

# Алгоритм Natasha-2

2 мая 2020 г.

## 1 Постановка задачи

$$\min_{x \in \mathbb{R}} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) = \mathbb{E}_i[f_i(x)],$$

где  $f_i$  — в общем случае невыпуклые функции.

Будет построен онлайн-алгоритм (т.е. не зависящий от  $n$ ), находящий локальный минимум с заданной точностью.

## 2 Необходимые предположения

- $f$  имеет ограниченную дисперсию, откуда  $\mathbb{E}_i[\|\nabla f_i(x) - \nabla f(x)\|] \leq \sigma$
- $L$ -Липшицевость градиента:  $\|\nabla f(x) - \nabla f(y)\| \leq L \cdot \|x - y\|$
- $L_2$ -Липшицевость гессиана:  $\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L_2 \cdot \|x - y\|$

## 3 Определения

- Локальный минимум с  $(\varepsilon, \delta)$ -точностью - точка  $x$ , такая что  $\|\nabla f(x)\| < \varepsilon$  и  $\nabla^2 f(x) \succeq -\delta \mathbf{I}$  (все собственные значения гессиана больше  $-\delta$ ).
- $\sigma$ -невыпуклая функция - функция, такая что минимальное собственное значение её гессиана превосходит  $-\sigma$

## 4 Описание алгоритма

Главная идея - использование гессиана для избегания седловых точек.

### 4.1 Carmon: Сведение к поиску стационарных точек

На каждой итерации в точке  $y_k$  сравниваем минимальное собственное значение  $\lambda_{\min}(\nabla^2 f(y_k))$  и  $-\delta$ .

- Если  $\exists v : v^T \nabla^2 f(y_k) v \leq -\frac{\delta}{2}$ , то пытаемся отдалиться от седла:  $y_{k+1} = y_k \pm \frac{\delta}{L_2} v$  (+ или - выбираем случайно)
- Иначе берём  $F^k(x) = f(x) + L(\max\{0, \|x - y_k\| - \frac{\delta}{L_2}\})^2$ , которая является  $5L$ -Липшицевой по градиенту и  $3\delta$ -невыпуклой; ищем для неё точку  $\hat{x} : \|\nabla F^k(\hat{x})\| \leq \varepsilon$  и переходим туда:  $y_{k+1} = \hat{x}$ . Как видно,  $F^k$  штрафует за выход из «безопасной зоны»  $\{x : \|x - y_k\| \leq \frac{\delta}{L_2}\}$

### 4.2 Natasha 1.5: Поиск стационарных точек для $\sigma$ -невыпуклых функций

Каждая эпоха делится на  $p = \Theta((\frac{\sigma}{\varepsilon L})^{\frac{2}{3}})$  суб-эпох со стартовым вектором  $\hat{x}$ .

Затем, при подсчёте градиентов, вместо  $f(x)$  фактически используется  $f(x) + \sigma\|x - \hat{x}\|^2$ , что должно стабилизировать алгоритм, немного замедляя его. Вместе с этим, используется пересчёт градиента из алгоритма SVRG

Псевдокод:

---

**Algorithm 1** Natasha1.5( $F, x^\varnothing, B, T', \alpha$ )

---

**Input:**  $f(\cdot) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ , starting vector  $x^\varnothing$ , epoch length  $B \in [n]$ , epoch count  $T' \geq 1$ , learning rate  $\alpha > 0$ .

```

1:  $\hat{x} \leftarrow x^\varnothing$ ;  $p \leftarrow \Theta((\sigma/\varepsilon L)^{2/3})$ ;  $m \leftarrow B/p$ ;  $X \leftarrow []$ ;
2: for  $k \leftarrow 1$  to  $T'$  do  $\diamond$   $T'$  epochs each of length  $B$ 
3:    $\tilde{x} \leftarrow \hat{x}$ ;  $\mu \leftarrow \frac{1}{B} \sum_{i \in S} \nabla f_i(\tilde{x})$  where  $S$  is a uniform random subset of  $[n]$  with  $|S| = B$ ;
4:   for  $s \leftarrow 0$  to  $p - 1$  do  $\diamond$   $p$  sub-epochs each of length  $m$ 
5:      $x_0 \leftarrow \tilde{x}$ ;  $X \leftarrow [X, \tilde{x}]$ ;
6:     for  $t \leftarrow 0$  to  $m - 1$  do
7:        $\tilde{\nabla} \leftarrow \nabla f_i(x_t) - \nabla f_i(\tilde{x}) + \mu + 2\sigma(x_t - \tilde{x})$  where  $i \in_R [n]$ 
8:        $x_{t+1} = x_t - \alpha \tilde{\nabla}$ ;
9:     end for
10:     $\hat{x} \leftarrow$  a random choice from  $\{x_0, x_1, \dots, x_{m-1}\}$ ;  $\diamond$  in practice, choose the average
11:  end for
12: end for
13:  $\hat{y} \leftarrow$  a random vector in  $X$ .  $\diamond$  in practice, simply return  $\hat{y}$ 

```

### 4.3 Оја: Быстрый поиск максимального собственного значения