

Стохастическая оптимизация (продолжение).

Координатный спуск

Методы оптимизации

Александр Безносиков

Московский физико-технический институт

7 декабря 2023



В прошлый раз

- Рассматривали постановку вида (оффлайн, ERM):

$$\min_{x \in \mathbb{R}^d} \left[f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right].$$

- Предполагаем, что вызывать полный градиент дорого, но можно, генерируя равномерно и независимо i_k , получить

$$\mathbb{E}_{i_k} [\nabla f_{i_k}(x^k)] = \nabla f(x^k).$$

В прошлый раз

- Рассматривали постановку вида (оффлайн, ERM):

$$\min_{x \in \mathbb{R}^d} \left[f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right].$$

- Предполагаем, что вызывать полный градиент дорого, но можно, генерируя равномерно и независимо i_k , получить

$$\mathbb{E}_{i_k} [\nabla f_{i_k}(x^k)] = \nabla f(x^k).$$

- Идея – взять метод на подобии SGD:

$$x^{k+1} = x^k - \gamma g^k,$$

где

$$g^k \rightarrow \nabla f(x^*) = 0, \quad \text{при} \quad x^k \rightarrow x^*.$$

Уже знакомы

Алгоритм 1 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 $y_i^{k+1} = \begin{cases} \nabla f_i(x^k), & \text{если } i = i_k \\ y_i^k, & \text{иначе} \end{cases}$

5: $x^{k+1} = x^k - \gamma g^k$

6: **end for**

Выход: x^K

- Идея – если я считал когда-то градиент для f_i , то зачем его забывать? Сохраним!

- Идея – если я считал когда-то градиент для f_i , то зачем его забывать? Сохраним!
- $\frac{1}{n} \sum_{j=1}^n y_j^k$ – «запаздывающая» версия $\nabla f(x^k)$.

- Идея – если я считал когда-то градиент для f_i , то зачем его забывать? Сохраним!
- $\frac{1}{n} \sum_{j=1}^n y_j^k$ – «запаздывающая» версия $\nabla f(x^k)$.
- $\mathbb{E} [g^k \mid x^k] = \nabla f(x^k)$.

- Идея – если я считал когда-то градиент для f_i , то зачем его забывать? Сохраним!
- $\frac{1}{n} \sum_{j=1}^n y_j^k$ – «запаздывающая» версия $\nabla f(x^k)$.
- $\mathbb{E}[g^k \mid x^k] = \nabla f(x^k)$.
- При $x^k \rightarrow x^*$ имеем, что $y_j^k \rightarrow \nabla f_j(x^*)$, и $\frac{1}{n} \sum_{j=1}^n y_j^k \rightarrow \nabla f(x^*) = 0$.
А значит $g^k \rightarrow 0$.

- Идея – если я считал когда-то градиент для f_i , то зачем его забывать? Сохраним!
- $\frac{1}{n} \sum_{j=1}^n y_j^k$ – «запаздывающая» версия $\nabla f(x^k)$.
- $\mathbb{E} [g^k \mid x^k] = \nabla f(x^k)$.
- При $x^k \rightarrow x^*$ имеем, что $y_j^k \rightarrow \nabla f_j(x^*)$, и $\frac{1}{n} \sum_{j=1}^n y_j^k \rightarrow \nabla f(x^*) = 0$.
А значит $g^k \rightarrow 0$.
- Из минусов: лишняя $\mathcal{O}(nd)$ память.

SAGA: доказательство

- Все f_i являются L -гладкими и выпуклыми, а $f - \mu$ - сильно выпуклой.

$$\|x^{k+1} - x^*\|_2^2 = \|x^k - x^*\|_2^2 - 2\gamma \langle g^k; x^k - x^* \rangle + \gamma^2 \|g^k\|_2^2$$

$$\mathbb{E}[\cdot | x^k]$$

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

$$1) \mathbb{E}[g^k | x^k] = \nabla f(x^k)$$

$$\left\{ \begin{aligned} \text{used: } \tilde{g}^k &= \frac{1}{n} \sum_{j=1}^n g_j^{k+1} \\ \tilde{g}^k &= \frac{1}{n} \sum_{j=1}^n g_j^k - \cancel{\frac{1}{n} (g_{i_k}^k + \nabla f_{i_k}(x^k))} \end{aligned} \right. \quad \text{from SAGA}$$

$$\mathbb{E}[\tilde{g}^k | x^k] =$$

$$= \mathbb{E}[\cancel{\frac{1}{n} (\nabla f_{i_k}(x^k) - g_{i_k}^k)} | x^k] + \frac{1}{n} \sum g^k$$

$$= \cancel{\frac{1}{n}} \sum_{j=1}^n \frac{1}{n} (\nabla f_j(x^k) - g_j^k) + \frac{1}{n} \sum g^k$$

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

$$2) \mathbb{E}[\|g^k - \nabla f(x^*)\|_2^2 | x^k] =$$

$$= \mathbb{E}[\| \underbrace{\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*)}_{a} + \underbrace{\nabla f_{i_k}(x^*) - g_{i_k}^k}_{b} + \frac{1}{n} \sum_{j=1}^n g_j^k - \nabla f(x^*) \|_2^2 | x^k]$$

$$\stackrel{\text{KBLU}}{\leq} 2 \mathbb{E}[\| \nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*) \|_2^2 | x^k] + 2 \mathbb{E}[\| \nabla f_{i_k}(x^*) - g_{i_k}^k + \frac{1}{n} \sum_{j=1}^n (g_j^k - \nabla f_j(x^*)) \|_2^2 | x^k]$$

$$\leq 2 \mathbb{E}[\| \nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*) \|_2^2 | x^k] + 2 \mathbb{E}[\| \nabla f_{i_k}(x^*) - g_{i_k}^k \|_2^2 | x^k]$$

$$= 2 \cdot \frac{1}{n} \sum_{i=1}^n \| \nabla f_i(x^k) - \nabla f_i(x^*) \|_2^2$$

$$+ 2 \cdot \frac{1}{n} \sum_{i=1}^n \| g_i^k - \nabla f_i(x^*) \|_2^2$$

$$\uparrow \quad \uparrow$$

$$\mathbb{E}[\| \xi - \mathbb{E}[\xi | x^k] \|_2^2 | x^k]$$

$$\leq \mathbb{E}[\| \xi \|_2^2 | x^k]$$

SAGA: доказательство

- Все f_i являются L -гладкими и выпуклыми, а $f - \mu$ - сильно выпуклой.
- Уже привычно:

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

SAGA: доказательство

- Все f_i являются L -гладкими и выпуклыми, а $f - \mu$ - сильно выпуклой.
- Уже привычно:

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

- Берем условное мат.ожидание по случайности только на итерации k :

$$\begin{aligned} \mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] &= \|x^k - x^*\|_2^2 - 2\gamma \langle \mathbb{E} [g^k \mid x^k], x^k - x^* \rangle \\ &\quad + \gamma^2 \mathbb{E} \left[\|g^k - \nabla f(x^*)\|_2^2 \mid x^k \right]. \end{aligned}$$

SAGA: доказательство

- Работаем с $\mathbb{E} [g^k \mid x^k]$:

$$\begin{aligned}\mathbb{E} [g^k \mid x^k] &= \mathbb{E} \left[\nabla f_{i_k}(x^k) - y_{i_k}^k + \frac{1}{n} \sum_{j=1}^n y_j^k \mid x^k \right] \\ &= \mathbb{E} \left[\nabla f_{i_k}(x^k) - y_{i_k}^k \mid x^k \right] + \frac{1}{n} \sum_{j=1}^n y_j^k \\ &= \frac{1}{n} \sum_{j=1}^n \left[\nabla f_j(x^k) - y_j^k \right] + \frac{1}{n} \sum_{j=1}^n y_j^k \\ &= \nabla f(x^k)\end{aligned}$$

SAGA: доказательство

- Теперь работаем с $\mathbb{E} [\|g^k - \nabla f(x^*)\|_2^2 \mid x^k]$:

$$\begin{aligned}\mathbb{E} [\|g^k - \nabla f(x^*)\|_2^2 \mid x^k] &= \mathbb{E} \left[\left\| \nabla f_{i_k}(x^k) - y_{i_k}^k + \frac{1}{n} \sum_{j=1}^n y_j^k - \nabla f(x^*) \right\|_2^2 \mid x^k \right] \\ &= \mathbb{E} \left[\left\| \nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*) + \nabla f_{i_k}(x^*) - y_{i_k}^k \right. \right. \\ &\quad \left. \left. + \frac{1}{n} \sum_{j=1}^n y_j^k - \nabla f(x^*) \right\|_2^2 \mid x^k \right] \\ &\leq 2\mathbb{E} \left[\left\| \nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*) \right\|_2^2 \mid x^k \right] \\ &\quad + 2\mathbb{E} \left[\left\| \nabla f_{i_k}(x^*) - y_{i_k}^k + \frac{1}{n} \sum_{j=1}^n y_j^k - \nabla f(x^*) \right\|_2^2 \mid x^k \right].\end{aligned}$$

SAGA: доказательство

- Пользуемся тем, что $\mathbb{D}\xi \leq \mathbb{E}[\xi^2]$:

$$\begin{aligned}\mathbb{E} \left[\|g^k - \nabla f(x^*)\|_2^2 \mid x^k \right] &\leq 2\mathbb{E} \left[\|\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*)\|_2^2 \mid x^k \right] \\ &\quad + 2\mathbb{E} \left[\left\| \nabla f_{i_k}(x^*) - y_{i_k}^k + \frac{1}{n} \sum_{j=1}^n y_j^k - \nabla f(x^*) \right\|_2^2 \mid x^k \right] \\ &\leq 2\mathbb{E} \left[\|\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*)\|_2^2 \mid x^k \right] \\ &\quad + 2\mathbb{E} \left[\|\nabla f_{i_k}(x^*) - y_{i_k}^k\|_2^2 \mid x^k \right]\end{aligned}$$

SAGA: доказательство

- Берем мат.ожидание, пользуемся гладкостью (с выпуклостью):

$$\begin{aligned}\mathbb{E} \left[\|g^k - \nabla f(x^*)\|_2^2 \mid x^k \right] &\leq 2\mathbb{E} \left[\|\nabla f_{i_k}(x^k) - \nabla f_{i_k}(x^*)\|_2^2 \mid x^k \right] \\ &\quad + 2\mathbb{E} \left[\|\nabla f_{i_k}(x^*) - y_{i_k}^k\|_2^2 \mid x^k \right] \\ &\leq 4L \cdot \frac{1}{n} \sum_{i=1}^n \left(f_i(x^k) - f_i(x^*) - \langle \nabla f_i(x^*), x^k - x^* \rangle \right) \\ &\quad + 2 \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2 \\ &= 4L \cdot (f(x^k) - f(x^*)) \\ &\quad + 2 \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2\end{aligned}$$

SAGA: доказательство

- Промежуточный итог:

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

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

$\| \nabla f(x^k) \|_2^2 + \| g^k - \nabla f(x^k) \|^2$

$$\mathbb{E} \left[\|g^k - \nabla f(x^*)\|_2^2 \mid x^k \right] \leq 4L \cdot (f(x^k) - f(x^*)) + 2 \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2$$

$y_i^k \rightarrow \nabla f_i(x^*)$

SAGA: доказательство

- Промежуточный итог:

$$\begin{aligned}\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] &= \|x^k - x^*\|_2^2 - 2\gamma \langle \mathbb{E} [g^k \mid x^k], x^k - x^* \rangle \\ &\quad + \gamma^2 \mathbb{E} \left[\|g^k - \nabla f(x^*)\|_2^2 \mid x^k \right].\end{aligned}$$

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

$$\mathbb{E} \left[\|g^k - \nabla f(x^*)\|_2^2 \mid x^k \right] \leq 4L \cdot (f(x^k) - f(x^*)) + 2 \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2$$

- Собираем вместе:

$$\begin{aligned}\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] &\leq \|x^k - x^*\|_2^2 - 2\gamma \langle \nabla f(x^k), x^k - x^* \rangle \\ &\quad + \gamma^2 \left(4L \cdot (f(x^k) - f(x^*)) + 2 \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2 \right)\end{aligned}$$

- Сильная выпуклость функции f :

$$\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] \leq (1 - \mu\gamma) \|x^k - x^*\|_2^2 - \underbrace{2\gamma(1 - 2\gamma L)(f(x^k) - f(x^*))}_{\text{negative term}} + \underbrace{2\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2}_{\text{variance term}}.$$

- Сильная выпуклость функции f :

$$\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] \leq (1 - \mu\gamma) \|x^k - x^*\|_2^2 - 2\gamma(1 - 2\gamma L)(f(x^k) - f(x^*)) \\ + 2\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2.$$

- Более формально пришли к тому, что если $y_i^k \rightarrow \nabla f_i(x^*)$, то дисперсия «убивается», а значит будет линейная сходимость. Покажем, как это можно строго оформить.

$$\mathbb{E} \left[\frac{1}{n} \sum \underbrace{\|y_i^{k+1} - \nabla f_i(x^*)\|_2^2}_{G_{k+1}} \mid x^k \right] =$$

$$= \frac{1}{n} \sum \mathbb{E} \left[\|y_i^{k+1} - \nabla f_i(x^*)\|_2^2 \mid x^k \right] =$$

$$= \frac{1}{n} \sum \left[\left(1 - \frac{1}{n}\right) \|y_i^k - \nabla f_i(x^*)\|_2^2 + \frac{1}{n} \|\nabla f_i(x^k) - \nabla f_i(x^*)\|_2^2 \right]$$

$$= \left(1 - \frac{1}{n}\right) \cdot \frac{1}{n} \sum \|y_i^k - \nabla f_i(x^*)\|_2^2 + \frac{1}{n} \cdot 2L (f(x^k) - f(x^*))$$

$$\mathbb{E}[G_{k+1} \mid x^k] = \left(1 - \frac{1}{n}\right) G_k + \frac{2L}{n} (f(x^k) - f(x^*))$$

SAGA: доказательство

- Рассмотрим, как ведет себя $\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2$:

$$\begin{aligned}\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|y_i^{k+1} - \nabla f_i(x^*)\|_2^2 \mid x^k \right] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\|y_i^{k+1} - \nabla f_i(x^*)\|_2^2 \mid x^k \right] \\ &= \left(1 - \frac{1}{n}\right) \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \\ &\quad + \frac{1}{n} \cdot \frac{1}{n} \sum_{i=1}^n \|f_i(x^k) - \nabla f_i(x^*)\|_2^2 \\ &\leq \left(1 - \frac{1}{n}\right) \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \\ &\quad + \frac{1}{n} \cdot 2L(f(x^k) - f(x^*)).\end{aligned}$$

SAGA: доказательство

- Итого (здесь сразу накинуто полное математическое ожидание):

$$\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \right] \leq (1 - \mu\gamma) \mathbb{E} \left[\|x^k - x^*\|_2^2 \right] - 2\gamma(1 - 2\gamma L) \mathbb{E} \left[f(x^k) - f(x^*) \right] + 2\gamma^2 \cdot \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2 \right]$$

mm. avg. $\delta \leq \gamma n$

$$\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|y_i^{k+1} - \nabla f_i(x^*)\|_2^2 \right] \leq \left(1 - \frac{1}{n} \right) \cdot \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \right] + \frac{1}{n} \cdot 2L \mathbb{E} \left[f(x^k) - f(x^*) \right]$$

mm avg

SAGA: доказательство

- Итого (здесь сразу накинута полное математическое ожидание):

$$\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \right] \leq (1 - \mu\gamma) \mathbb{E} \left[\|x^k - x^*\|_2^2 \right] - 2\gamma(1 - 2\gamma L) \mathbb{E} \left[f(x^k) - f(x^*) \right] \\ + 2\gamma^2 \cdot \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*) - y_i^k\|_2^2 \right]$$

$$\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|y_i^{k+1} - \nabla f_i(x^*)\|_2^2 \right] \leq \left(1 - \frac{1}{n} \right) \cdot \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \right] \\ + \frac{1}{n} \cdot 2L \mathbb{E} \left[f(x^k) - f(x^*) \right].$$

- Получилось две «сходящиеся» последовательности, осталось их аккуратно «сшить».

SAGA: доказательство

- Пусть $M > 0$:

$$\begin{aligned} & \mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 + \underbrace{M\gamma^2}_{\text{red circle}} \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^{k+1} - \nabla f_i(x^*)\|_2^2 \right] \quad \left(1 + \frac{2}{n} - \frac{1}{n}\right) < 1 \\ & \leq (1 - \mu\gamma) \mathbb{E} \left[\|x^k - x^*\|_2^2 \right] \quad \text{blue circle around } (1 - \mu\gamma) \\ & \quad + \left(1 + \frac{2}{\underbrace{M}_{\text{blue circle}}} - \frac{1}{n}\right) \mathbb{E} \left[\underbrace{M\gamma^2}_{\text{red underline}} \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \right] \\ & \quad - 2\gamma \left(1 - 2\gamma L - \frac{\gamma ML}{n}\right) \mathbb{E} \left[f(x^k) - f(x^*) \right] \end{aligned}$$

$M = 4n$

SAGA: доказательство

- Возьмем $M = 4n$:

$$\begin{aligned} & \mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 + 4n\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^{k+1} - \nabla f_i(x^*)\|_2^2 \right] \\ & \leq (1 - \mu\gamma) \mathbb{E} \left[\|x^k - x^*\|_2^2 \right] \\ & \quad + \left(1 - \frac{1}{2n} \right) \mathbb{E} \left[4n\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \right] \\ & \quad - \cancel{2\gamma(1 - 6\gamma L) \mathbb{E} [f(x^k) - f(x^*)]} \end{aligned}$$

$$\gamma \leq \frac{1}{6L}$$

SAGA: доказательство

- Возьмем $M = 4n$:

$$\begin{aligned} & \mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 + 4n\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^{k+1} - \nabla f_i(x^*)\|_2^2 \right] \\ & \leq \underbrace{(1 - \mu\gamma)}_{\text{blue underline}} \mathbb{E} \left[\|x^k - x^*\|_2^2 \right] \\ & \quad + \underbrace{\left(1 - \frac{1}{2n}\right)}_{\text{blue underline}} \mathbb{E} \left[4n\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \right] \\ & \quad - 2\gamma(1 - 6\gamma L) \mathbb{E} \left[f(x^k) - f(x^*) \right] \end{aligned}$$

- Теперь $\gamma \leq \frac{1}{6L}$:

$$\begin{aligned} & \mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 + 4n\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^{k+1} - \nabla f_i(x^*)\|_2^2 \right] \\ & \leq \max \left\{ (1 - \mu\gamma); \left(1 - \frac{1}{2n}\right) \right\} \mathbb{E} \left[\|x^k - x^*\|_2^2 + 4n\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2 \right] \end{aligned}$$

SAGA: сходимость

- Получили сходимость, но по необычному критерию. Суть критерия в отражении физики, как сходимости $x^k \rightarrow x^*$, так и $y_i^k \rightarrow \nabla f_i(x^*)$, что и закладывали в метод.

Теорема сходимость SAGA

Пусть задача безусловной стохастической оптимизации вида конечной суммы с L -гладкими, выпуклыми функциями f_i и μ -сильно выпуклой целевой функцией f решается с помощью SAGA с $\gamma \leq \frac{1}{6L}$. Тогда справедлива следующая оценка сходимости

$$\mathbb{E}[\|x^k - x^*\|_2^2] \leq \mathbb{E}[V_k] \leq \max \left\{ (1 - \mu\gamma); \left(1 - \frac{1}{2n}\right) \right\}^k \mathbb{E}[V_0],$$

$$\text{где } V_k = \|x^k - x^*\|_2^2 + 4n\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2.$$

$$\|x^k - x^*\|_2^2 \leq V_k \rightarrow 0$$

SAGA: сходимость

- Получили сходимость, но по необычному критерию. Суть критерия в отражении физики, как сходимости $x^k \rightarrow x^*$, так и $y_i^k \rightarrow \nabla f_i(x^*)$, что и закладывали в метод.

Теорема сходимость SAGA

Пусть задача безусловной стохастической оптимизации вида конечной суммы с L -гладкими, выпуклыми функциями f_i и μ -сильно выпуклой целевой функцией f решается с помощью SAGA с $\gamma \leq \frac{1}{6L}$. Тогда справедлива следующая оценка сходимости

$$\mathbb{E}[V_k] \leq \max \left\{ (1 - \mu\gamma); \left(1 - \frac{1}{2n}\right) \right\}^k \mathbb{E}[V_0],$$

где $V_k = \|x^k - x^*\|_2^2 + 4n\gamma^2 \cdot \frac{1}{n} \sum_{i=1}^n \|y_i^k - \nabla f_i(x^*)\|_2^2$.

- Легко заметить, что из сходимости по $\mathbb{E}[V_k]$ следует и сходимость по $\mathbb{E}[\|x^k - x^*\|_2^2]$: $\mathbb{E}[\|x^k - x^*\|_2^2] \leq \mathbb{E}[V_k]$

SAGA: сходимость

- **Вопрос:** почему не взять M огромным, тогда сходимость будет лучше?

SAGA: сходимость

- **Вопрос:** почему не взять M огромным, тогда сходимость будет лучше? M еще влияет на критерий сходимости, который так же будет расти с ростом M . При этом сходимости лучше, чем $(1 - \frac{1}{n})$ не добиться.

SAGA: сходимость

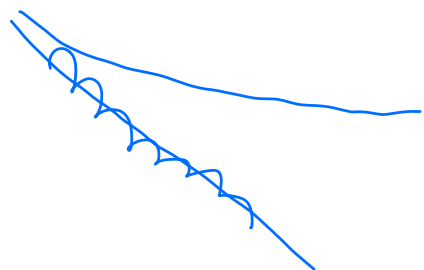
- **Вопрос:** почему не взять M огромным, тогда сходимость будет лучше? M еще влияет на критерий сходимости, который так же будет расти с ростом M . При этом сходимости лучше, чем $(1 - \frac{1}{n})$ не добиться.
- Получаем следующую (уже анонсированную в прошлый раз) оценку на число итераций для достижения точности ε :

$$\frac{\left(1 - \frac{\mu}{6L}\right)}{\left(1 - \frac{1}{2n}\right)}$$

$$\begin{aligned} & \mathcal{O}\left(\left[n + \frac{L}{\mu}\right] \log \frac{1}{\varepsilon}\right) \\ & \max\left(2n; \frac{6L}{\mu}\right) \log \frac{1}{\varepsilon} \\ & \leq \left(2n + \frac{6L}{\mu}\right) \log \frac{1}{\varepsilon} \\ & \mathcal{O}\left(n + \frac{L}{\mu}\right) \log \frac{1}{\varepsilon} \end{aligned}$$

SAGA: сходимость

- **Вопрос:** почему не взять M огромным, тогда сходимость будет лучше? M еще влияет на критерий сходимости, который так же будет расти с ростом M . При этом сходимости лучше, чем $(1 - \frac{1}{n})$ не добиться.
- Получаем следующую (уже анонсированную в прошлый раз) оценку на число итераций для достижения точности ε :



$$\mathcal{O} \left(\left[n + \frac{L}{\mu} \right] \log \frac{1}{\varepsilon} \right)$$

$$f(x) - f(x^*) + \|x^k - x^*\|_2^2$$

- У классического градиентного спуска оценка:

$$\mathcal{O} \left(n \frac{L}{\mu} \log \frac{1}{\varepsilon} \right), \quad n \left(1 + \frac{L}{\mu} \right)$$

но оракульная сложность (подсчет градиентов ∇f_i) у градиентного спуска в n раз больше.

Поворот на 90 градусов

- Простой пример:

$$\min_{x \in \mathbb{R}^d} \left[f(x) = \frac{1}{n} \sum_{i=1}^n \|a_i^T x - b_i\|_2^2 \right],$$

где $\{a_i, b_i\}_{i=1}^n$ – обучающая выборка.

Handwritten notes illustrating the matrix notation for the optimization problem:

$$A = \begin{pmatrix} a_1 \\ \vdots \\ a_i \\ \vdots \\ a_n \end{pmatrix}$$
$$\nabla g(x) = [Bx - C]_{i_k}$$

Diagram illustrating the dimensions of the matrices and vectors:

- B is a matrix of size $n \times d$ (represented by a rectangle with n rows and d columns).
- x is a vector of size d (represented by a vertical rectangle with d elements).
- C is a vector of size n (represented by a horizontal rectangle with n elements).
- The result $[Bx - C]_{i_k}$ is a scalar value (represented by a circle with i_k inside).

The dimensions d^2 and d are also indicated next to the matrix B and vector x respectively.

Поворот на 90 градусов

- Простой пример:

$$\min_{x \in \mathbb{R}^d} \left[f(x) = \frac{1}{n} \sum_{i=1}^n \|a_i^T x - b_i\|_2^2 \right],$$

где $\{a_i, b_i\}_{i=1}^n$ – обучающая выборка.

- До этого, мы брали не всю выборку для подсчета градиента, чтобы быть более эффективными. Т.е. использовали только часть строк матрицы A , составленной из a_i **Вопрос:** а как по-другому можно добиться эффективности?

Поворот на 90 градусов

- Простой пример:

$$\min_{x \in \mathbb{R}^d} \left[f(x) = \frac{1}{n} \sum_{i=1}^n \|a_i^T x - b_i\|_2^2 \right],$$

где $\{a_i, b_i\}_{i=1}^n$ – обучающая выборка.

- До этого, мы брали не всю выборку для подсчета градиента, чтобы быть более эффективными. Т.е. использовали только часть строк матрицы A , составленной из a_i **Вопрос:** а как по-другому можно добиться эффективности? если до этого был выбор строк матрицы A , то теперь можно попробовать как-то завязаться на столбцы. **Вопрос:** а что означает «выбор столбцов»?

Поворот на 90 градусов

- Простой пример:

$$\min_{x \in \mathbb{R}^d} \left[f(x) = \frac{1}{n} \sum_{i=1}^n \|a_i^T x - b_i\|_2^2 \right],$$

где $\{a_i, b_i\}_{i=1}^n$ – обучающая выборка.

- До этого, мы брали не всю выборку для подсчета градиента, чтобы быть более эффективными. Т.е. использовали только часть строк матрицы A , составленной из a_i **Вопрос:** а как по-другому можно добиться эффективности? если до этого был выбор строк матрицы A , то теперь можно попробовать как-то завязаться на столбцы. **Вопрос:** а что означает «выбор столбцов»? Выбор координат вектора x .

Производная по направлению

- Часто для более сложных задач к подсчету производных по координатам/направлениям прибегают не из-за удешевления процесса, а из-за доступности только оракула нулевого порядка (значения функции). Потому что производную по направлению $e \in \{e \in \mathbb{R}^d \mid \|e\|_2 \leq 1\}$ можно аппроксимировать через конечную разность:

$$[\nabla f(x)]_e \approx \frac{f(x + \tau e) - f(x - \tau e)}{2\tau} \quad \tau \rightarrow 0$$

(таким образом можно «собрать» и весь «градиент»).

Координатный метод

Алгоритм 2 Координатный метод

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 из $[d]$
- 3: Вычислить $[\nabla f(x^k)]_{i_k}$
- 4: $x^{k+1} = x^k - \gamma \cdot d[\nabla f(x^k)]_{i_k} e_{i_k}$
- 5: **end for**

Выход: x^K

Здесь e_{i_k} – i_k -ый базисный вектор

Координатный метод

Алгоритм 3 Координатный метод

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 из $[d]$
- 3: Вычислить $[\nabla f(x^k)]_{i_k}$
- 4: $x^{k+1} = x^k - \gamma \cdot d[\nabla f(x^k)]_{i_k} e_{i_k}$
- 5: **end for**

Выход: x^K

Здесь e_{i_k} – i -ый базисный вектор

- Зачем в шаге метода есть домножение на d ?

Координатный метод

Алгоритм 4 Координатный метод

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 из $[d]$
- 3: Вычислить $[\nabla f(x^k)]_{i_k}$
- 4: $x^{k+1} = x^k - \gamma \cdot d[\nabla f(x^k)]_{i_k} e_{i_k}$
- 5: **end for**

Выход: x^K

Здесь e_{i_k} — i -ый базисный вектор

- Зачем в шаге метода есть домножение на d ? Для несмещенности того, что мы используем вместо градиента.

Координатный метод: доказательство

- f является L -гладкой и μ - сильно выпуклой.

Координатный метод: доказательство

- f является L -гладкой и μ - сильно выпуклой.
- Уже привычно:

$$\|x^{k+1} - x^*\|_2^2 = \|x^k - x^*\|_2^2 - 2\gamma \langle d[\nabla f(x^k)]_{i_k} e_{i_k}, x^k - x^* \rangle + \gamma^2 \|d[\nabla f(x^k)]_{i_k} e_{i_k}\|_2^2.$$

$$\mathbb{E}[d[\nabla f(x^k)]_{i_k} e_{i_k} | x_k] = \underbrace{d}_{\text{скаляр}} \cdot \sum_{j=1}^d \frac{1}{d} [\nabla f(x^k)]_j e_j =$$

$$\begin{aligned} \mathbb{E}[\|d[\nabla f(x^k)]_{i_k} e_{i_k}\|_2^2 | x^k] &= d^2 \cdot \sum_{j=1}^d \frac{1}{d} \|\underbrace{e_j}_{\text{вектор}}\|_2^2 \|\nabla f(x^k)\|_2^2 \\ &= d \sum_{j=1}^d [\nabla f(x^k)]_j^2 = d \|\nabla f(x^k)\|_2^2 \end{aligned}$$

Координатный метод: доказательство

- f является L -гладкой и μ - сильно выпуклой.
- Уже привычно:

$$\begin{aligned}\|x^{k+1} - x^*\|_2^2 &= \|x^k - x^*\|_2^2 - 2\gamma \langle d[\nabla f(x^k)]_{i_k} e_{i_k}, x^k - x^* \rangle \\ &\quad + \gamma^2 \|d[\nabla f(x^k)]_{i_k} e_{i_k}\|_2^2.\end{aligned}$$

- Берем условное мат.ожидание по случайности только на итерации k :

$$\begin{aligned}\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] &= \|x^k - x^*\|_2^2 - 2\gamma \langle \mathbb{E} \left[d[\nabla f(x^k)]_{i_k} e_{i_k} \mid x^k \right], x^k - x^* \rangle \\ &\quad + \gamma^2 \mathbb{E} \left[\|d[\nabla f(x^k)]_{i_k} e_{i_k}\|_2^2 \mid x^k \right].\end{aligned}$$

Координатный метод: доказательство

- Работаем с $\mathbb{E} [d[\nabla f(x^k)]_{i_k} e_{i_k} \mid x^k]$:

$$\begin{aligned}\mathbb{E} [d[\nabla f(x^k)]_{i_k} e_{i_k} \mid x^k] &= \frac{1}{d} \sum_{j=1}^d d[\nabla f(x^k)]_j e_j \\ &= \nabla f(x^k)\end{aligned}$$

Координатный метод: доказательство

- Теперь работаем с $\mathbb{E} [\|d[\nabla f(x^k)]\|_{i_k}^2 \mid x^k]$:

$$\begin{aligned}\mathbb{E} [\|d[\nabla f(x^k)]\|_{i_k}^2 \mid x^k] &= \mathbb{E} [\|d[\nabla f(x^k)]\|_{i_k}^2 \mid x^k] \\ &= d^2 \mathbb{E} [\|[\nabla f(x^k)]_{i_k}\|_2^2 \mid x^k] \\ &= d^2 \cdot \frac{1}{d} \sum_{j=1}^d \|[\nabla f(x^k)]_j e_j\|_2^2 \\ &= d \|\nabla f(x^k)\|_2^2\end{aligned}$$

Координатный метод: доказательство

- Промежуточный итог:

$$\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] = \|x^k - x^*\|_2^2 - 2\gamma \langle \mathbb{E} \left[d[\nabla f(x^k)]_{i_k} e_{i_k} \mid x^k \right], x^k - x^* \rangle + \gamma^2 \mathbb{E} \left[\|d[\nabla f(x^k)]_{i_k} e_{i_k}\|_2^2 \mid x^k \right].$$

$$\mathbb{E} \left[d[\nabla f(x^k)]_{i_k} e_{i_k} \mid x^k \right] = \underline{\nabla f(x^k)}$$

$$\mathbb{E} \left[\|d[\nabla f(x^k)]_{i_k} e_{i_k}\|_2^2 \mid x^k \right] = \underline{d \|\nabla f(x^k)\|_2^2}$$

Координатный метод: доказательство

- Промежуточный итог:

$$\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] = \|x^k - x^*\|_2^2 - 2\gamma \langle \mathbb{E} [d[\nabla f(x^k)]_{i_k} e_{i_k} \mid x^k], x^k - x^* \rangle + \gamma^2 \mathbb{E} \left[\|d[\nabla f(x^k)]_{i_k} e_{i_k}\|_2^2 \mid x^k \right].$$

$$\mathbb{E} \left[d[\nabla f(x^k)]_{i_k} e_{i_k} \mid x^k \right] = \nabla f(x^k)$$

$$\mathbb{E} \left[\|d[\nabla f(x^k)]_{i_k} e_{i_k}\|_2^2 \mid x^k \right] = d \|\nabla f(x^k)\|_2^2$$

- Собираем вместе:

$$\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] \leq \underbrace{\|x^k - x^*\|_2^2} + \underbrace{- 2\gamma \langle \nabla f(x^k), x^k - x^* \rangle + d\gamma^2 \|\nabla f(x^k)\|_2^2}_{\text{blue underlines}}.$$

Координатный метод: доказательство

- Сильная выпуклость и гладкость функции f :

$$\mathbb{E} \left[\|x^{k+1} - x^*\|_2^2 \mid x^k \right] \leq \underbrace{(1 - \mu\gamma)}_{\text{blue}} \underbrace{\|x^k - x^*\|_2^2}_{\text{blue}} - \underbrace{2\gamma(1 - d\gamma L)}_{\text{blue}} (f(x^k) - f(x^*)).$$

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

$$\frac{2}{L}$$

$$\min \left(\begin{matrix} x_1^2 + x_2^2 \\ + x_d^2 \end{matrix} \right)$$

$$\frac{(1, 0, \dots, 0)}{d \begin{pmatrix} 2x_1 \\ 0 \end{pmatrix}}$$

Координатный метод: доказательство

- Сильная выпуклость и гладкость функции f :

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

- Пусть $\gamma \leq \frac{1}{dL}$:

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

Координатный метод: сходимость

$$\frac{1}{d}$$

$$x_1^2 + x_2^2 + \dots + x_d^2$$

$$\begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Теорема сходимость SAGA

Пусть задача безусловной оптимизации с L -гладкой и μ -сильно выпуклой целевой функцией f решается с помощью координатного метода с $\gamma \leq \frac{1}{dL}$. Тогда справедлива следующая оценка сходимости

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

$$\frac{dL}{\mu} \log \frac{1}{\varepsilon}$$

$$\left(1 - \frac{\mu}{dL}\right)^k$$

Координатный метод: сходимость

- Подставив $\gamma = \frac{1}{dL}$, получаем следующую итерационную сложность

$$\mathcal{O} \left(\frac{dL}{\mu} \log \frac{1}{\varepsilon} \right).$$

Координатный метод: сходимость

- Подставив $\gamma = \frac{1}{dL}$, получаем следующую итерационную сложность

$$\mathcal{O} \left(\frac{dL}{\mu} \log \frac{1}{\varepsilon} \right).$$

Вопрос: есть ли улучшения по сравнению с обычным градиентным спуском?

Координатный метод: сходимость

- Подставив $\gamma = \frac{1}{dL}$, получаем следующую итерационную сложность

$$\mathcal{O} \left(\frac{dL}{\mu} \log \frac{1}{\varepsilon} \right).$$

Вопрос: есть ли улучшения по сравнению с обычным градиентным спуском? В общем случае нет. Это доказуемо так.

Координатный метод: сходимость

- Подставив $\gamma = \frac{1}{dL}$, получаем следующую итерационную сложность

$$\mathcal{O} \left(\frac{dL}{\mu} \log \frac{1}{\varepsilon} \right).$$

Вопрос: есть ли улучшения по сравнению с обычным градиентным спуском? В общем случае нет. Это доказуемо так.

- Если есть дополнительная информация о задаче (например, свойства констант Липшица градиента по направлению), то улучшения можно получить.

Координатный метод: сходимость

- Подставив $\gamma = \frac{1}{dL}$, получаем следующую итерационную сложность

$$\mathcal{O} \left(\frac{dL}{\mu} \log \frac{1}{\varepsilon} \right).$$

Вопрос: есть ли улучшения по сравнению с обычным градиентным спуском? В общем случае нет. Это доказуемо так.

- Если есть дополнительная информация о задаче (например, свойства констант Липшица градиента по направлению), то улучшения можно получить.
- Еще координатный метод часто хорошо себя проявляет на практике.

Координатный метод: сходимость

- Подставив $\gamma = \frac{1}{dL}$, получаем следующую итерационную сложность

$$\mathcal{O} \left(\frac{dL}{\mu} \log \frac{1}{\varepsilon} \right).$$

$$d \int \frac{L}{\mu}$$

Вопрос: есть ли улучшения по сравнению с обычным градиентным спуском? В общем случае нет. Это доказуемо так.

- Если есть дополнительная информация о задаче (например, свойства констант Липшица градиента по направлению), то улучшения можно получить.
- Еще координатный метод часто хорошо себя проявляет на практике.
- Результат обобщается и на случай выбора нескольких координат.
- Возможно ускорение с помощью двух моментумов.

SEGA

$$y \in \mathbb{R}^d$$

$$y^0 = 0$$

i_k - random index.

$$[\nabla f(x^k)]_{i_k}$$

$$y^{k+1} = y^k - y_{i_k}^k + [\nabla f(x^k)]_{i_k}$$

$$g^{k+1} = y^k - d(y_{i_k}^k + [\nabla f(x^k)]_{i_k}) \leftarrow \text{SEGA}$$

$$x^{k+1} = x^k - \alpha g^{k+1}$$

$$x^{k+1} = x^k - \alpha y^{k+1} \leftarrow \text{JAGUAR}$$

$$\mathbb{E}[y^{k+1} | x^k] \neq \nabla f(x^k)$$

$$\mathbb{E}[g^{k+1} | x^k] = y^k - d \mathbb{E}[y_{i_k}^k - [\nabla f(x^k)]_{i_k} | x^k]$$

$$= y^k - d \sum_{j=1}^d \frac{1}{d} (y_j - [\nabla f(x^k)]_j) e_j$$

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

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

$$\nabla f_i(x^*) \rightarrow 0$$

$$\frac{dL}{\mu}$$