

# Модельные методы обучения с подкреплением часть 2

Осминин Константин

13 августа 2019

**Tinkoff.ru**



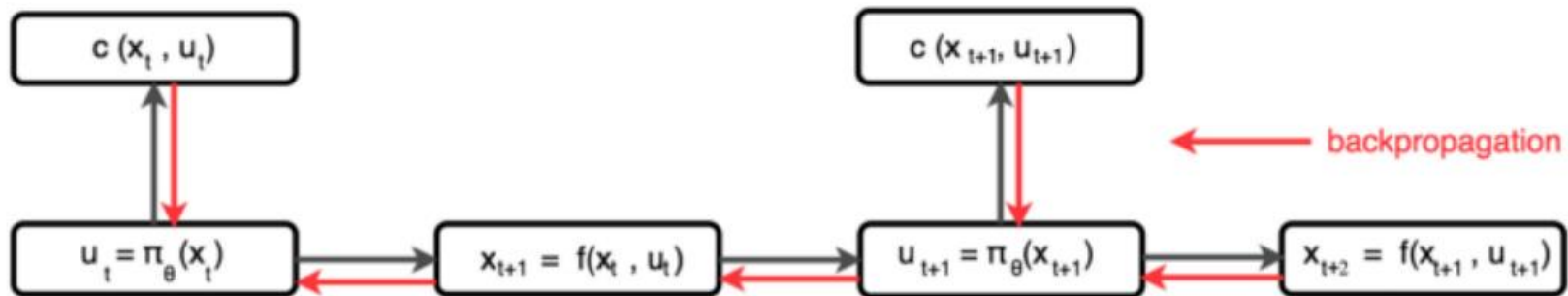
- ✓ Выучивание динамики среды
- ✓ Гауссовский процесс
- ✓ PILCO



| Обучение с подкреплением | Теория управления |
|--------------------------|-------------------|
| Состояние $s$            | Состояние $x$     |
| Действие $a$             | Действие $u$      |
| Вознаграждение $r$       | Кост $c$          |
| Агент                    | Контроллер        |



# Выучивание динамики среды



1. run base policy  $\pi_0(\mathbf{a}_t|\mathbf{s}_t)$  (e.g., random policy) to collect  $\mathcal{D} = \{(\mathbf{s}, \mathbf{a}, \mathbf{s}')_i\}$
2. learn dynamics model  $f(\mathbf{s}, \mathbf{a})$  to minimize  $\sum_i \|f(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{s}'_i\|^2$
3. backpropagate through  $f(\mathbf{s}, \mathbf{a})$  into the policy to optimize  $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$
4. run  $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ , appending the visited tuples  $(\mathbf{s}, \mathbf{a}, \mathbf{s}')$  to  $\mathcal{D}$



# Гауссовский процесс

- ✓ **Определение:** Случайный процесс  $\{f_x\}, x \in \mathbb{R}^n$  является гауссовским тогда и только тогда, когда для любого конечного множества индексов  $(x_1, \dots, x_m)$ :

$(f_{x_1}, \dots, f_{x_m})$  есть многомерная гауссовская случайная величина.

- ✓ Если  $\forall x: \mathbb{E} f_x = 0$ , то ГП полностью определяется  $\text{cov}(f_x, f_y)$ .

- ✓ Для аппроксимации ковариации используются ядра:

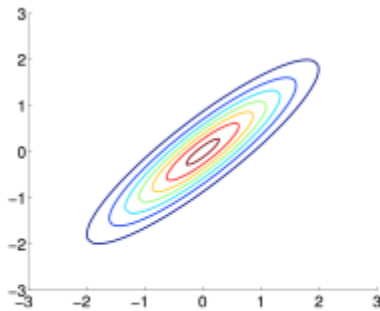
- Константа:  $K_C(x, x') = C$
- Линейная функция:  $K_L(x, x') = x^T x'$
- Гауссовский шум:  $K_{GN}(x, x') = \sigma^2 \delta_{x, x'}$
- Квадратичная экспоненциальная функция:  $K_{SE}(x, x') = \exp\left(-\frac{\|d\|^2}{2\ell^2}\right)$
- Функция Орнштейна-Уленбека:  $K_{OU}(x, x') = \exp\left(-\frac{|d|}{\ell}\right)$
- Периодическая функция:  $K_P(x, x') = \exp\left(-\frac{2 \sin^2\left(\frac{d}{2}\right)}{\ell^2}\right)$



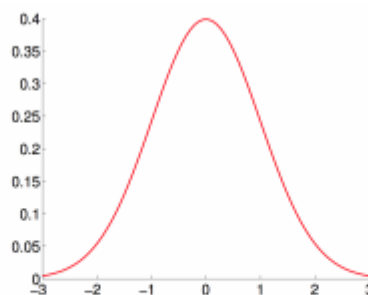
# Свойства гауссовского распределения

- ✓ **Факт:** Маргинальное и условное распределения гауссовского распределения также гауссовы.

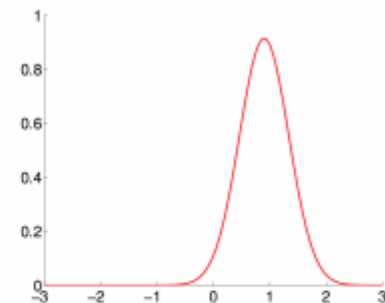
$$\begin{pmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{bmatrix} \right), \quad p(\mathbf{f}_1) = \int p(\mathbf{f}_1, \mathbf{f}_2) d\mathbf{f}_2 = \mathcal{N}(\mathbf{f}_1 | \mu_1, \Sigma_{11})$$
$$p(\mathbf{f}_1 | \mathbf{f}_2) = \mathcal{N}(\mathbf{f}_1 | \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{f}_2 - \mu_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12}^T)$$



$$p(\mathbf{f}_1, \mathbf{f}_2) \sim \mathcal{N}(\mathbf{f}_1, \mathbf{f}_2 | \mu, \Sigma)$$



$$p(\mathbf{f}_1)$$



$$p(\mathbf{f}_1 | \mathbf{f}_2)$$

Source: [DeepBayes2018 Burnaev GP lecture](#)



# Регрессия гауссовского процесса

✓ **Задача регрессии:** Дано  $(\mathbf{X}, \mathbf{y}) = \{(x_i, y_i)\}_{i=1}^N$ ,  $x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$ .

Нужно найти  $y^*$  по  $x^*$ .

✓ **Предположение:**  $y = f(x) + \varepsilon$ , где  $f(x)$  - гауссовский процесс со средним 0 и вариацией  $k(x, x')$ .  $\varepsilon \sim \mathcal{N}(0, \sigma_n^2)$ .

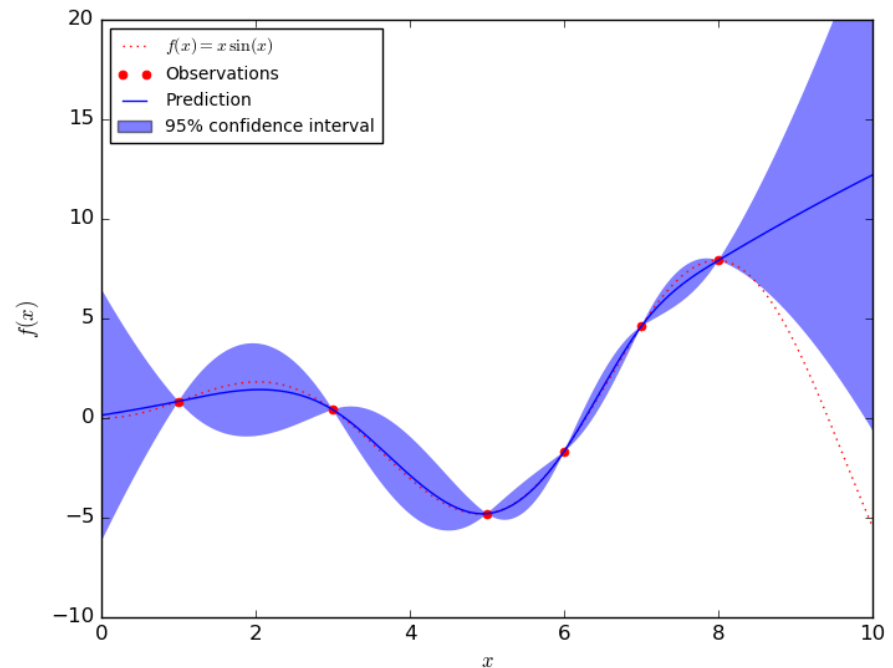
✓ Тогда  $y^* | \mathbf{X}, \mathbf{y}, x^* \sim \mathcal{N}(m(x^*), \sigma(x^*))$ .

✓  $m(x^*) = \mathbf{k}^T \mathbf{K}_y^{-1} \mathbf{y}$ ,

✓  $\sigma(x^*) = k(x^*, x^*) - \mathbf{k}^T \mathbf{K}_y^{-1} \mathbf{k}$

✓  $\mathbf{k} = (k(x^*, x_1), \dots, k(x^*, x_N))$ ,

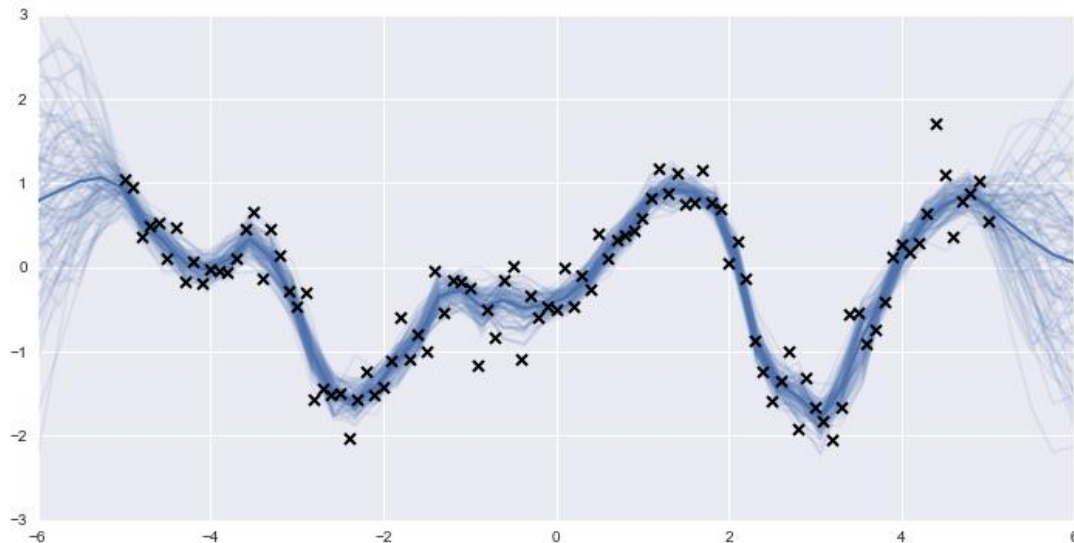
✓  $\mathbf{K}_y = \left\| k(x_i, x_j) \right\|_{i,j=1}^N + \sigma_n^2 \mathbf{I}$





# Свойства регрессии ГП

- ✓ Линейна по  $y$ .
- ✓ Условна – нет предположений о распределении  $X$ .
- ✓ Эффективна при малых  $N$  и неэффективна при больших.
- ✓ Дает не только оценку значения  $y^*$ , но и степень неуверенности модели  $\sigma(x^*)$ .
- ✓ (Почти) непараметрична.







- ✓ PILCO (2011) – probabilistic inference for learning control.
- ✓ Моделирует динамику среды гауссовским процессом.
  - ✓ Пусть  $x_{t-1}, u_{t-1}, x_t, \dots$  - траектория.
  - ✓ Тогда  $\Delta_t = x_t - x_{t-1} + \varepsilon$  – приращение с шумом  $\varepsilon \sim \mathcal{N}(0, \sigma_n)$ .
  - ✓ Обозначим  $\tilde{x} = [x_t \ u_t]$ .
  - ✓ Считаем, что  $\Delta_t$  - гауссовский процесс от  $\tilde{x}$  с

$$k(\tilde{x}, \tilde{x}') = \alpha^2 \exp(-0.5 (\tilde{x} - \tilde{x}')^T \Lambda^{-1} (\tilde{x} - \tilde{x}')), \quad \Lambda = \text{diag}(l_1^2, \dots, l_D^2).$$

- ✓ Параметры гауссовского процесса  $\phi = [\sigma_i, \alpha, l_i]$  определяются из максимизации правдоподобия на обучающей выборке  $\tilde{X}$ :

$$\log p(f(x) | \phi, x) = -\frac{1}{2} f(x)^T K(\phi, x, x')^{-1} f(x) - \frac{1}{2} \log \det(K(\phi, x, x')) - \frac{N}{2} \log 2\pi$$



# Probabilistic inference

Дано  $\tilde{x}_{t-1} \sim \mathcal{N}(\mu'_{t-1}, \Sigma'_{t-1})$ . Требуется найти аппроксимацию  $x_t \sim \mathcal{N}(\mu_t, \Sigma_t)$ .

$$\begin{aligned} \mu_\Delta &= \mathbb{E}_{\tilde{x}_{t-1}} m(\tilde{x}_{t-1}) = \int m(\tilde{x}_{t-1}) \cdot \mathbf{f}_{\mathcal{N}(\mu'_{t-1}, \Sigma'_{t-1})}(\tilde{x}_{t-1}) d\tilde{x}_{t-1} = \\ &= [\mathbf{k} = k(\tilde{\mathbf{X}}, \tilde{x}_{t-1}), \quad \tilde{\mathbf{X}} - \text{обуч. в - ка}, \quad \mathbf{K}_y = k(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \sigma_n^2 \mathbf{I}, \quad \mathbf{y} = \Delta_X] = \\ &= \int \mathbf{k}^T \mathbf{K}_y^{-1} \mathbf{y} \cdot \mathbf{f}_{\mathcal{N}(\mu'_{t-1}, \Sigma'_{t-1})}(\tilde{x}_{t-1}) d\tilde{x}_{t-1} = \boldsymbol{\beta}^T \cdot \mathbf{q} \end{aligned}$$

$$\begin{aligned} \boldsymbol{\beta} &= \mathbf{K}_y^{-1} \mathbf{y}, \quad \mathbf{q}_i = \int k(\tilde{x}_i, \tilde{x}_{t-1}) \cdot \mathbf{f}_{\mathcal{N}(\mu'_{t-1}, \Sigma'_{t-1})}(\tilde{x}_{t-1}) d\tilde{x}_{t-1} = \\ &= \frac{\alpha^2}{\sqrt{\Sigma_{t-1} \Lambda_i^{-1} + I}} \exp(-0.5(\tilde{x}_i - \mu'_{t-1})^T (\Sigma_{t-1} + \Lambda_i)^{-1} (\tilde{x}_i - \mu'_{t-1})) \end{aligned}$$

$$\mu_t = \mu_{t-1} + \mu_\Delta$$

$\Sigma_t$  найти посложнее, но тоже можно.



# Целевая функция

- ✓  $J(\pi) = \sum_t \mathbb{E}_{x_t} c(x_t)$ , при условии
$$x_t - x_{t-1} \sim \text{Gauss Proc}(x_{t-1}, \pi_\theta(x_{t-1})), \quad x_0 \sim \mathcal{N}(\mu_0, \sigma_0)$$
- ✓ Считаем  $c(x) = 1 - \exp(-|x - x_{\text{target}}|^2 / \sigma_c^2)$
- ✓ Контроллер  $\pi_\theta(x)$  - произвольный.  $\theta = \operatorname{argmin} J(\pi)$ .
- ✓  $\mathbb{E}_{x_t} c(x_t) = \int c(x_t) \cdot f_{\mathcal{N}(\mu_t, \Sigma_t)}(x_t) dx_t$
- ✓ 
$$\frac{d\mathbb{E}_{x_t} c(x_t)}{d\theta} = \frac{\partial \mathbb{E}_{x_t} c(x_t)}{\partial \mu_t} \cdot \frac{d\mu_t}{d\theta} + \frac{\partial \mathbb{E}_{x_t} c(x_t)}{\partial \Sigma_t} \cdot \frac{d\Sigma_t}{d\theta}$$
- ✓ 
$$\frac{d\mu_t}{d\theta} = \frac{\partial \mu_t}{\partial \mu_{t-1}} \frac{d\mu_{t-1}}{d\theta} + \frac{\partial \mu_t}{\partial \Sigma_{t-1}} \frac{d\Sigma_{t-1}}{d\theta} + \frac{\partial \mu_t}{\partial \theta}$$
- ✓ 
$$\frac{\partial \mu_t}{\partial \theta} = \frac{\partial \mu_\Delta}{\partial \mu_u} \frac{d\mu_u}{d\theta} + \frac{\partial \mu_\Delta}{\partial \Sigma_u} \frac{d\Sigma_u}{d\theta}$$



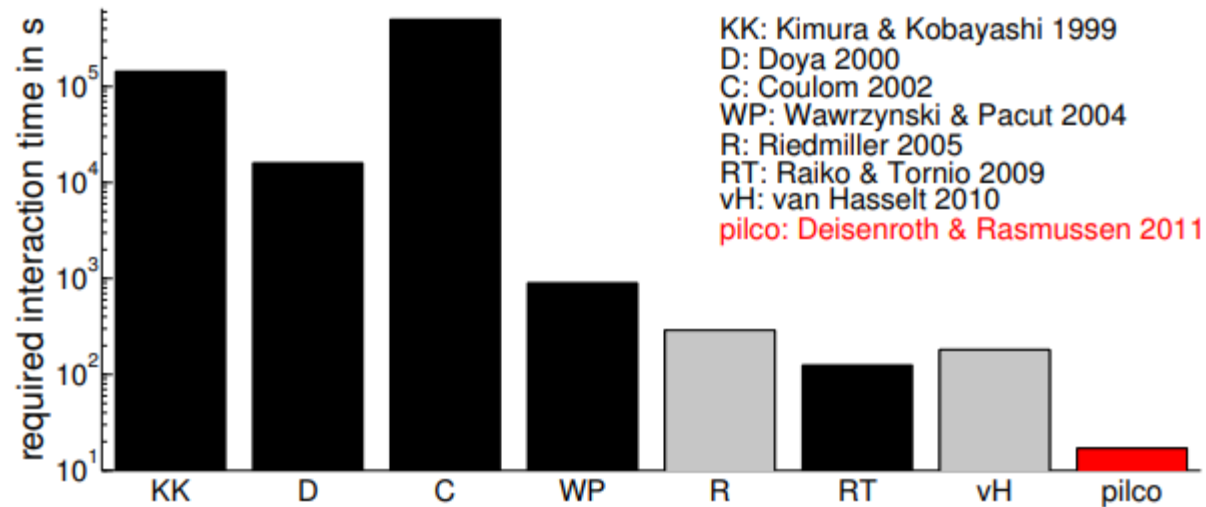
# PILCO алгоритм

---

**Algorithm 1** PILCO

---

- 1: **init:** Sample controller parameters  $\theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .  
Apply random control signals and record data.
- 2: **repeat**
- 3:   Learn probabilistic (GP) dynamics model,
- 4:   Model-based policy search,
- 5:   **repeat**
- 6:     Approximate inference for policy evaluation,  
      get  $J^\pi(\theta)$ , get  $dJ^\pi(\theta)/d\theta$
- 7:     Gradient-based policy improvement,
- 8:     Update parameters  $\theta$  (e.g., CG ).
- 9:   **until** convergence; **return**  $\theta^*$
- 10:   Set  $\pi^* \leftarrow \pi(\theta^*)$ .
- 11:   Apply  $\pi^*$  to system (single trial/episode) and  
      record data.
- 12: **until** task learned



*Figure 5.* Data efficiency for learning the cart-pole task in the absence of expert knowledge. The horizontal axis chronologically orders the references according to their publication date. The vertical axis shows the required interaction time with the cart-pole system on a log-scale.



- ✓ Гауссовский процесс
  - ✓ [ВИКИ](#)
  - ✓ [Лекция DeepBayes2018](#)
  
- ✓ PILCO
  - ✓ [статья 2011](#)
  - ✓ [сайт](#)



Регрессия гауссовского процесса

[https://github.com/bayesgroup/deepbayes-2018/blob/master/day5\\_gp/gp\\_basic.ipynb](https://github.com/bayesgroup/deepbayes-2018/blob/master/day5_gp/gp_basic.ipynb)



Спасибо