

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

$$\Downarrow$$

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

Пример

$$(*) \min_{x \in \mathbb{R}^d} \mathbb{E}_{(a,b) \sim D} [L(g(x, a); b)]$$

(здесь ML-суп- D с параметром g)
 (здесь ML-суп- D с параметром g)
 (здесь ML-суп- D с параметром g)
 (здесь ML-суп- D с параметром g)

Проблема: D не given, не можем $f, \nabla f$

Можем сэмплировать из D

- $\xi = (a, b)$ извлекаем ξ по случайной

$$\nabla f(x, \xi) = \nabla_x [L(g(x, a), b)]$$

непрямое:

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

- выборка - операция given: $\{a_i, b_i\}_{i=1}^n \sim D$

$$(**) \min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n L(g(x, a_i), b_i) \leftarrow \text{ERM}$$

(мин. эм. риска)

$(**)$ с вероятн n производится $(*)$
 $\sim \frac{1}{\sqrt{n}}$

(Морне - Купро аналог.)

Далее имея всю выборку уже не сумасно

- дешево и просто

- обладает способностью и стохастичности

$$\nabla f(x, \xi) = \nabla f_{\xi_i}(x) = \nabla_x [L(g(x, a_i), b_i)]$$

(выбираю один из $1 \dots n$ объектов)

$$\mathbb{E}_{\xi} [\nabla f(x, \xi)] = \mathbb{E}_{\xi} [\nabla_x L(g(x, a_{\xi}), b_{\xi})]$$
$$= \sum_{i=1}^n \frac{1}{n} \nabla_x L(g(x, a_i), b_i)$$

ξ генер. с.в.
 $[1 \dots n]$

p_i

$$p = \frac{1}{n}$$

$$= \nabla f(x) \leftarrow \text{траг ERM}$$

Метод со стох. траг.

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** $k = 0, 1, \dots, K-1$ **do**
- 2: Сгенерировать независимо ξ^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

гип. н.о.

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

↑
σ-алгебра, порожд.
 $\mathcal{X}^0, \xi^0, \xi^1, \dots, \xi^{k-1}$

$$E[E[X|Y]] = E[X]$$

Свойства:

• \mathcal{F} — L -матрица, μ -матрица бонграс

• $E_{\mathcal{F}}[\nabla f(x, \xi)] = \nabla f(x)$

$$E_{\mathcal{F}}[\|\nabla f(x, \xi) - \nabla f(x)\|_2^2] \leq \sigma^2$$

$\gamma_k = \gamma$

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

$E[\cdot | \mathcal{X}^k]$ c.l.

$$E[\langle \nabla f(x^k, \xi^k), x^k - x^* \rangle | \mathcal{X}^k] = \langle E[\nabla f(x^k, \xi^k) | \mathcal{X}^k], x^k - x^* \rangle$$

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

$$E[\nabla f(x^k, \xi^k) | x^k] \underset{\text{reg. } \xi^k \text{ on } x^k}{=} E_{\xi^k}[\nabla f(x^k, \xi^k)] = \nabla f(x^k)$$

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

$$E[\|\nabla f(x^k, \xi^k)\|_2^2 | x^k] = (D_\xi = E_{\xi^2} - (E_{\xi})^2)$$

$$= E[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|_2^2 | x^k] + \|\nabla f(x^k)\|_2^2$$

$$\leq 6^2 + \|\nabla f(x^k)\|_2^2$$

↑
of op.
given.

$$E[\|x^{k+1} - x^*\|_2^2 | x^k] \leq \|x^k - x^*\|_2^2 - 2\gamma \langle \nabla f(x^k); x^k - x^* \rangle + \gamma^2 \|\nabla f(x^k)\|_2^2 + \gamma^2 6^2$$

новое

L - Lipschitz, μ - strong convex

$$E[\|x^{k+1} - x^*\|_2^2 | x^k] \leq \|x^k - x^*\|_2^2 - 2\gamma \left(\frac{\mu}{2} \|x^k - x^*\|_2^2 + f(x^k) - f(x^*) \right) + 2\gamma^2 (f(x^k) - f(x^*)) + \gamma^2 6^2$$

$$= (1 - \gamma\mu) \|x^k - x^*\|_2^2 + \gamma^2 6^2 - 2\gamma(1 - \gamma L) (f(x^k) - f(x^*))$$

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

$$\leq (1 - \gamma\mu) \|x^k - x^*\|_2^2 + \gamma^2 6^2$$

$\mathbb{E}[\mathbb{E}[X|Y]]$ от условия к условию

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

Теорема сходимости SGD в случае ограниченной дисперсии

Пусть задача безусловной стохастической оптимизации с L -гладкой, μ -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

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

Заменим шаг $\gamma_k = \gamma$

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

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

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

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

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

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

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

$$= \underbrace{(1 - \gamma\mu)^k R_0}_{\text{слож. шаг}} + \underbrace{\frac{\gamma^2 \sigma^2}{\gamma\mu}}_{\text{огрещенность сходимости}}$$

слож. шаг
сходимости

огрещенность
сходимости

Требуется с орг. шагом:

- $\gamma_k = \frac{1}{k+1}$ $\gamma_k = \frac{1}{\sqrt{k+1}}$

- уменьшаем σ^2

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

где S^k - выб. набор
элементов
размера b

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

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

$$+ \sum_{\xi \neq \eta} \underbrace{\mathbb{E} [\langle \nabla f(x, \xi) - \nabla f(x), \nabla f(x, \eta) - \nabla f(x) \rangle]}_{\substack{\mathbb{E}_\xi \\ \mathbb{E}_\eta}}$$

$$\leq \frac{1}{b^2} (b \cdot \sigma^2 + 0) = \boxed{\frac{\sigma^2}{b}}$$

SGD:

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

можно получить том же

с разными шагами γ_k :

1) $\gamma_k = \frac{1}{L}$ 2) $\gamma_k \sim \frac{1}{k}$

$$R_k \leq (1 - \gamma \mu)^k R_0 + \frac{\sigma^2 \gamma}{\mu b}$$

$$\gamma = \min\left(\frac{1}{L}; \frac{1}{\mu k} \ln(\dots)\right)$$

$$\leq \left(1 - \frac{\mu}{L}\right)^k R_0 + \exp\left(-\frac{\mu k}{\mu k} \ln(\dots)\right) R_0$$

$$+ \frac{\sigma^2}{\mu^2 b k}$$

$$\frac{\sigma^2}{\mu^2 b k R_0}$$

Accer SGD:

(Nesterov)

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

↑
ускорение

↑
ускорение шаг.
комм.

при quasi. γ - сходимости к опр.:

γ шаг. сходим.

$$\gamma \nabla f(x^k)$$

$$x^k \rightarrow x^*$$

$$\nabla f(x^k) \rightarrow 0$$

$$\gamma \nabla f(x^k) \rightarrow 0$$

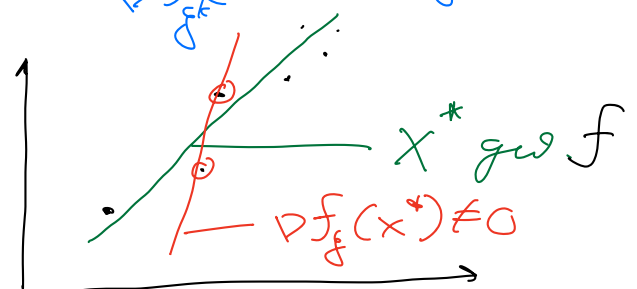
можно считать.

γ SGD:

$$\gamma \nabla f_k(x^k)$$

$$x^k \rightarrow x^*$$

$$\nabla f_k(x^k) \rightarrow \nabla f_k(x^*) \neq 0$$



"Другой" способ градиент: $x^{k+1} = x^k - \gamma g^k$ ← "способ градиент"

но $g^k \rightarrow \nabla f(x^*) = 0$ $x^k \rightarrow x^*$

но возмущения

$$\mathbb{E}[g^k | x^k] = \nabla f(x^k)$$

в онлайн постановке невозможно
а вот офлайн:

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$$

- y_i^k — ^{можно} память: $y_i^0 = 0$
- класть в $y_{i_k}^k = \nabla f_{i_k}(x^k)$
ост. y_i^k не менять
- $\frac{1}{n} \sum_{i=1}^n y_i^k$ — ^{запасов.} градиент $\approx \nabla f(x^k)$

Алгоритм 2 SAGA

Вход: размер шага $\gamma > 0$, стартовая точка $x^0 \in \mathbb{R}^d$, значения памяти $y_i^0 = 0$ для всех $i \in [n]$, количество итераций K

- 1: **for** $k = 0, 1, \dots, K-1$ **do**
- 2: Сгенерировать независимо i_k
- 3: Вычислить $g^k = \nabla f_{i_k}(x^k) - y_{i_k}^k + \frac{1}{n} \sum_{j=1}^n y_j^k$
- 4: Обновить $y_i^{k+1} = \begin{cases} \nabla f_{i_k}(x^k), & \text{если } i = i_k \\ y_i^k, & \text{иначе} \end{cases}$
- 5: $x^{k+1} = x^k - \gamma g^k$
- 6: **end for**

Выход: x^K
