# Задача стохастической оптимизации

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

## Пример из ML

ф. потерь     веса

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

$(\xi_a, \xi_b)$    природа    модель    объект    метка

(неизв)

$f, \triangledown f$ вычислить нельзя

## Что делать?

1) Онлайн метод

$$\triangledown_x f(x, \xi) = \triangledown_x L\left( g(x, \xi_a), \xi_b \right)$$   градиент по тоже знаем

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

2) Оффлайн подход

есть выборка $\{\xi_i\}_{i=1}^n$

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

$\tilde{f}$ — Монте-Карло оценка исходной ф. $f$

$\tilde{f} \neq f$ — другая задача

$$\tilde{f} \approx f \text{ — при большом } n \text{ (аппроксимация}$$
<span style="color:blue">выпуска)</span>
$$\sim \frac{1}{\sqrt{n}}$$

можно считать $\triangledown \tilde{f}$:

$\ominus$ <span style="color:red">дорого</span>

$\ominus$ <span style="color:red">переобучение ($\tilde{f} \neq f$)</span>

$\oplus$ <span style="color:green">вычислять не $\triangledown f$, а град. по части выборки</span>

## Стохастический градиентный спуск

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

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

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

$\xi^k$ — независимо из $D/$ равномерно

- независимость

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

- равномерно

$$\mathbb{E}_\xi \left[ \triangledown f(x, \xi) \right] = \text{(оправдаем постановку)}$$

<span style="color:blue">$\tilde{f}$ для простоты $\tilde{f} \Rightarrow f$</span>

$$= \sum_{i=1}^{n} \mathbb{P}\{\xi = \xi_i\} \triangledown f(x, \xi_i)$$

$$\underset{\frac{1}{n}}{}$$

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

## Условное математическое ожидание

$$\mathbb{E}[\,\cdot\,|X^{k}] = \mathbb{E}[\,\cdot\,|\mathcal{F}_k]$$

$$\mathcal{F}_k = \text{б-алгебра, порожд. } X^0, \xi^0, \xi^1 \ldots \xi^{k-1}$$

фиксирует всю случайность, которая
произошла до $X^k$ (включительно)

tower property:

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

---

## Док-во сходимости:

- $f - \mu$-сильно выпуклая
- $f(\cdot, \xi) - L$-гладкая ($L_{max}$ по выборке)
- $\mathbb{E}_\xi[\triangledown f(x,\xi)] = \triangledown f(x) \qquad \mathbb{E}_\xi[f(x,\xi)] = f(x)$

### Док-во:

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

$$= \|x^k - x^*\|_2^2 - 2\gamma\langle\triangledown f(x^k, \xi^k); x^k - x^*\rangle$$

$$+ \gamma^2\|\triangledown f(x^k, \xi^k)\|_2^2$$

Полное м.о. от обеих частей:

$$\mathbb{E}\left[\|x^{k+1} - x^*\|_2^2\right] = \mathbb{E}\left[\|x^k - x^*\|_2^2\right]$$

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

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

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



$x^*$ ( не
        был подобран
        $\nabla f(x^*)$ )

$\nabla f(x^*, \xi) \neq 0$

$x^*_\xi$

$\neq 0$

к.э.ш.: $\|a+b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$

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

$L$ - гладкость $f(\cdot, \xi)$

$$\leq \quad 4L\,\mathbb{E}\left[f(x^k, \xi^k) - f(x^*, \xi^k) + \langle \nabla f(x^*, \xi^k), x^* - x^k \rangle\right]$$
$$+ 2\,\mathbb{E}\left[\|\nabla f(x^*, \xi^k)\|_2^2\right] \longleftarrow \sigma_*^2 \quad \text{(}\leq\text{)}$$

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

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

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

$$= \mathbb{E}\left[ \left\langle \underbrace{\mathbb{E}[\triangledown f(x^*, \xi^k) | x^k]}_{\triangledown f(x^*) = 0} ; x^k - x^* \right\rangle \right] = 0$$

$$\sigma_*^2 = \mathbb{E}_\xi \left[ \| \triangledown f(x^*, \xi) \|_2^2 \right]$$

$$= (\text{одределение}) = \frac{1}{n} \sum_{i=1}^n \| \triangledown f(x^*, \xi_i) \|_2^2$$

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

**Используем** $(**)$ **в** $(*)$

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

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

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

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

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

Подставим $(\divideontimes\divideontimes\divideontimes\divideontimes)$ в $(\divideontimes\divideontimes\divideontimes)$

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

$\mu$-сильной выпуклости

$$\mathbb{E}\left[\|x^{k+1}-x^*\|_2^2\right] \leq \mathbb{E}\left[\underline{\|x^k-x^*\|_2^2}\right]$$
$$-2\gamma\,\mathbb{E}\left[\underset{\sim}{f(x^k)}-f(x^*)+\frac{\mu}{2}\underline{\|x^k-x^*\|_2^2}\right]$$
$$+4L\gamma^2\mathbb{E}\left[\underset{\sim}{f(x^k)}-f(x^*)\right]+2\gamma^2\sigma_*^2$$
$$\leq (1-\gamma\mu)\,\mathbb{E}\left[\|x^k-x^*\|_2^2\right]$$
$$-2\gamma\underbrace{(1-2\gamma L)}_{\geq 0}\,\mathbb{E}\left[\underbrace{f(x^k)-f(x^*)}_{\geq 0}\right]$$
$$+2\gamma^2\sigma_*^2$$

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

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

то же самое, что и для GD

не было

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

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

Запустим рекурсию:

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

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

$$\cdots$$

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

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

$$\leq (1-\gamma\mu)^{k+1}R_0^2 + \frac{\gamma\sigma_*^2}{\mu}$$

Сходимость SGD с постоянным шагом

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

⊕ сходимость линейна, как у GD

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


гр. потерь

⊕ простота + хорошие ML-метрики

# Как бороться со сходимостью к окрестности?

## 1) $\gamma$ уменьшать:

⊕ скорость лучше
⊖ медленее

$$\gamma \to \gamma_k \sim \frac{1}{k}; \frac{1}{\sqrt{k}} \quad \text{без линейной сход.}$$



ф. потерь

к-пост. шаг

уменьш. шаг

$k$

## 2) $\sigma_*^2$ уменьшить

$$\nabla f(x, \delta) \to \frac{1}{b} \sum_{\delta \in S} \nabla f(x, \delta)$$

$S$ — батч (выбор объектов и меток) размера $b$

$$\sigma_*^2 \to \frac{\sigma_*^2}{b} \quad \text{Эффект батчирования}$$

⊕ окрестность уменьшится
⊖ стоимость вычислений растет