

• Проблема разнор. оптимизации - экстремальное комбинирование

Все локальные экстремумы экстремальны за всем переводом наоборот, но всегда интересно

Оптимизационный процесс - разнор. по норме, но решение

Parallel GD/Fed Avg/Local GD:

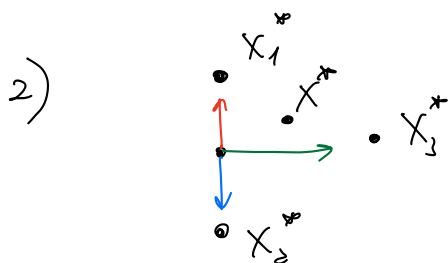
$$X_m^{k+1} = X_m^k - \gamma \nabla S_m(X_m^k) \quad \leftarrow \text{локальный макс по } S_m$$

при t итерации $X^k = \frac{1}{M} \sum_{m=1}^M X_m^k$, $X_m^k = X^k \quad \leftarrow \text{гетерогенность}$

⊕ существует разнор. задача

⊖ нельзя создать метод по определению

Рисунок:



чем нужно, все идет в свое направление
 все близко локальных макс,
 не хуже

Теорема

$$\|\bar{X}^k - X^*\|_2^2 \leq (1 - \gamma\mu)^k \|X^0 - X^*\|_2^2 + \frac{\gamma \sigma_{\text{opt}}^2 (t-1)}{\mu}$$

$\bar{X}^k = \frac{1}{M} \sum_{m=1}^M X_m^k$

$\sigma_{\text{opt}}^2 = \frac{1}{M} \sum_{m=1}^M \|\nabla S_m(X^*)\|_2^2$

нельзя по X^*

создано к опр.

Хотим метод, который сходится к глоб. мин.

Идея: локальное мин = переобучение
регуляризатор = функция от переобучения

Локализация: минимизация $f_m(x)$
 \Downarrow регуляризатор.
минимизация $\tilde{f}_m(x) = f_m(x) + \frac{\lambda}{2} \|x - v\|_2^2$
регуляризатор

Fed Prox:

$$\begin{aligned} x_m^{k+1} &= x_m^k - \gamma (\nabla f_m(x_m^k) + \lambda (x_m^k - v^k)) \\ v^{k+1} &= v^k \\ \text{при } t \text{ итерации: } x^k &= \frac{1}{M} \sum_{m=1}^M x_m^k \quad x_m^k = x^k \quad v^k = x^k \end{aligned}$$

Fed Prox:

$$x_m^{k+1} = x_m^k - \gamma (\nabla f_m(x_m^k) + \lambda x_m^k - \lambda v^k)$$

Scaffold:

$$x_m^{k+1} = x_m^k - \gamma (\nabla f_m(x_m^k) - \underbrace{C_m^k + C^k}_{\text{регуляризатор}})$$

⊕ плюс локального GD

⊕ сходимость по решению

Проблема локальных минимумов:

они не нужны, тем GD

но в случае, когда жесткое локальное неиспользование
(f_m)
минимумов гораздо лучше, чем GD

- Как оценить нормальность?

$$1) \quad \|\nabla f_m(x) - \nabla f(x)\|_2 \leq \Delta \quad \forall x$$

это эквив. : $\|A_m x - A x\|_2 \sim \|A_m - A\|_2 \|x\|_2$ $x \in \mathbb{R}^d$,
но Δ не const.

$$2) \quad \|\nabla^2 f_m(x) - \nabla^2 f(x)\|_2 \leq \delta \quad \forall x$$

это эквив. : $\|A_m - A\|_2 \leq \delta$

Оценим δ :

$$\bullet \quad f(x) = \frac{1}{N} \sum_{m=1}^M f_m(x) = \frac{1}{M} \sum_{m=1}^M \left[\frac{1}{N} \sum_{i=1}^N \mathcal{L}(g(x, a_i^m), b_i^m) \right]$$

где b_i^m и a_i^m — const.

- выборка на f_m берется независимо от всей выборки в f

- кер. в. теорема (nonparametric)

$\{X_i\}_{i=1}^N$ — i.i.d. samples размер d

$$X_i^T = X_i \quad \mathbb{E}[X_i] = 0 \quad X_i^2 \preceq A^2$$

матрица с вероятностью $1-p$

$$\left\| \frac{1}{N} \sum X_i \right\|_2 \leq \sqrt{\frac{8 \|A\|_2^2}{N} \ln \frac{d}{p}}$$

↑
можно $\left\| \frac{1}{N} \sum X_i \right\|_2 \leq \frac{1}{N} \sum \|X_i\|_2 \leq \frac{1}{N} \sum \|A\|_2 = \|A\|_2$
↑
нужно.

- теорема для рекурсивных

$$f_m(x) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(x, c_i^m)$$

$$\triangleleft \quad \left\| \nabla^2 f_m(x) - \nabla^2 f(x) \right\|_2$$

• $\mathbb{E} = 0$ (из-за независимости)

• \mathcal{L} — L -гладкая, то $\nabla^2 \mathcal{L} \preceq L I$

- $X_i = \nabla^2 t - \nabla^2 f$

- $EX_i = 0, \quad X_i^2 \leq A^2 = 4L^2 \cdot I$

$$\Downarrow$$

$$\|\nabla^2 f_m(x) - \nabla^2 f(x)\|_2 \leq \delta \sim \frac{L}{\sqrt{N}}$$

δ малое, если
выборка на y -ве
большая

• Как в методе градиентного спуска?

Метод зрительного инкремента:

$$x^{[t+1]} = \operatorname{argmin}_{x \in \mathbb{R}^d} \{ \gamma \langle \nabla f(x^k); x \rangle + V(x, x^k) \}$$

$$V(x, y) = \varphi(x) - \varphi(y) - \langle \nabla \varphi(y); x - y \rangle$$

где φ задан с помощью лок. функции:

$$\varphi(x) = f_1(x) + \frac{\gamma}{2} \|x\|_2^2$$

argmin становится регрессией = нужно считать
меньше argmin на градиенте с f_1
всего 1 вычисление за итерацией

Скорость

$$\gamma = 1$$

$$V(x^*, x^k) \leq \left(1 - \frac{\mu}{\mu + 2\delta}\right)^k V(x^*, x^0)$$

$\sim \varepsilon$

малое δ уменьшает f

где δ заданное значение ε нужно

$$k = \frac{2\delta + \mu}{\mu} \log \frac{V(x^*, x^0)}{\varepsilon}$$

gr GD

$$k = \frac{L}{\mu} \log \frac{\|x^0 - x^*\|^2}{\varepsilon}$$

$$\delta \sim \frac{L}{\sqrt{N}}$$

⊕ с новой границей разрыва можем найти число итераций

⊕ с новой границей возврата вращении

• на сервере $\frac{L}{\mu} \log \frac{1}{\varepsilon}$ (так же GD)

• на сум. же-лат. $\frac{\delta}{\mu} \log \frac{1}{\varepsilon}$ (уже же GD)

⊕ вращение на сервере

NB говорим, что \mathcal{F}_1 и \mathcal{F} δ -мощны

⊙ основное внимание см. в презентации.