Policy gradients methods, $\pi(a|s) = \pi(a|s,\Theta)$, $a \in \mathbb{R}^d$



$\mu_1, \sigma_1$ ... $\mu_d, \sigma_d$

$$\pi(a|s,\Theta) = \prod_{j=1}^{d} N(a_j | \mu_j(s,\Theta), \sigma_j^2(s,\Theta))$$

$$a \in \{1...k\}, \quad \pi(a|s,\Theta) = sm(output(s,\Theta))$$

$$F(\Theta) = \mathbb{E}_{p(s)} \mathbb{E}_{\pi(a|s,\Theta)} Q^{\pi_\Theta}(s,a) \longrightarrow \max_{\Theta}$$

$$p(\tau|\Theta) = p(s_0) \prod_{j=0}^{\infty} p(s_{j+1} | s_j, a_j) \pi(a_j | s_j, \Theta)$$

$$\nabla_\Theta F(\Theta) = \nabla_\Theta \overbrace{\mathbb{E}_{p(\tau|\Theta)}}^{?} f(\tau), \quad log-der-trick$$

$$\nabla_\Theta F(\Theta) = \nabla_\Theta \int p(\tau|\Theta) f(\tau) d\tau = \int \nabla_\Theta p(\tau|\Theta) f(\tau) d\tau =$$

$$= \left\{ \nabla_\Theta \log p(\tau|\Theta) = \frac{1}{p(\tau|\Theta)} \nabla_\Theta p(\tau|\Theta) \right\} = \int p(\tau|\Theta) \nabla_\Theta \log p(\tau|\Theta) f(\tau) d\tau =$$

$$= \mathbb{E}_{p(\tau|\Theta)} \nabla_\Theta \log p(\tau|\Theta) f(\tau)$$

$$\nabla_\Theta \log p(\tau|\Theta) = \nabla_\Theta \log \left( p(s_0) \prod_{t=0}^{\infty \nearrow T} p(s_{t+1}|s_t, a_t) \pi(a_t | s_t, \Theta) \right) =$$

$\supset \pi(s_t,\Theta) \leftarrow$ дет.
$\parallel \qquad \swarrow$ случ.

$$= \sum_{t=0}^{\infty \nearrow T} \nabla_\Theta \log \pi(a_t | s_t, \Theta)$$

## REINFORCE

Алгоритм

Иниц. $\Theta$

повторять

огромная дисперсия
из-за log-der-trick

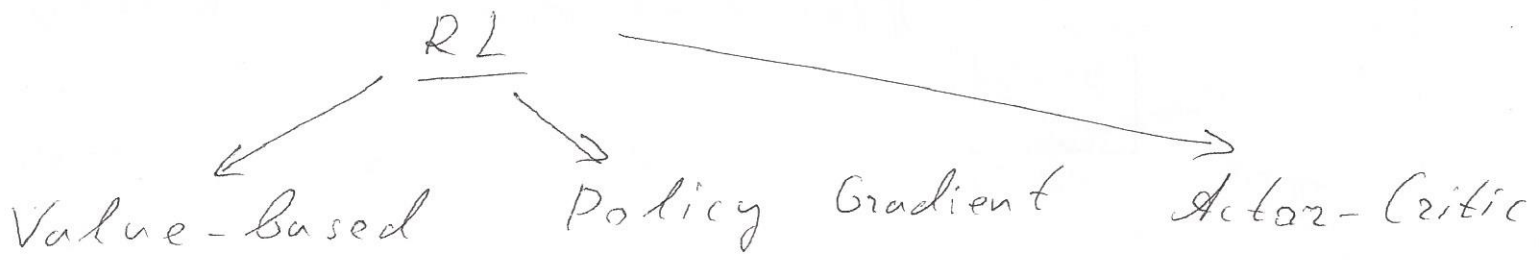$$\tau = \{s_0, a_0, s_1, a_1, ..., s_T\}$$

$$\nabla F(\Theta) = \left( \sum_{t=0}^{T} \nabla_\Theta \log \pi(a_t, s_t | \Theta) \right) \left( \sum_{t=0}^{T} \gamma^t z(s_t, a_t) \right)$$

$$\Theta \leftarrow \Theta + \alpha \nabla_\Theta F(\Theta)$$

$$\mathop{\mathbb{E}}_{p(\tau|\theta)} \nabla_\theta \log p(\tau|\theta) = \int p(\tau|\theta) \frac{1}{p(\tau|\theta)} \nabla_\theta p(\tau|\theta) \, d\tau = 0$$

Baseline $B$

$$RL$$

Value-based $\qquad$ Policy Gradient $\qquad$ Actor-Critic

Actor: $\pi(a|s,\theta)$

Critic: $Q(s,a|w) \approx Q^\pi(s,a)$

$$F(\theta) = \mathop{\mathbb{E}}_{p(s)} \mathop{\mathbb{E}}_{\pi(a|s,\theta)} Q(s,a|w) \xrightarrow[\theta]{} \max$$

$$Q^\pi(s,a) = z(s,a) + \gamma \mathop{\mathbb{E}}_{p(s'|s,a)} \mathop{\mathbb{E}}_{\pi(a'|s')} Q^\pi(s',a')$$

$$G(w) = \frac{1}{|S||A|} \sum_{s,a} \left( Q(s,a|w) - [z(s,a) + \gamma \mathop{\mathbb{E}}_{p(s'|s,a)} \mathop{\mathbb{E}}_{\pi(a'|s',\theta)} Q(s',a'|w)] \right)^2 \xrightarrow[w]{} \min$$

Алгоритм QAC

Иниц. $\theta, w$
повторять
сэмпл $(s,a,z,s')$ $\leftarrow$
$a' \sim \pi(a'|s',\theta)$
$y = z(s,a) + \gamma Q(s',a'|w)$
$\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi(a|s,\theta) Q(s,a|w)$
$w \leftarrow w - \beta \cdot 2 (Q(s,a|w) - y) \nabla_w Q(s,a|w)$

1) memory replay
2) parallel learning

Baseline: $\nabla_\theta F(\theta) = \nabla_\theta \mathop{\mathbb{E}}_{p(s)} \mathop{\mathbb{E}}_{\pi(a|s,\theta)} (Q(a,s|w) - B(s))$

$\nabla_\theta \mathop{\mathbb{E}}_{p(s)} \mathop{\mathbb{E}}_{\pi(a|s,\theta)} B(s) = \nabla_\theta \int p(s) \int \pi(a|s,\theta) B(s) \, da \, ds =$

$= \int p(s) B(s) \left( \nabla_\theta \underbrace{\int \pi(a|s,\theta) \, da}_{1} \right) ds = 0$

2

$$A^{\pi}(s,a) = Q^{\pi}(s,a) - \underbrace{V^{\pi}(s)}_{B(s)}$$

advantage

$$Q^{\pi}(s,a) = r(s,a) + \gamma \underset{p(s'|s,a)}{\mathbb{E}} V^{\pi}(s'')$$

## Схема A2C

Инициализация $\Theta, w$

повторять
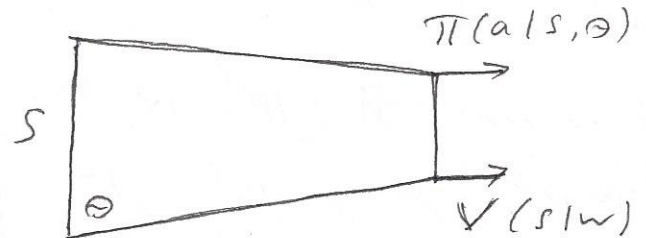
  сэмплим $(s, a, r, s')$

$$y = r(s,a) + \gamma V(s'|w)$$

$$A(s,a) = y - V(s|w)$$

$$\Theta \leftarrow \Theta - \alpha \nabla_{\Theta} \log \pi(a|s, \Theta) A(s,a)$$

$$w \leftarrow w - \beta 2(V(s|w) - y) \nabla_w V(s|w)$$



$\pi(a|s, \Theta)$

$s$

$V(s|w)$

## Процесс Дирихле

$$\begin{cases} f(x) \sim GP(0, k(\cdot, \cdot)) \\ y_i \mid f(x_i) \sim p(y|f) \end{cases}$$

$$\Rightarrow p(f|y, x)$$

$$p(y_{test}|x_{test}) = \int p(y_{test}|f(x_{test})) \cdot p(f|y,x)\,df$$

непараметрическая модель, зависит от числа объектов

$$\begin{cases} \pi \sim Dir(\pi|\alpha) \\ \Theta_1, \dots \Theta_k \sim p(\Theta) \\ z_1, \dots z_N \sim Discrete(z|\pi) \\ x_i \mid z_i, \Theta \sim p(x|\Theta_{z_i}) \end{cases}$$

$$p(X, z, \Theta, \pi) = \underbrace{p(\Theta)p(\pi)}p(z|\pi)p(x|z,\Theta)$$

непараметрическая

модель ?

(здесь нет, но нужно сделать)