A(s, a)|$$

$$\begin{aligned} J(\pi_{old}) &= L(\pi_{old}) \\ \nabla J \Big|_{\pi=\pi_{old}} &= \nabla L \Big|_{\pi=\pi_{old}} \end{aligned} \quad \boxed{\text{ДЗ}}$$

Это означает, что максимизируя правую часть неравенства, мы автоматически максимизируем левую часть



$$L(\tilde{\pi}) = J(\pi_{old}) + \sum_s \rho_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

$$J(\tilde{\pi}) \geq L(\tilde{\pi}) - CD_{KL}^{\max}(\pi, \tilde{\pi})$$

Гарантированное улучшение политики:

$$\pi = \arg \max_{\tilde{\pi}} [L(\tilde{\pi}) - CD_{KL}^{\max}(\pi, \tilde{\pi})]$$



Trusted Region Policy Optimization(TRPO)

На практике используется средняя дивергенция вместо максимальной

$$\underset{\theta}{\text{maximize}} L(\theta) - C \cdot \overline{D}_{KL}(\theta_{old}, \theta)$$

Trust region оптимизационная задача:

Приближение L

Приближение KL

$$\underset{\theta}{\text{maximize}} g \cdot (\theta - \theta_{old}) - \frac{c}{2} (\theta - \theta_{old})^T F (\theta - \theta_{old})$$

$$\text{where, } g = \frac{\partial}{\partial \theta} L(\theta) |_{\theta=\theta_{old}}, \quad F = \frac{\partial^2}{\partial^2 \theta} \overline{D}_{KL}(\theta_{old}, \theta) |_{\theta=\theta_{old}}$$

Где член первого порядка ряда Тейлора у дивергенции и нулевого порядка у обоих членов?



Trusted Region Policy Optimization(TRPO)

$$\max_{\theta} g \cdot (\theta - \theta_{old}) - \frac{c}{2} (\theta - \theta_{old})^T F (\theta - \theta_{old})$$

$$\text{where, } g = \frac{\partial}{\partial \theta} L(\theta) |_{\theta=\theta_{old}}, \quad F = \frac{\partial^2}{\partial^2 \theta} \overline{D_{KL}}(\theta_{old}, \theta) |_{\theta=\theta_{old}}$$

Решение

$$\theta - \theta_{old} = \frac{1}{c} F^{-1} g$$

Не хочется обращаться к гессиану

Вспользуемся тем, что мы хотим посчитать произведение $F^{-1}g$

Оно равняется решению уравнения $Fx = g$

Эффективно решается с помощью метода сопряженных градиентов



Trusted Region Policy Optimization(TRPO)

Как выбирать C?

- Unconstrained problem: $\underset{\theta}{\text{maximize}} L(\theta) - C \cdot \overline{D_{KL}}(\theta_{old}, \theta)$
- Constrained problem: $\underset{\theta}{\text{maximize}} L(\theta)$ subject to $C \cdot \overline{D_{KL}}(\theta_{old}, \theta) \leq \delta$
- δ is a hyper-parameter, remains fixed over whole learning process
- Solve constrained quadratic problem: compute $F^{-1}g$ and then rescale step to get correct KL
 - $\underset{\theta}{\text{maximize}} g \cdot (\theta - \theta_{old})$ subject to $\frac{1}{2}(\theta - \theta_{old})^T F(\theta - \theta_{old}) \leq \delta$
 - Lagrangian: $\mathcal{L}(\theta, \lambda) = g \cdot (\theta - \theta_{old}) - \frac{\lambda}{2}[(\theta - \theta_{old})^T F(\theta - \theta_{old}) - \delta]$
 - Differentiate wrt θ and get $\theta - \theta_{old} = \frac{1}{\lambda} F^{-1}g$
 - We want $\frac{1}{2} s^T F s = \delta$
 - Given candidate step $s_{unscaled}$ rescale to $s = \sqrt{\frac{2\delta}{s_{unscaled}^T F s_{unscaled}}} s_{unscaled}$



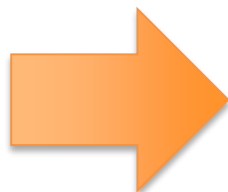
Из-за

- перехода от \mathbf{C} к δ
- приближенных вычислений (сопряженные градиенты, Тейлор)
- семплирования вместо матожиданий
- ...

данный шаг может оказаться слишком большим

Проверим необходимые условия и, если они не выполняются, уменьшим длину, не меняя направления

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g$$



$$\theta_{k+1} = \theta_k + \alpha^j \sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g$$

j – мин. целое при котором:

$$\begin{aligned} \mathcal{L}(\theta_k, \theta_{k+1}) &> 0 \\ \bar{D}_{KL}(\theta_{k+1} || \theta_k) &\leq \delta \end{aligned}$$



Trusted Region Policy Optimization(TRPO)

Algorithm 1 Trust Region Policy Optimization

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: Hyperparameters: KL-divergence limit δ , backtracking coefficient α , maximum number of backtracking steps K
- 3: **for** $k = 0, 1, 2, \dots$ **do**
- 4: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 5: Compute rewards-to-go \hat{R}_t .
- 6: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 7: Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta_k} \hat{A}_t.$$

- 8: Use the conjugate gradient algorithm to compute

$$\hat{x}_k \approx \hat{H}_k^{-1} \hat{g}_k,$$

where \hat{H}_k is the Hessian of the sample average KL-divergence.

- 9: Update the policy by backtracking line search with

$$\theta_{k+1} = \theta_k + \alpha^j \sqrt{\frac{2\delta}{\hat{x}_k^T \hat{H}_k \hat{x}_k}} \hat{x}_k,$$

where $j \in \{0, 1, 2, \dots, K\}$ is the smallest value which improves the sample loss and satisfies the sample KL-divergence constraint.

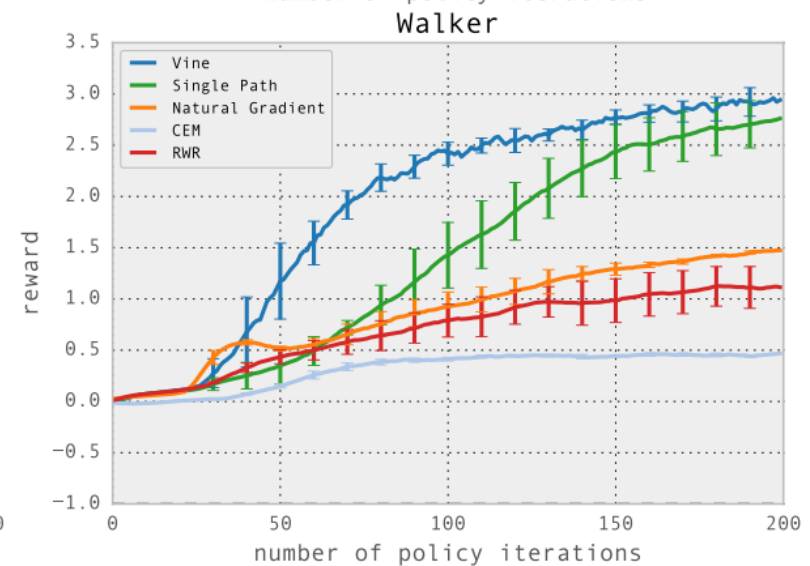
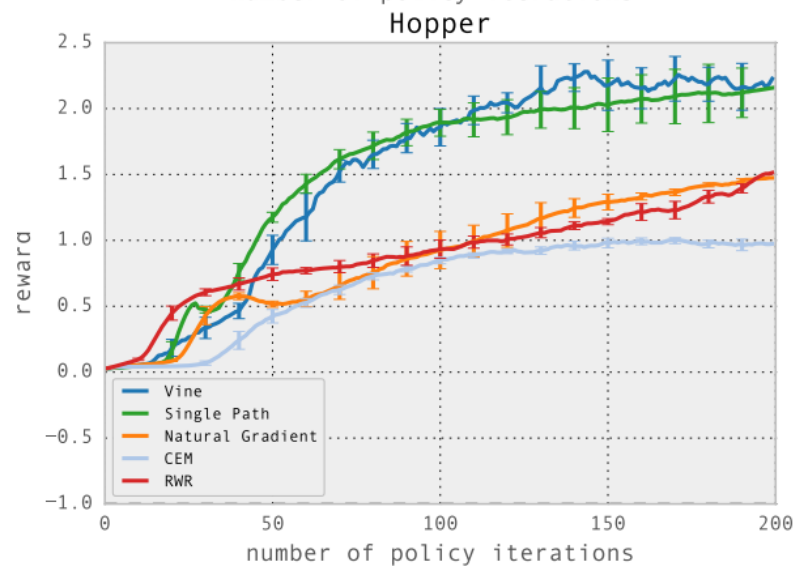
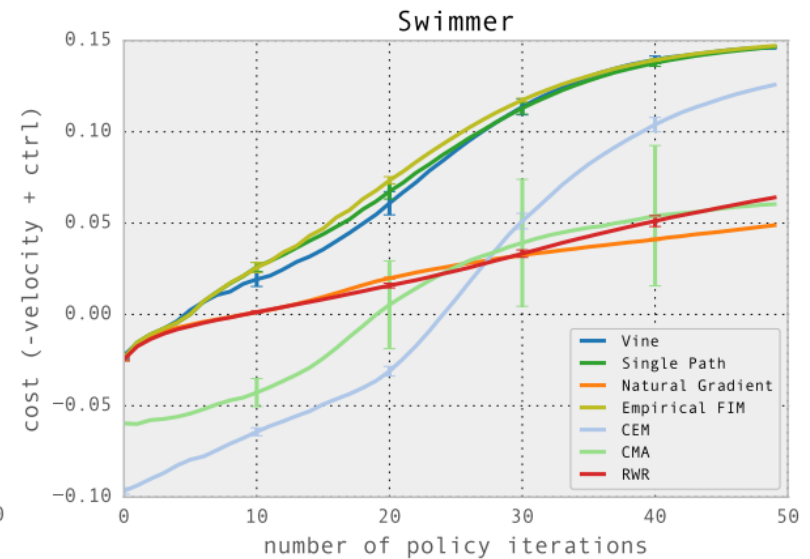
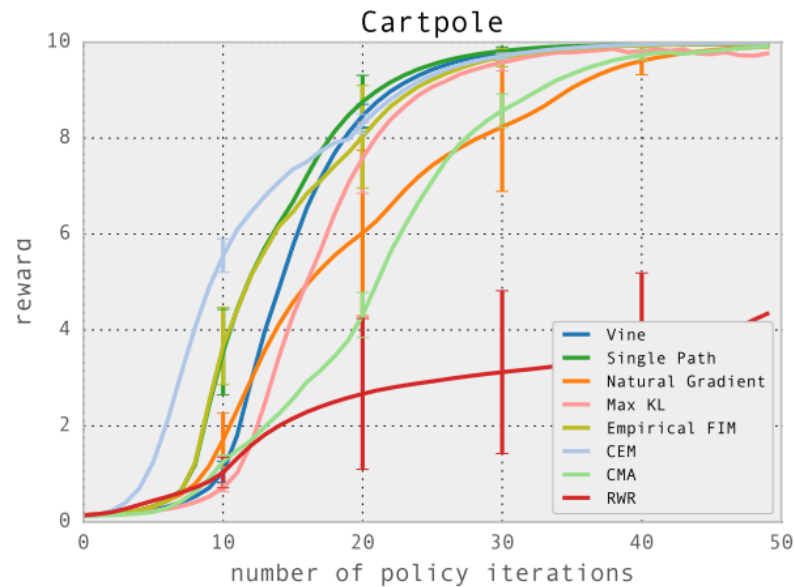
- 10: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 11: **end for**
-

Результаты



Результаты



	<i>B. Rider</i>	<i>Breakout</i>	<i>Enduro</i>	<i>Pong</i>	<i>Q*bert</i>	<i>Seaquest</i>	<i>S. Invaders</i>
Random	354	1.2	0	−20.4	157	110	179
Human (Mnih et al., 2013)	7456	31.0	368	−3.0	18900	28010	3690
Deep Q Learning (Mnih et al., 2013)	4092	168.0	470	20.0	1952	1705	581
UCC-I (Guo et al., 2014)	5702	380	741	21	20025	2995	692
TRPO - single path	1425.2	10.8	534.6	20.9	1973.5	1908.6	568.4
TRPO - vine	859.5	34.2	430.8	20.9	7732.5	788.4	450.2



- Using several epochs of minibatch SGD, optimize the KL-penalized objective

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

- Compute $d = \hat{\mathbb{E}}_t[\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_\theta(\cdot | s_t)]]$
 - If $d < d_{\text{targ}}/1.5$, $\beta \leftarrow \beta/2$
 - If $d > d_{\text{targ}} \times 1.5$, $\beta \leftarrow \beta \times 2$

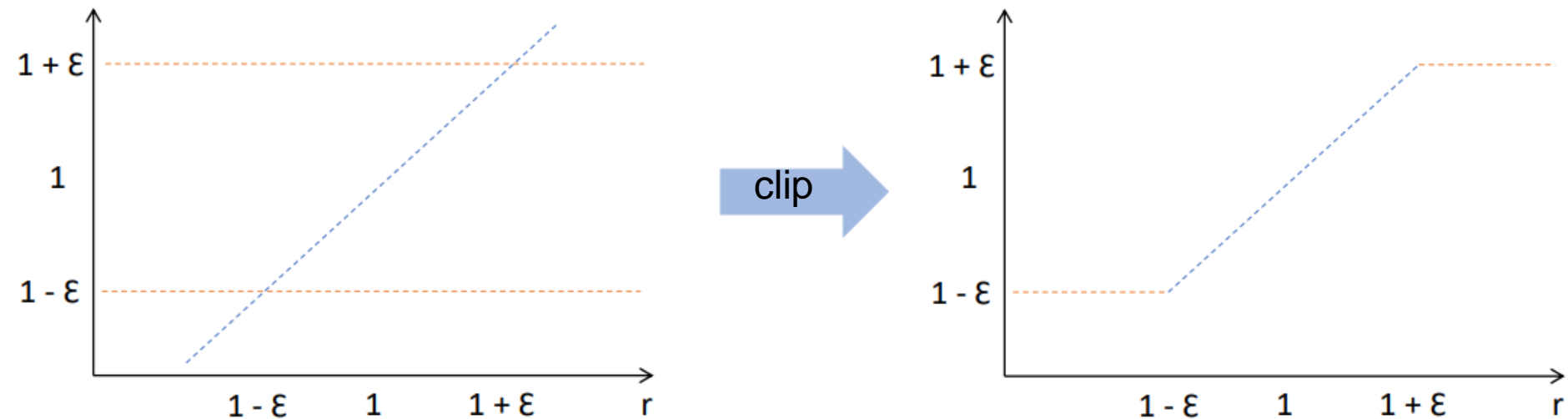


Proximal Policy Optimization(PPO)

По-прежнему будем максимизировать приближение L, но теперь проще

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

$$r_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$$



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

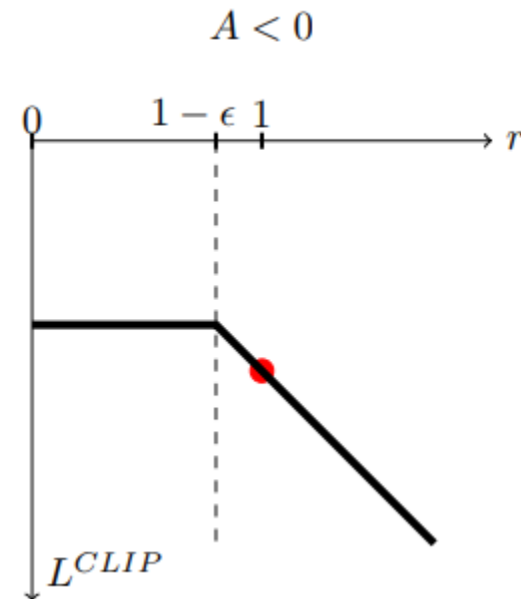
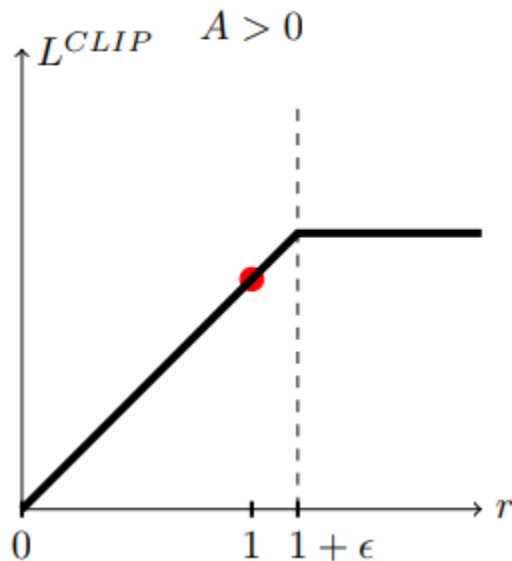


Proximal Policy Optimization(PPO)

Будем ограничивать шаги в сторону улучшения политики

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

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





algorithm	avg. normalized score
No clipping or penalty	-0.39
Clipping, $\epsilon = 0.1$	0.76
Clipping, $\epsilon = 0.2$	0.82
Clipping, $\epsilon = 0.3$	0.70
Adaptive KL $d_{\text{targ}} = 0.003$	0.68
Adaptive KL $d_{\text{targ}} = 0.01$	0.74
Adaptive KL $d_{\text{targ}} = 0.03$	0.71
Fixed KL, $\beta = 0.3$	0.62
Fixed KL, $\beta = 1.$	0.71
Fixed KL, $\beta = 3.$	0.72
Fixed KL, $\beta = 10.$	0.69

Table 1: Results from continuous control benchmark. Average normalized scores (over 21 runs of the algorithm, on 7 environments) for each algorithm / hyperparameter setting . β was initialized at 1.



TRPO

- Works for smaller models
- Second-order optimization

PPO

- Works for big models
- First-order optimization