$$\min_{x \in \mathbb{R}^d} f(x)$$

$$\Downarrow$$

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

Пример:

$$\min_{x \in \mathbb{R}^d} \mathbb{E}_{\xi \sim D}\left[l\left(g(x, \xi_a), \xi_b\right)\right] \qquad (1)$$

ф. потерь ← (over $l$)
объект ← (over $\xi_a$ area)
метка ← (over $\xi_b$ area)
$(\xi_a, \xi_b)$ ← (under $\xi \sim D$)
модель ← (under $g$)
веса модели ← (under $x$)

не знаем (красным, под $\xi \sim D$)

Можем считать не $\nabla f(x)$, а $\nabla_x f(x, \xi)$

$$\nabla_x\left(l\left(g(x, \xi_a), \xi_b\right)\right)$$

Предположение:

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

Часто имеем дело с ограниченной выборкой:

$$\min_{x \in \mathbb{R}^d}\left[\frac{1}{n}\sum_{i=1}^{n} l\left(g(x, \xi_a^i), \xi_b^i\right)\right] \qquad (2)$$

у нас есть $n$ объектов с извест. ответами

(2) Монте - Карло приближение (1)

для (2) считать $\nabla f(x)$, но не хочу

- дорого
- добавить случайности/колебаний для того, чтобы быть ближе к (1)

Поэтому будем использовать $\nabla f(x, \xi^i)$

Если $i$ выбирается равномерно, то

$$\mathbb{E}_i \left[ \nabla f(x, \xi^i) \right] = \sum_{j=1}^{n} \mathbb{P}\{\xi_i = j\} \nabla f(x, \xi^j)$$

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

$$= \frac{1}{n} \sum_{j=1}^{n} \nabla_x \left[ L(g(x, \xi_a^j), \xi_b^j) \right]$$

$$= \nabla_x \left[ \frac{1}{n} \sum_{j=1}^{n} L(g(x, \xi_a^j), \xi_b^j) \right]$$

$$= \nabla f(x)$$

Метод стох градиента

---

**Алгоритм 1** Стохастический градиентный спуск (SGD)

**Вход:** размеры шагов $\{\gamma_k\}_{k=0} > 0$, стартовая точка $x^0 \in \mathbb{R}^d$, количество итераций $K$
1: **for** $k = 0, 1, \ldots, K-1$ **do**
2:     Сгенерировать независимо $\xi^k$
3:     Вычислить стохастический градиент $\nabla f(x^k, \xi^k)$
4:     $x^{k+1} = x^k - \gamma_k \nabla f(x^k, \xi^k)$
5: **end for**
**Выход:** $x^K$

---

приближает $\nabla f(x)$

Справка

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

$\mathcal{F}_k$ — $\sigma$ алгебра, задан. сл.в. $X^0, \xi^0, \xi^1 ... \xi^{k-1}$

Суть — "фиксируем" всю случайность до текущей итерации

tower property

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

<u>Док-ва сходимости:</u>

- $f$ — $L$-гладкой, $\mu$ — сильно-выпуклой
- предположение на случайность

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

$$\mathbb{E}_\xi\Big[\|\triangledown f(x, \xi) - \triangledown f(x)\|_2^2\Big] \le \sigma^2$$

Док-во:

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

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

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

$$\color{blue}{\mathbb{E}[\cdot \mid x^k]}$$

$$\color{blue}{\text{это не с.в. значит } \mathbb{E}[\mid x^k]}$$

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

$$- 2\gamma \mathbb{E}\Big[\langle \triangledown f(x^k, \xi^k); x^k - x^* \rangle \mid x^k\Big]$$

$$+ \gamma^2 \mathbb{E}\Big[\|\triangledown f(x^k, \xi^k)\|_2^2 \mid x^k\Big]$$

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

$$= \langle \mathbb{E}\left[\nabla f(x^k, \xi^k) \mid x^k\right]; x^k - x^* \rangle \qquad \text{(по теореме}$$
$$\text{о несмещ.}$$
$$\nabla f(x, \xi)$$
$$\text{и нез. } \xi^k)$$

$$= \langle \nabla f(x^k); x^k - x^* \rangle$$

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

$$-2\gamma \langle \nabla f(x^k); x^k - x^* \rangle \qquad \textcircled{+}$$

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

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

$$\mathbb{E}\left[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|_2^2 + \|\nabla f(x^k)\|_2^2 \right.$$
$$\left. -2\langle \nabla f(x^k); \nabla f(x^k, \xi^k) - \nabla f(x^k)\rangle \mid x^k\right] \;\longrightarrow\; 0$$

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

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

$$\text{(по 2 предположению)}$$

$$\leq \sigma^2 + \|\nabla f(x^k)\|_2^2$$

$$\mathbb{E}\left[\|x^{(k+1)} - x^*\|_2^2 \mid x^k\right] \leq \|x^k - x^*\|_2^2 \quad \textcircled{+}$$

$$-2\gamma \langle \nabla f(x^k); x^k - x^* \rangle \; \textcircled{+}$$

$$+\gamma^2 \|\nabla f(x^k)\|_2^2 + \gamma^2 \sigma^2$$
$$\textcircled{+}$$

$$\leq \|x^k - x^*\|_2^2$$

$$- \mu\gamma \|x^k - x^*\|_2^2 - 2\gamma\left(f(x^k) - f(x^*)\right)$$

$$+ \gamma^2 \cdot 2L\left(f(x^k) - f(x^*)\right) + \gamma^2\sigma^2$$

$$= (1 - \gamma\mu)\|x^k - x^*\|_2^2$$

$$- \gamma(1 - \gamma L)\left(f(x^k) - f(x^*)\right)$$

$$+ \gamma^2\sigma^2$$

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

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

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

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

---

**Теорема сходимость SGD в случае ограниченной дисперсии**

Пусть задача безусловной стохастической оптимизации с $L$-гладкой, $\mu$-сильно выпуклой целевой функцией $f$ решается с помощью SGD с $\gamma_k \leq \frac{1}{L}$ в условиях несмещенности и ограниченности дисперсии стохастического градиента. Тогда справедлива следующая оценка сходимости

$$\mathbb{E}\left[\|x^{k+1} - x^*\|^2\right] \leq (1 - \gamma_k\mu)\mathbb{E}\left[\|x^k - x^*\|^2\right] + \gamma_k^2\sigma^2.$$
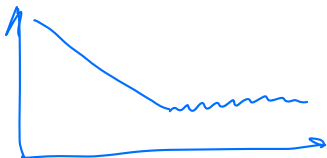
Запустим рекурсию

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

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

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

$$\cdots$$
$$\leq (1-\gamma\mu)^{k+1}\,\mathbb{E}\left[\|x^0 - x^*\|_2^2\right]$$
$$+ \gamma^2\sigma^2 \sum_{i=0}^{k}(1-\gamma\mu)^i$$

$$\textcolor{blue}{\leq \sum_{i=0}^{\infty}(1-\gamma\mu)^i = \frac{1}{\gamma\mu}}$$

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

<span style="color:red">окрестность</span>

<span style="color:green">линейная сходимость к решению</span>

<span style="color:blue">Характерная сходимость:</span>

Борьба за окрестность:

- рестарт метода с меньшим шагом

- $\gamma_k = \dfrac{1}{k+1}$ ; $\dfrac{1}{\sqrt{k+1}}$



$\dfrac{1}{2}$

$\dfrac{1}{\sqrt{k+1}}$

- не иск, иле $\sigma^2$

$$\nabla f(x^k, \xi^k) \rightarrow \frac{1}{b} \sum_{\xi \in S^k} \nabla f(x^k, \xi)$$

$|S^k| = b$

батч, набор из независимых $\xi$

$$\mathbb{E}\left[ \left\| \frac{1}{b} \sum_{\xi \in S^k} \nabla f(x^k, \xi) - \nabla f(x^k) \right\|_2^2 \Big| x^k \right] \leq \frac{\sigma^2}{b}$$

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

$$= \mathbb{E}\left[ \frac{1}{b^2} \sum_{\xi \in S^k} \underbrace{\left\| \nabla f(x^k, \xi) - \nabla f(x^k) \right\|_2^2}_{\leq \sigma^2} \Big| x^k \right]$$

$$+ \mathbb{E}\left[ \frac{1}{b^2} \sum_{\substack{\xi, \eta \in S^k \\ \xi \neq \eta}} \langle \nabla f(x^k, \xi) - \nabla f(x^k); \nabla f(x^k, \eta) - \nabla f(x^k) \rangle \Big| x^k \right]$$

$0$ ($\xi, \eta$ - независимы)

$$\leq \quad \frac{1}{b^2} \sum_{\xi \in S^k} \sigma^2$$

$$\underbrace{\phantom{xxxxx}}$$

$$\frac{\sigma^2 \cdot b}{b^2} = \frac{\sigma^2}{b}$$

Pytorch : HB Polyak => Stoch HB Polyak

$$g^k = \beta g^{k-1} + (1-\beta)\, \nabla S(x^k, \xi^k)$$

$$x^k \approx x^{k-1}$$

есть эффект батширования
(доп. эффект стохаст. HB)

для метода Нестерова:

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

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxx}}_{\text{ускорение}} \qquad \underbrace{\phantom{xx}}$$

сублинейн.
сход.
из-за стох

для град. спуска

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