

Стохастические безградиентные методы для седловых задач

Садиев Абдурахмон

Московский физико-технический институт
Фихтех-школа прикладной математики и информатики
Кафедра интеллектуальных систем

Научный руководитель д.ф.-м.н. А. В. Гасников

Москва,
2020 г.

- Для стохастической седловой задачи предложить метод, использующий оракул нулевого порядка, то есть имеется доступ только к значению функции в точке.
- Сравнить предложенный метод с градиентным аналогом.

- Ben-Tal A., Nemirovski A. Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications. 2019.
- Shamir O. An Optimal Algorithm for Bandit and Zero-Order Convex Optimization with Two-Point Feedback. // Journal of Machine Learning Research. 2017. Vol. 18, no. 52. P. 1–11
- Gasnikov A. Universal gradient descent // arXiv preprint arXiv:1711.00394. 2017.

- Седловая задача:

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} \varphi(x, y)$$

- $\mathcal{X} \subset \mathbb{R}^{n_x}$, $\mathcal{Y} \subset \mathbb{R}^{n_y}$ - выпуклые компактные множества.
- Функция $\varphi(\cdot, y)$ - выпуклая на \mathcal{X}
- Функция $\varphi(x, \cdot)$ - вогнутая на \mathcal{Y}
- Функция $\varphi(x, y)$ - M -липшицево непрерывная

- Данную задачу можно решать алгоритмом зеркального спуска для седловых задач (Mirror Descent, А. Немировский)¹.
- Метод использует оракул первого порядка (выдает значение градиента в точке).
- Количество итераций N , необходимых для нахождения решения с точностью ε : $\mathcal{O}\left(\frac{\Omega^2 M^2}{\varepsilon^2}\right)$, где Ω - диаметр Брегмана² множества $\mathcal{X} \times \mathcal{Y}$.

¹ Ben-Tal A., Nemirovski A. Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications. 2019.

² Определение будет дано ниже.

- Неточный стохастический оракул нулевого порядка:

$$\tilde{\varphi}(x, y, \xi) = \varphi(x, y, \xi) + \delta(x, y),$$

$$\mathbb{E}_{\xi}[\tilde{\varphi}(x, y, \xi)] = \tilde{\varphi}(x, y), \quad \mathbb{E}_{\xi}[\varphi(x, y, \xi)] = \varphi(x, y),$$

где случайная переменная ξ отвечает за несмещенный стохастический шум, а $\delta(x, y)$ – за детерминистический шум.

- Ограничения:

$$\|\nabla\varphi(x, y, \xi)\|_2 \leq M(\xi), \quad \mathbb{E}[M^2(\xi)] = M^2$$

- Обозначим $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, тогда $z \in \mathcal{Z}$ означает $z = (x, y)$, где $x \in \mathcal{X}$, $y \in \mathcal{Y}$
- $\varphi(z) = \varphi(x, y)$, и $\varphi(z, \xi) = \varphi(x, y, \xi)$.
- Оценка градиента:

$$g(z, \xi, \mathbf{e}) = \frac{n}{2\tau} (\tilde{\varphi}(z + \tau\mathbf{e}, \xi) - \tilde{\varphi}(z - \tau\mathbf{e}, \xi)) \begin{pmatrix} \mathbf{e}_x \\ -\mathbf{e}_y \end{pmatrix},$$

где $\mathbf{e} \in \mathcal{RS}_2^n(1)$ (случайный вектор, равномерно распределенный на евклидовой сфере) и некоторая положительная константа τ .

- Сглаженная версия функции $\varphi(z)$:

$$\hat{\varphi}(z) = \mathbb{E}_{\mathbf{e}} [\varphi(z + \tau \mathbf{e})]$$

- Свойства:
 - 1 Сглаженная версия функции является так же выпукло-вогнутой функцией
 - 2 Сглаженная версия функции является непрерывно дифференцируемой функцией

Определение

- Функция $d(z) : \mathcal{Z} \rightarrow \mathbb{R}$ называется прокс-функцией, если $d(z)$ является 1-сильно выпуклой по отношению к $\|\cdot\|$ -норме и дифференцируемой на \mathcal{Z} функцией.
- Дивергенция Брегмана:
$$V_z(w) = d(z) - d(w) - \langle \nabla d(w), z - w \rangle.$$
- Прокс-оператор: $\text{prox}_x(\xi) = \arg \min_{y \in \mathcal{Z}} (V_x(y) + \langle \xi, y \rangle)$
- Диаметр Брегмана: $\Omega_{\mathcal{Z}}$ множества \mathcal{Z} по отношению к $V_{z_1}(z_2)$:

$$\Omega_{\mathcal{Z}} = \max\{\sqrt{2V_{z_1}(z_2)} : z_1, z_2 \in \mathcal{Z}\}$$

Algorithm 1 Zeroth-Order Saddle-Point Algorithm (zoSPA)

Input: Iteration limit N .

Let $z_1 = \arg \min_{z \in \mathcal{Z}} d(z)$.

for $k = 1, 2, \dots, N$ **do**

 Sample \mathbf{e}_k, ξ_k independently.

 Initialize γ_k .

$z_{k+1} = \text{prox}_{z_k}(\gamma_k g(z_k, \xi_k, \mathbf{e}_k))$

end for

Output: \bar{z}_N ,

где

$$\bar{z}_N = \frac{1}{\Gamma_N} \left(\sum_{k=1}^N \gamma_k z_k \right), \quad \Gamma_N = \sum_{k=1}^N \gamma_k.$$

Лемма 1

Для $g(z, \xi, \mathbf{e})$ выполнены следующие условия:

$$\mathbb{E} [\|g(z, \xi, \mathbf{e})\|_q^2] \leq 2 \left(cnM^2 + \frac{n^2 \Delta^2}{\tau^2} \right) a_q^2,$$

где c - некоторая положительная константа (независимо от n), а a_q определяется следующим образом:

$$a_q^2 = \min\{2q - 1, 32 \log n - 8\} n^{\frac{2}{q}-1}, \quad \forall n \geq 3$$

Лемма 2

Определим $\Delta_k = g(z_k, \xi_k, \mathbf{e}_k) - \mathbb{E}_{\mathbf{e}_k} [g(z_k, \xi_k, \mathbf{e}_k)]$. Пусть $D(u) = \sum_{k=1}^N \gamma_k \langle \Delta_k, u - z_k \rangle$. Тогда мы имеем

$$\mathbb{E} \left[\max_{u \in \mathcal{Z}} D(u) \right] \leq \Omega^2 + \frac{\Delta \Omega n a_q}{\tau} \sum_{k=1}^N \gamma_k + M_{all}^2 \sum_{k=1}^N \gamma_k^2.$$

Теорема

Пусть шаг Алгоритма 1 $\gamma_k = \frac{\Omega}{M_{all}\sqrt{N}}$. Тогда скорость сходимости Алгоритма 1:

$$\mathbb{E} [\varepsilon_{sad}(\bar{z}_N)] \leq \frac{3M_{all}\Omega}{\sqrt{N}} + \frac{\Delta\Omega na_q}{\tau} + 2\tau M,$$

где Ω есть диаметр множества \mathcal{Z} , $M_{all}^2 = 2 \left(cnM^2 + \frac{n^2\Delta^2}{\tau^2} \right) a_q^2$ и

$$\varepsilon_{sad}(\bar{z}_N) = \max_{y' \in \mathcal{Y}} \varphi(\bar{x}_N, y') - \min_{x' \in \mathcal{X}} \varphi(x', \bar{y}_N),$$

Следствие

При допущениях Теоремы пусть ε будет точность решения седловой задачи, полученная с помощью алгоритма 1.

Предположим, что

$$\tau = \Theta\left(\frac{\varepsilon}{M}\right), \quad \Delta = \mathcal{O}\left(\frac{\varepsilon^2}{M\Omega n a_q}\right),$$

тогда количество итераций N для нахождения ε -решения

$$N = \mathcal{O}\left(\frac{\Omega^2 M^2 n^{2/q}}{\varepsilon^2} C^2(n, q)\right),$$

где $C(n, q) = \min\{2q - 1, 32 \log n - 8\}$.

$p, (1 \leq p \leq 2)$	$q, (2 \leq q \leq \infty)$	$N, \text{ Количество итераций}$
$p = 2$	$q = 2$	$\mathcal{O}\left(\frac{\Omega^2 M^2}{\varepsilon^2} n\right)$
$p = 1$	$q = \infty$	$\mathcal{O}\left(\frac{\Omega^2 M^2}{\varepsilon^2} \log^2(n)\right)$

Сводка оценок сходимости для негладкого случая: $p = 2$ и $p = 1$.

- Количество итераций для зеркального спуска (Mirror Descent, А. Немировский)³: $\mathcal{O}\left(\frac{\Omega^2 M^2}{\varepsilon^2}\right)$

³ Ben-Tal A., Nemirovski A. Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications. 2019.

Постановка задачи

Решается классическая седловая задача

$$\min_{x \in \Delta_n} \max_{y \in \Delta_k} [y^T C x],$$

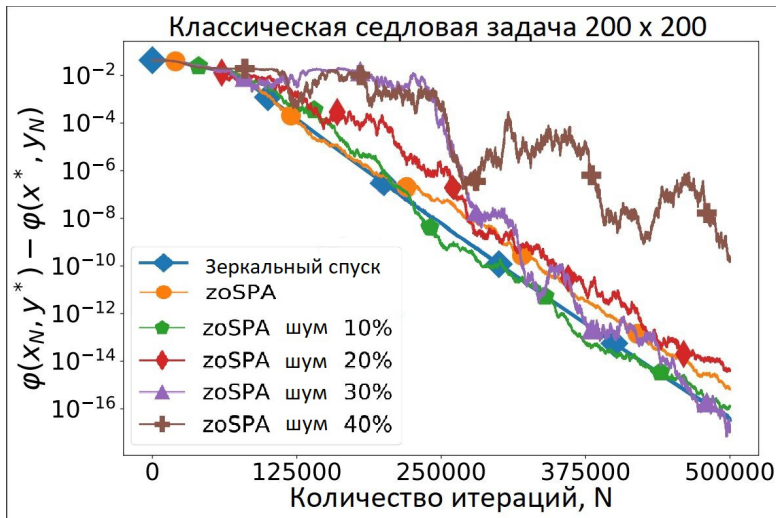
где $\Delta_n = \{w \in \mathbb{R}^n : \forall i \rightarrow w_i \geq 0, \sum_{i=1}^n w_i = 1\}$ - вероятностный симплекс.

Проксимальная настройка

Дивергенция Брегмана для данной задачи:

$$V_y(x) = \sum_{i=1}^n x_i \log x_i / y_i$$

расстояние Кульбака — Лейблера.



zoSPA с 0 - 40 % шума и Зеркальный спуск, примененные для решения седловой задачи.

Результаты

- 1 Представлен новый метод для решения негладкой седловой задачи.
- 2 Данный алгоритм использует оракул нулевого порядка со стохастическим и ограниченным детерминистическим шумом.
- 3 Показано, что количество итераций предложенного метода необходимых для нахождения решения с точностью ε отличается в $Const(n, q)$ от градиентного аналога (Mirror Descent).

Публикации

Beznosikov A., Sadiev A., Gasnikov A.
Gradient-Free Methods for Saddle-Point Problem. 2020.
arXiv:math.OA/2005.05913.