

Методы оптимизации. Повторение лекций.

Корнилов Никита Максимович

Московский физико-технический институт

Декабрь 2025г

Скорости сходимости

- Сублинейная: $\|x^k - x^*\|_2 \leq \frac{C}{k^\alpha}$, $C > 0$, $\alpha > 0$.
- Линейная: $\|x^k - x^*\|_2 \leq Cq^k$, $C > 0$, $0 < q < 1$.
- Сверхлинейная: $\|x^k - x^*\|_2 \leq Cq^{kp}$, $C > 0$, $0 < q < 1$, $p > 1$.
- Квадратичная: $\|x^k - x^*\|_2 \leq Cq^{2^k}$, $C > 0$, $0 < q < 1$. Или
 $\|x^{k+1} - x^*\| \leq C\|x^k - x^*\|^2$, $C > 0$.

Свойства L -гладких и μ -сильно выпуклых функций

Пусть f — μ -сильно выпуклая, тогда

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \mu \|x - y\|_2^2.$$

Пусть f — L -гладкая, тогда

$$|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq \frac{L}{2} \|x - y\|_2^2.$$

Пусть f — L -гладкая и μ -сильно выпуклая, тогда

$$LI \succeq \nabla^2 f(x) \succeq \mu I.$$

Пусть f — L -гладкая и μ -сильно выпуклая, тогда

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{\mu L}{\mu + L} \|x - y\|_2^2 + \frac{1}{\mu + L} \|\nabla f(x) - \nabla f(y)\|_2^2.$$

Theorem

Пусть дана выпуклая непрерывно дифференцируемая на \mathbb{R}^d функция $f : \mathbb{R}^d \rightarrow \mathbb{R}$. Если для некоторой точки $x^* \in \mathbb{R}^d$ верно, что $\nabla f(x^*) = 0$, то x^* – глобальный минимум f на всем \mathbb{R}^d .

Theorem

Пусть дана выпуклая непрерывно дифференцируемая на \mathbb{R}^d функция $f : \mathbb{R}^d \rightarrow \mathbb{R}$ и выпуклое множество \mathcal{X} . Тогда $x^* \in \mathcal{X}$ – глобальный минимум f на \mathcal{X} тогда и только тогда, когда для всех $x \in \mathcal{X}$ выполнено

$$\langle \nabla f(x^*), x - x^* \rangle \geq 0.$$

Начало безусловной оптимизации. Градиентный спуск

$$x^{k+1} = x^k - \gamma_k \nabla f(x^k)$$

Шаг: оптимальное значение $\gamma_k = \frac{1}{2L}$.

Интуиция: решение системы $\frac{dx_t}{dt} = (x_t)dt$ или минимизация ограничивающей параболы

$$x^{k+1} = \arg \min_x \left\{ f(x^k) + \langle f(x^k), x - x^k \rangle + \frac{L}{2} \|x - x^k\|_2^2 \right\}$$

Сходимость для L -гладких функций:

$$f(x^K) - f(x^*) \leq \frac{2L\|x^0 - x^*\|_2^2}{K}.$$

Сходимость для L -гладких и μ -сильно выпуклых функций:

$$\|x^K - x^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^K \|x^0 - x^*\|_2^2.$$

Метод тяжелого шарика

$$x^{k+1} = x^k - \gamma_k \nabla f(x^k) + \tau_k(x^k - x^{k-1})$$

Интуиция: использовать инерцию траектории

Сходимость: не лучше градиентного спуска в теории

Ускоренный метод Нестерова

$$x^{k+1} = x^k - \gamma_k \nabla f(x^k + \tau_k(x^k - x^{k-1})) + \tau_k(x^k - x^{k-1})$$

Сходимость с $\gamma_k = \frac{1}{L}$, $\tau_k = \frac{\sqrt{L}-\sqrt{\mu}}{\sqrt{L}+\sqrt{\mu}}$:

$$f(x^K) - f(x^*) \leq \left(1 - \sqrt{\frac{\mu}{L}}\right)^K \cdot L \|x^0 - x^*\|_2^2.$$

Сходимость с $\gamma_k = \frac{1}{L}$, $\tau_k = \frac{k}{k+3}$:

$$f(x^K) - f(x^*) \leq \frac{4L\|x^0 - x^*\|_2^2}{(K+2)^2}.$$

Особенность: нужно подбирать два параметра, совпадает с нижними оценками

Метод сопряженных градиентов

Используется для решения СЛАУ: $Ax = b, A \succ 0$.

Идея: разложить решение в базис сопряженных относительно A направлений $x^* = \sum_{i=0}^{d-1} \lambda_i p_i$, восстанавливая на каждой итерации p^k, λ_k .

По сопряженности:

$$\lambda_j = \frac{p_j^T b}{p_j^T A p_j}, \quad x^{k+1} = x^k + \alpha_k p_k.$$

По индукции доказываем, что след направление - сопряженное всем

$$r_k = Ax^k - b = \nabla f(x^k), \quad p_k = -r_k + \beta_k p_{k-1}.$$

Сопряженность p_{k-1}, p_k :

$$\beta_k = \frac{p_{k-1}^T A r_k}{p_{k-1}^T A p_{k-1}}.$$

Метод сопряженных градиентов: особенности

- С квадратной положительно определенной матрицей размера d находит точное решение за не более чем d итераций.
- Сходимость с ошибкой

$$\|x^k - x^*\|_A \leq 2 \left(\frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^k \|x^0 - x^*\|_A.$$

Здесь $\|x\|_A^2 = x^T A x$ и $\kappa(A) = \lambda_{\max}(A)/\lambda_{\min}(A)$.

- Обобщения для произвольных функций: $r_k = \nabla f(x^k)$, α_k - правило подбора шага, $\beta_{k+1} = \frac{\langle \nabla f(x^{k+1}), \nabla f(x^{k+1}) \rangle}{\langle \nabla f(x^k), \nabla f(x^k) \rangle}$ (Флетчер - Ривс) или $\beta_{k+1} = \frac{\langle \nabla f(x^{k+1}), \nabla f(x^{k+1}) - \nabla f(x^k) \rangle}{\langle \nabla f(x^k), \nabla f(x^k) \rangle}$ (Полак - Рибьер). Полезно использовать рестарты.

Метод Ньютона

$$x^{k+1} = x^k - \left(\nabla^2 f(x^k) \right)^{-1} \nabla f(x^k)$$

Интуиция:

$$x^{k+1} = \arg \min_{y \in \mathbb{R}^d} \left[f(x^k) + \langle \nabla f(x^k), y - x^k \rangle + \frac{1}{2} \langle \nabla^2 f(x^k)(y - x^k), y - x^k \rangle \right]$$

Сходимость: Квадратичная, но локальная скорость

$$\|x^{k+1} - x^*\|_2 \leq \frac{M}{2\mu} \|x^k - x^*\|_2^2.$$

Дорогая итерация $O(d^3)$, вне области сходимости может расходиться, но демпинг $\gamma_k \neq 1$ помогает

Квазиньютоновские методы (Итерация за $O(d^2)$)

Идея: приблизительный, но быстрый пересчёт обратного гессиана.

Квазиньютоновское урав: $s^k = x^{k+1} - x^k$ и $y^k = \nabla f(x^{k+1}) - \nabla f(x^k)$:

$$s^k = H_{k+1}y^k.$$

Симметричность $H_{k+1} = H_{k+1}^\top$.

SR-1:

$$H_{k+1} = H_k + \mu_k q^k (q^k)^T,$$

где $\mu_k \in \mathbb{R}$ и $q^k \in \mathbb{R}^d$ нужно подобрать.

BFGS:

$$H_{k+1} = \arg \min_{H \in \mathbb{R}^{d \times d}} \|H - H_k\|^2$$

$$s.t. \quad s^k = Hy^k$$

$$H^T = H$$

Начало оптимизации на выпуклых множествах. Проекционный GD

Ставится задача минимизации на множестве $\min_{x \in X} f(x)$.

Проекция на выпуклое замкнутое мн-во:

$$\Pi(x) := \arg \min_{y \in X} \frac{1}{2} \|x - y\|_2^2, \quad \|\Pi_X(x_1) - \Pi_X(x_2)\|_2 \leq \|x_1 - x_2\|_2.$$

GD с проекцией:

$$x^{k+1} = \Pi_X \left[x^k - \gamma_k \nabla f(x^k) \right].$$

Итерационная сходимость как у GD: $\|x^K - x^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^K \|x^0 - x^*\|_2^2$.

Примеры проекций с готовыми решениями:

- ℓ_2 -шар радиуса R с центром в 0: $\Pi_X(x) = \min \left\{ 1, \frac{R}{\|x\|_2} \right\} x$.
- $X = \{y \in \mathbb{R}^d \mid Ay = b\}$: $\Pi_X(x) = x - A^T(AA^T)^{-1}(Ax - b)$

Метод Франка Вульфа

$$s^k = \arg \min_{s \in X} \langle s, \nabla f(x^k) \rangle$$

$$\gamma_k = \frac{2}{k+2}$$

$$x^{k+1} = (1 - \gamma_k)x^k + \gamma_k s^k$$

Сходимость итерационная:

$$f(x^K) - f(x^*) \leq \frac{2 \max\{L \operatorname{diam}(X)^2, f(x^0) - f(x^*)\}}{K+2},$$

где $\operatorname{diam}(X) := \max_{x,y \in X} \|x - y\|_2$ – диаметр множества X .

Примеры подзадачи с готовыми решениями:

- ℓ_1 -шар радиуса R с центром в 0 :
 $y^* = -R \operatorname{sign}(x_i) \mathbf{e}_i, i = \arg \max_j |x_j|,$
- Симплекс $\Delta = \left\{ y \in \mathbb{R}^d \mid y_i \geq 0, \sum_{i=1}^d y_i = R \right\}$:

$$y^* = R \mathbf{e}_i, \text{ где } i = \arg \min x_i.$$

Зеркальный спуск

Ставится задача минимизации на множестве $\min_{x \in X} f(x)$.

Идея: обобщить понятия расстояний и проекции, использовав геометрию задачи и выиграв в константе сложности

Дивергенцией Брэгмана (аналог метрики)

Пусть дана дифференцируемая 1-сильно выпуклая относительно нормы $\|\cdot\|$ на множестве X функция d . Дивергенцией Брэгмана $V(x, y) : X \times X \rightarrow \mathbb{R}$ такая, что для любых $x, y \in X$

$$V(x, y) = d(x) - d(y) - \langle \nabla d(y), x - y \rangle.$$

Примеры:

$$d(x) = \frac{1}{2}\|x\|_2^2 \implies V(x, y) = \frac{1}{2}\|x - y\|_2^2.$$

$$d(x) = \sum_{i=1}^d x_i \log x_i \implies V(x, y) = \sum_{i=1}^d x_i \log \frac{x_i}{y_i}.$$

Зеркальный спуск 2

$$x^{k+1} = \arg \min_{x \in X} \{ \langle \gamma \nabla f(x^k), x \rangle + V(x, x^k) \}$$

Эквивалентная запись

$$x^{k+1} = P_{V(\cdot, \cdot)} \left[(\nabla d)^{-1} (\nabla d(x^k) - \gamma \nabla f(x^k)) \right].$$

Сходимость: на выпуклом множестве X с L -гладкой относительно нормы $\|\cdot\|$, выпуклой целевой функцией f и шагом $\gamma \leq \frac{1}{L}$

$$f \left(\frac{1}{K} \sum_{k=1}^K x^k \right) - f(x^*) \leq \frac{V(x^*, x^0)}{\gamma K}.$$

На единичном симплексе:

$$x_i^{k+1} = x_i^* = \frac{x_i^k \exp(-\gamma [\nabla f(x^k)]_i)}{\sum_{i=1}^d x_i^k \exp(-\gamma [\nabla f(x^k)]_i)}.$$

Негладкая задача

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

где f выпуклая и M -Липшицева.

Субградиентный метод:

$$g^k \in \partial f(x^k), \quad x^{k+1} = x^k - \gamma g^k.$$

Сходимость:

$$f\left(\frac{1}{K} \sum_{k=0}^{K-1} x^k\right) - f(x^*) \leq \frac{M \|x^0 - x^*\|_2}{\sqrt{K}}.$$

Оптимальная оценка (но медленнее GD для гладких функций) и маленький шаг $\gamma = \frac{\|x^0 - x^*\|_2}{M\sqrt{K}}$.

Проксимальный оператор

$$\text{prox}_r(x) = \arg \min_{\tilde{x} \in \mathbb{R}^d} \left(r(\tilde{x}) + \frac{1}{2} \|x - \tilde{x}\|^2 \right).$$

Примеры:

- $r(x) = \lambda \|x\|_1$, тогда $[\text{prox}_r(x)]_i = [|x_i| - \lambda]_+ \cdot \text{sign}(x_i)$,
- $r(x) = \frac{\lambda}{2} \|x\|_2^2$, тогда $\text{prox}_r(x) = \frac{x}{1+\lambda}$.

Свойства:

- $\text{prox}_r(x) = y \iff x - y \in \partial r(y)$.
- $\langle x - y, \text{prox}_r(x) - \text{prox}_r(y) \rangle \geq \|\text{prox}_r(x) - \text{prox}_r(y)\|_2^2$,
- $\|\text{prox}_r(x) - \text{prox}_r(y)\|_2 \leq \|x - y\|_2$.

Композитная задача и проксимальный метод

$\min_{x \in \mathbb{R}^d} [f(x) + r(x)]$ – композитная задача,

где f является L -гладкой выпуклой функцией, r выпуклой (необязательно гладкой, но проксимально дружественной функцией).

Проксимальный метод:

$$x^{k+1} = \text{prox}_{\gamma r}(x^k - \gamma \nabla f(x^k)).$$

Альтернативная запись:

$$x^{k+1} = x^k - \gamma(\nabla f(x^k) + \partial r(x^{k+1})).$$

Сходимость: Композитная задача с L -гладкой, μ -сильно выпуклой целевой функцией f и выпуклой r при $\gamma_k = \frac{1}{L}$:

$$\|x^K - x^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^K \|x^0 - x^*\|_2^2.$$

Точка x^* — неподвижная точка для $\text{prox}_{\gamma r}(\cdot - \gamma \nabla f(\cdot))$

Начало оптимизации с ограничениями. Штрафная функция

Рассмотрим следующую задачу с ограничениями:

$$\begin{aligned} & \min_{x \in \mathbb{R}^d} f(x), \\ \text{s.t. } & h_i(x) = 0, \quad i = 1, \dots, m, \\ & g_i(x) \leq 0, \quad j = 1, \dots, n. \end{aligned}$$

Аугментация:

$$\min_{x \in \mathbb{R}^d} \left[f_\rho(x) = f(x) + \rho \cdot \frac{1}{2} \sum_{i=1}^m h_i^2(x) + \rho \cdot \frac{1}{2} \sum_{j=1}^n (g_j^+)^2(x) \right].$$

Пусть все функции являются непрерывными и $\{x \in \mathbb{R}^d \mid f(x) \leq f(x^*)\}$ ограничено. Тогда для любого $e > 0$ существует $\rho(e) > 0$ такое, что множество решений штрафной задачи X_ρ^* для любых $\rho \geq \rho(e)$ содержится в

$$X_e^* = \{x \in \mathbb{R}^d \mid \exists x^* \in X^* : \|x - x^*\|_2 \leq e\}.$$

Метод штрафных функций

Алгоритм:

- ① решить задачу для текущего ρ ,
- ② увеличить ρ ,
- ③ использовать предыдущее решение как начальную точку.

Особенности:

- Условная задача превращена в безусловную.
- Увеличение ρ приближает к исходной задаче.
- При большом ρ наблюдается нарушение ограничений, что подходит не для всех задач.
- Увеличение ρ влечет за собой увеличение обусловленности задачи (константа Липшица градиента будет сильно расти). А значит задачу будет сложнее решать.

ADMM

$$\begin{aligned} & \min_{x \in \mathbb{R}^{d_x}, y \in \mathbb{R}^{d_y}} f(x) + g(y), \\ & \text{s.t. } Ax + By = c, \end{aligned}$$

где $A \in \mathbb{R}^{n \times d_x}$, $B \in \mathbb{R}^{n \times d_y}$, $c \in \mathbb{R}^n$. Аугментация

$$\begin{aligned} & \min_{x \in \mathbb{R}^{d_x}, y \in \mathbb{R}^{d_y}} f(x) + g(y) + \frac{\rho}{2} \|Ax + By - c\|_2^2, \\ & \text{s.t. } Ax + By = c, \end{aligned}$$

Алгоритм:

$$\begin{aligned} x^{k+1} &= \arg \min_{x \in \mathbb{R}^{d_x}} L_\rho(x, y^k, \lambda^k) \\ y^{k+1} &= \arg \min_{y \in \mathbb{R}^{d_y}} L_\rho(x^{k+1}, y, \lambda^k) \\ \lambda^{k+1} &= \lambda^k + \rho (Ax^{k+1} + By^{k+1} - c) \end{aligned}$$

Вернуть $\frac{1}{K} \sum_{k=1}^K x^k$, $\frac{1}{K} \sum_{k=1}^K y^k$, $\frac{1}{K} \sum_{k=1}^K \lambda^k$.

Theorem

Если функции f и g являются выпуклыми и дружественными с точки зрения вычислений $\arg \min$, то ADMM имеет следующую оценку сходимости для любого $x \in \mathbb{R}^{d_x}$, $y \in \mathbb{R}^{d_y}$, $\lambda \in \mathbb{R}^n$

$$L_0 \left(\frac{1}{K} \sum_{k=1}^K x^k, \frac{1}{K} \sum_{k=1}^K y^k, \lambda \right) - L_0 \left(x, y, \frac{1}{K} \sum_{k=1}^K \lambda^k \right) \leq \frac{1}{2K} \|z^0 - z\|_P^2,$$

где L_0 – Лагранжиан без аугментации, $P = \begin{pmatrix} \rho A^T A & 0 & -A^T \\ 0 & 0 & 0 \\ -A & 0 & \frac{1}{\rho} I \end{pmatrix}$,

$$z^0 = \begin{pmatrix} x^0 \\ y^0 \\ \lambda^0 \end{pmatrix}$$

Барьерная функция

$$\begin{aligned} & \min_{x \in \mathbb{R}^d} f(x), \\ & \text{s.t. } g_i(x) \leq 0, \quad i = 1, \dots, m. \end{aligned}$$

Барьером будем называть функцию $F : \mathbb{R}^d \rightarrow \mathbb{R}$:

- F непрерывно дифференцируема на $\text{int}G$;
- Для любой последовательности $\{x_i\} \in \text{int}G$ такой, что $x_i \rightarrow x \in \partial G$ (граница множества G), выполнено $F(x_i) \rightarrow +\infty$.

Примеры:

- Барьер Кэррола: $F(x) = -\sum_{i=1}^m \frac{1}{g_i(x)}$;
- Логарифмический барьер: $F(x) = -\sum_{i=1}^m \ln(-g_i(x))$.

$$\min_{x \in \text{int}(G)} \left[F_\rho(x) = f(x) + \frac{1}{\rho} F(x) \right].$$

Метод внутренней точки

Сходимость: Для любого $\epsilon > 0$ существует $\rho(\epsilon) > 0$ такое, что множество решений барьерной задачи X_ρ^* для любых $\rho \geq \rho(\epsilon)$ содержится в

$$X_\epsilon^* = \{x \in G \mid \exists x^* \in X^* : \|x - x^*\|_2 \leq \epsilon\},$$

где X^* – множество решений исходной задачи оптимизации с ограничениями вида неравенств.

Алгоритм:

- ① Увеличить $\rho_k > \rho_{k-1}$
- ② С помощью некоторого метода решить численно задачу безусловной оптимизации с целевой функцией F_{ρ_k} и стартовой точкой x_k . Гарантировать, что выход метода x_{k+1} будет близок к реальному решению $x^*(\rho_k)$.

Всегда соблюдаем ограничения неравенства !

Седловая задача и Экстреградиент

$$\min_{x \in \mathbb{R}^d} \max_{\lambda \in \mathbb{R}^d} L(x, \lambda),$$

где L непрерывно дифференцируема по обеим группам переменных, выпукла-вогнута: выпукла по x (для любого фиксированного λ) и вогнута по λ (для любого фиксированного x)

Обобщения GD расходятся или сходятся неоптимально ($L(x, y) = xy$).

Экстраградиент

$$x^{k+1/2} = x^k - \gamma \nabla_x L(x^k, \lambda^k)$$

$$\lambda^{k+1/2} = \lambda^k + \gamma \nabla_\lambda L(x^k, \lambda^k)$$

$$x^{k+1} = x^{k+1/2} - \gamma \nabla_x L(x^{k+1/2}, \lambda^{k+1/2})$$

$$\lambda^{k+1} = \lambda^{k+1/2} + \gamma \nabla_\lambda L(x^{k+1/2}, \lambda^{k+1/2})$$

Вернуть $\frac{1}{K} \sum_{k=0}^{K-1} x^{k+1/2}, \frac{1}{K} \sum_{k=0}^{K-1} \lambda^{k+1/2}$

Экстраградиент: сходимость

Для любого $u \in \mathbb{R}^d \times \mathbb{R}^n$ и для любого $\gamma \leq \frac{1}{L}$:

$$\left(L\left(\frac{1}{K} \sum_{k=0}^{K-1} x^{k+1/2}, u_\lambda\right) - L\left(u_x, \frac{1}{K} \sum_{k=0}^{K-1} \lambda^{k+1/2}\right) \right) \leq \frac{\|z^0 - u\|_2^2}{2\gamma K}$$

Метрика для решения на компактах: $\max_\lambda L(x^k, \lambda) - \min_x L(x, \lambda^k)$

- Можно добавить проекции и решать седловую задачу на множествах $X \neq \mathbb{R}^d$ и $\Lambda \neq \mathbb{R}^n$.
- Можно получить линейную сходимость для сильно выпуклых–сильно вогнутых задач.
- Часто применяют для решения двойственной задачи.

Стохастическая оптимизация

Онлайн-постановка:

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

Оффлайн-постановка:

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

SGD:

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

Предположения:

$$\mathbb{E}_\xi[\nabla f(x, \xi)] = \nabla f(x), \quad \mathbb{E}_\xi[\|\nabla f(x, \xi) - \nabla f(x)\|_2^2] \leq \sigma^2.$$

SGD Сходимость

Задача безусловной стохастической оптимизации с L -гладкой, μ -сильно выпуклой функцией f решается с помощью SGD с $\gamma_k \leq \frac{1}{L}$:

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

Постоянный шаг:

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

Уменьшающийся шаг: $\gamma_k = \frac{1}{k+1}$ или $\gamma_k = \frac{1}{\sqrt{k+1}}$. Плюс: точнее сходимость, минус: потеря линейной сходимости в начале.

Batching:

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

Уменьшение дисперсии с σ до σ/\sqrt{b} .

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

SAGA

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

Идея: уменьшаем с каждой итерацией дисперсию, чтобы
 $g^k \rightarrow \nabla f(x^*) = 0$, при $x^k \rightarrow x^*$.

Алгоритм SAGA:

Сгенерировать независимо i_k

$$g^k = \nabla f_{i_k}(x^k) - y_{i_k}^k + \frac{1}{n} \sum_{j=1}^n y_j^k$$

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

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

$$\frac{1}{n} \sum_{j=1}^n y_j^k - \text{«запаздывающая» версия } \nabla f(x^k), \mathbb{E}[g^k | x^k] = \nabla f(x^k).$$

О сходимости SAGA

При $x^k \rightarrow x^*$ имеем, что $y_j^k \rightarrow \nabla f_j(x^*)$, и $\frac{1}{n} \sum_{j=1}^n y_j^k \rightarrow \nabla f(x^*) = 0$,
 $\nabla f_{i_k}(x^k) \rightarrow \nabla f_j(x^*)$, значит $g^k \rightarrow 0$.

Сходимость. 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 получается следующая итерационная сложность:

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