# Задача стохастической оптимизации

$$\min_{x \in \mathbb{R}^d} \left[ f(x) := \mathbb{E}_{\xi \sim D} \left[ f(x, \xi) \right] \right]$$

Пример: в ML

$$f(x) := \mathbb{E}_{\xi \sim D} \left[ L\left( \underset{\text{модель}}{g}(\overset{\text{вес}}{x}, \xi_a), \xi_b \right) \right]$$

$(\xi_a, \xi_b)$ ←

природа (неизв.) ←

$\Downarrow$

Вычислить $f$, $\nabla f$ невозможно

Что делать?

1) Online подход:

$$\nabla f(x, \xi) = \nabla_x L(g(x, \xi_a), \xi_b) \quad \leftarrow \text{ градиент по отдел. случаю}$$

$$\mathbb{E}_{\xi \sim D} \left[ \nabla f(x, \xi) \right] = \nabla f(x)$$

2) Offline подход:
есть выборка $\{ \xi_i \}_{i=1}^{n}$

Аппроксимация $\mathbb{E}_{\xi \sim D}$ по Монте-Карло

$$\min_{x \in \mathbb{R}^d} \left[ \tilde{f}(x) := \frac{1}{n} \sum_{i=1}^{n} L(g(x, \xi_{i,a}), \xi_{i,b}) \right]$$

это другая задача: $\tilde{f} \approx f$ при больших $n$

можно считать $\nabla \tilde{f}$:

$\ominus$ дорого

$\ominus$ переобучение ($\tilde{f} \neq f$)

$\oplus$ вычислять не полный град., а град. по части выборки — Batch

# Стохастический градиентный спуск

$$x^{k+1} = x^k - \gamma \, \nabla f(x^k, \xi^k)$$

<span style="color:blue">стох. град.</span>

<span style="color:blue">по 1ому объекту выборки</span>

$$\xi^k - \text{независимо} \quad \text{и} \quad \text{равномерно}$$

- Независимость:

$$\mathbb{E}_\xi \left[ \nabla f(x^k, \xi^k) \right]$$

- Равномерность:

$$\mathbb{E}_\xi \left[ \nabla f(x^k, \xi^k) \right] = \text{(определение матожидания)}$$

$$= \sum_{i=1}^{n} \mathbb{P}\{\xi^k = \xi_i\} \, \nabla f(x^k, \xi_i) = \text{равномерность}$$

$$= \sum_{i=1}^{n} \frac{1}{n} \, \nabla f(x^k, \xi_i) = \frac{1}{n} \sum_{i=1}^{n} \nabla f(x^k, \xi_i) = \nabla f(x^k)$$

---

# Условное математическое ожидание

$$\mathbb{E}[\cdot \mid x^k] = \mathbb{E}[\cdot \mid \mathcal{F}_k]$$

$$\mathcal{F}_k - \sigma\text{-алгебра, порожд. } x^0, \xi^0, \ldots \xi^{k-1}$$

<span style="color:blue">фиксируем всё, что произошло до $x^k$ (включительно)</span>

tower property:

$$\mathbb{E}\left[ \mathbb{E}[X \mid Y] \right] = \mathbb{E}[X]$$

# Предположения

- $f$ — $\mu$ — сильно выпукла
- $f(\cdot, \xi)$ — $L$ — гладкая (можно $L = L_{max}$ по всем)
- $\mathbb{E}_{\xi}\left[\nabla f(x, \xi)\right] = \nabla f(x)$    $\mathbb{E}_{\xi}\left[f(x, \xi)\right] = f(x)$

## Док-во сходимости:

$$\|x^{(k+1)} - x^*\|_2^2 = \|x^k - \gamma \nabla f(x^k, \xi^k) - x^*\|_2^2$$
$$= \|x^k - x^*\|_2^2 - 2\gamma \langle \nabla f(x^k, \xi^k); x^k - x^* \rangle$$
$$+ \gamma^2 \|\nabla f(x^k, \xi^k)\|_2^2$$

Полное м.о. от обеих частей:
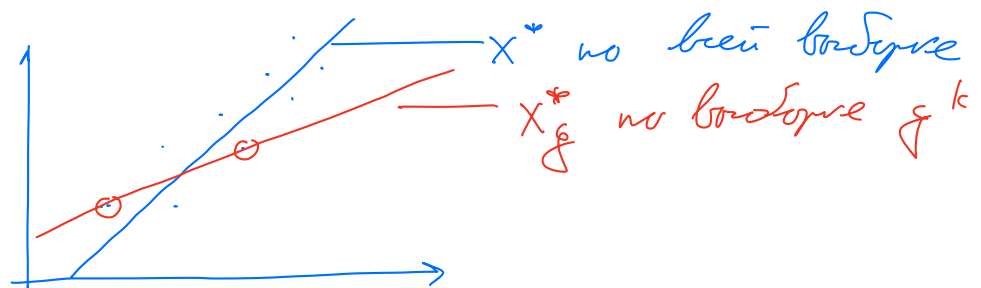
$$\mathbb{E}\left[\|x^{(k+1)} - x^*\|_2^2\right] = \mathbb{E}\left[\|x^k - x^*\|_2^2\right] - 2\gamma \mathbb{E}\left[\langle \nabla f(x^k, \xi^k); x^k - x^* \rangle\right]$$
$$+ \gamma^2 \mathbb{E}\left[\|\nabla f(x^k, \xi^k)\|_2^2\right] \qquad (\ast)$$

$\mathbb{E}\left[\|\nabla f(x^k, \xi^k)\|_2^2\right]$:

$$\mathbb{E}\left[\|\nabla f(x^k, \xi^k)\|_2^2\right] = \mathbb{E}\left[\|\underbrace{\nabla f(x^k, \xi^k) - \nabla f(x^*, \xi^k)}_{\neq 0} + \underbrace{\nabla f(x^*, \xi^k)}_{\neq 0}\|_2^2\right]$$

$$\nabla f(x^*) = 0 \quad \not\Rightarrow \quad \nabla f(x^*, \xi^k) = 0$$



$x^*$ по всей выборке

$x_{\xi}^*$ по выборке $\xi^k$

КБШ

$$\leq 2\mathbb{E}\left[\|\nabla f(x^k, \xi^k) - \nabla f(x^*, \xi^k)\|_2^2\right] + 2\mathbb{E}\left[\|\nabla f(x^*, \xi^k)\|_2^2\right]$$

$$\leq 4L\, \mathbb{E}\left[ f(x^t, \xi^k) - f(x^*, \xi^k) + \langle \nabla f(x^*, \xi^k); x^* - x^t \rangle \right]$$
$$+ 2\,\mathbb{E}\left[ \|\nabla f(x^*, \xi^k)\|_2^2 \right] \ominus \quad \neq 0$$

Tower property: $\quad \mathbb{E}\left[\mathbb{E}\left[\ \cdot\ | x^k\right]\right] = \mathbb{E}\left[\ \right]$

$$\mathbb{E}\left[\mathbb{E}\left[ f(x^t, \xi^k) | x^k \right]\right] = \mathbb{E}\left[ f(x^k) \right]$$

$$\mathbb{E}\left[ f(x^*, \xi^k) \right] = f(x^*)$$

$$\mathbb{E}\left[\mathbb{E}\left[ \langle \nabla f(x^t, \xi^k); x^* - x^k \rangle | x^k \right]\right] =$$

$$= \mathbb{E}\left\langle \mathbb{E}\left[ \nabla f(x^*, \xi^k) | x^k \right]; x^* - x^t \right\rangle$$

$$= \mathbb{E}\left[\langle \nabla f(x^*); x^* - x^k \rangle\right] = 0$$

следует
$$\ominus \quad 4L\, \mathbb{E}\left[ f(x^t) - f(x^*) \right] + 2\mathbb{E}\left[ \|\nabla f(x^*, \xi^k)\|_2^2 \right]$$

$$\mathbb{E}\left[ \|\nabla f(x^*, \xi^k)\|_2^2 \right] = \text{равномерность}$$
$$= \frac{1}{n}\sum_{i=1}^{n} \|\nabla f(x^*, \xi_i)\|_2^2 \ :\!= \sigma_*^2$$

$$= 4L\, \mathbb{E}\left[ f(x^t) - f(x^*) \right] + 2\sigma_*^2 \qquad (**)$$

Подставим (★★) в (★)

$$\mathbb{E}\left[\|x^{k+1}-x^*\|_2^2\right] \le \mathbb{E}\left[\|x^k-x^*\|_2^2\right] - 2\gamma\,\mathbb{E}\left[\langle \nabla f(x^k,\xi^k); x^k-x^*\rangle\right]$$
$$+ \gamma^2\left(4L\,\mathbb{E}\left[f(x^k)-f(x^*)\right] + 2\sigma_*^2\right) \qquad (★★★)$$

$\mathbb{E}\left[\langle \nabla f(x^k,\xi^k); x^k-x^*\rangle\right]$, tower property:

$$\mathbb{E}\left[\mathbb{E}\left[\langle \nabla f(x^k,\xi^k); x^k-x^*\rangle \mid x^k\right]\right]$$

$$= \mathbb{E}\left[\langle \mathbb{E}\left[\nabla f(x^k,\xi^k)\mid x^k\right]; x^k-x^*\rangle\right]$$

$$= \mathbb{E}\left[\langle \nabla f(x^k); x^k-x^*\rangle\right] \qquad (★★★★)$$

Подставляем (★★★★) в (★★★)

$$\mathbb{E}\left[\|x^{k+1}-x^*\|_2^2\right] \le \mathbb{E}\left[\|x^k-x^*\|_2^2\right] - 2\gamma\,\mathbb{E}\left[\langle \nabla f(x^k); x^k-x^*\rangle\right]$$
$$+ \gamma^2\left(4L\,\mathbb{E}\left[f(x^k)-f(x^*)\right] + 2\sigma_*^2\right)$$

$\mu$-сильная выпуклость для $f$

$$\mathbb{E}\left[\|x^{k+1}-x^*\|_2^2\right] \le \mathbb{E}\left[\|x^k-x^*\|_2^2\right] - 2\gamma\,\mathbb{E}\left[\frac{\mu}{2}\|x^k-x^*\|_2^2 + f(x^k)-f(x^*)\right]$$
$$+ \gamma^2\left(4L\,\mathbb{E}\left[f(x^k)-f(x^*)\right] + 2\sigma_*^2\right)$$

$$= (1-\gamma\mu)\,\mathbb{E}\left[\|x^k-x^*\|_2^2\right]$$
$$- 2\gamma(1-2\gamma L)\,\underbrace{\mathbb{E}\left[f(x^k)-f(x^*)\right]}_{\ge 0}$$
$$+ 2\gamma^2\sigma_*^2$$

$$\gamma \le \frac{1}{2L}$$

$$\le (1-\gamma\mu)\,\mathbb{E}\left[\|x^{k}-x^{*}\|_2^2\right] + 2\gamma^2\sigma_*^2$$

$$\boxed{\mathbb{E}\left[\|x^{k+1}-x^{*}\|_2^2\right] \le (1-\gamma\mu)\,\mathbb{E}\left[\|x^{k}-x^{*}\|_2^2\right] + 2\gamma^2\sigma_*^2}$$

то же самое было в GD, но без $\mathbb{E}$ и без

$$R_k^2 = \mathbb{E}\left[\|x^{k}-x^{*}\|_2^2\right]$$

$$R_{k+1}^2 \le (1-\gamma\mu)\,R_k^2 + 2\gamma^2\sigma_*^2$$

Запустим рекурсию:

$$\le (1-\gamma\mu)\left((1-\gamma\mu)R_{k-1}^2 + 2\gamma^2\sigma_*^2\right) + 2\gamma^2\sigma_*^2$$

$$= (1-\gamma\mu)^2\,R_{k-1}^2 + 2\gamma^2\sigma_*^2\left[1 + (1-\gamma\mu)\right]$$

$$\cdots$$

$$\le (1-\gamma\mu)^{k+1}\,R_0^2 + 2\gamma^2\sigma_*^2\sum_{i=0}^{k}(1-\gamma\mu)^i$$

$$\le (1-\gamma\mu)^{k+1}\,R_0^2 + 2\gamma^2\sigma_*^2\sum_{i=0}^{\infty}(1-\gamma\mu)^i \underbrace{\qquad}_{=\frac{1}{\gamma\mu}}$$
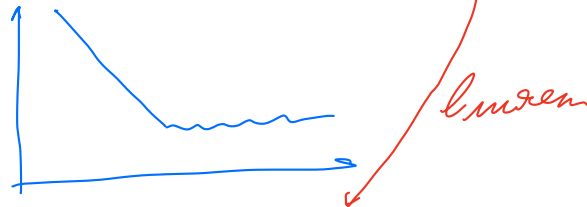
$$\le (1-\gamma\mu)^{k+1}\,R_0^2 + \frac{2\gamma\sigma_*^2}{\mu}$$

# Сходимость SGD с постоян. шагом

$$\mathbb{E}\left[\|X^{k+1}-X^*\|_2^2\right] \leq (1-\gamma\mu)^{k+1}\mathbb{E}\left[\|X^0-X^*\|_2^2\right] + \frac{2\gamma\sigma_*^2}{\mu}$$

⊕ сходимость линейная, как у GD

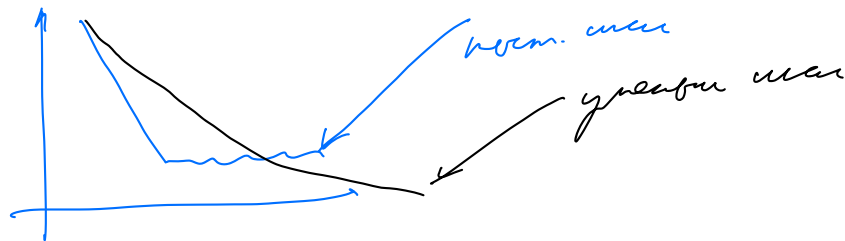⊖ сходимость до окрестности: $\sim\gamma$, $\sim\sigma_*^2$



влияет

⊕ простота + ML-метрики

---

# Как бороться со сходимостью к окрестности?

1) $\gamma$ уменьшить : ⊕ сходимость глубже
   ⊖ медленее

$\gamma \to \gamma_k \sim \frac{1}{k} ; \frac{1}{\sqrt{k}}$ : ⊕ сходимость до решения
   ⊖ вообще пропадает или слож.



пост. шаг

уменьш шаг

2) $\sigma_*^2$ уменьшить

$$\nabla f(X^t, g^k) \to \frac{1}{b}\sum_{g\in S^k}\nabla f(X^t, g)$$

$S^k$ - батч, набор объектов из об. выборки, размера $b$
(все объекты берутся независимо и равномерно)

$$\bullet \; \mathbb{E}_{S^k}\left[ \frac{1}{b} \sum_{\xi \in S^k} \triangledown f(x^k, \xi) \right] =$$

$$= \frac{1}{b} \sum_{\xi \in S^k} \mathbb{E}_{\xi}\left[ \triangledown f(x^k, \xi) \right] = \text{(нез. и одинак.)}$$

$$= \frac{1}{b} \sum_{\xi \in S^k} \triangledown f(x^k) = b \cdot \frac{1}{b} \triangledown f(x^k) = \triangledown f(x^k)$$

$$\bullet \; \mathbb{E}_{S^k}\left[ \left\| \frac{1}{b} \sum_{\xi \in S^k} \triangledown f(x^*, \xi) \right\|_2^2 \right] =$$

$$= \mathbb{E}_{S^k}\left[ \frac{1}{b^2} \sum_{\xi \in S^k} \left\| \triangledown f(x^*, \xi) \right\|_2^2 \right]$$

$$+ \mathbb{E}_{S^k}\left[ \frac{1}{b^2} \sum_{\xi \neq \eta} \langle \triangledown f(x^*, \xi); \triangledown f(x^*, \eta) \rangle \right]$$

$$= \frac{1}{b^2} \sum_{\xi \in S^k} \mathbb{E}_{\xi}\left[ \left\| \triangledown f(x^*, \xi) \right\|_2^2 \right] = \frac{1}{b^2} \cdot b \cdot \underbrace{\frac{1}{n} \sum_{i=1}^{n} \left\| \triangledown f(x^*, \xi_i) \right\|_2^2}_{\sigma_{\phi}^2}$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\frac{\sigma^2_{\phi}}{b}}$$

$$+ \frac{1}{b^2} \sum_{\xi \neq \eta} \mathbb{E}_{\xi, \eta}\left[ \langle \triangledown f(x^*, \xi); \triangledown f(x^*, \eta) \rangle \right]$$

$$\xi, \eta - \text{независимые}$$

$$\langle \mathbb{E}_{\xi}\left[ \triangledown f(x^*, \xi) \right]; \mathbb{E}_{\eta}\left[ \triangledown f(x^*, \eta) \right] \rangle$$

$$\|$$

$$\langle \triangledown f(x^*); \triangledown f(x^*) \rangle = 0$$

$$= \boxed{\frac{\sigma^2_{\phi}}{\textcircled{b}}!} \;\leftarrow\; \text{Эффект батчирования}$$

$\oplus$ окрестность увеличилась в $b$ раз

$\ominus$ стоимость расчет

На практике:

- $\triangleright f_1 \rightarrow \triangleright f_2 \rightarrow \triangleright f_3 \dots \rightarrow \triangleright f_n$

не стох. метод.

- зашикловать данные $\rightarrow \triangleright f_1 \dots \triangleright f_n$ более равномерно 1 раз

- шикловать каждую эпоху $\rightarrow$ еще меньше следствии

механизм Shuffling