

最优化方法 第十四讲

Weiwen Wang(王伟文)

暨南大学

2025 年 6 月 25 日

目录

① 有限和目标函数

② 方差缩减方法

- SGD_{*}
- SGD_{*} 的收敛性分析
- 随机平均梯度 (Stochastic Average Gradient, SAG)
- Stochastic Average Gradient Amélioré: SAGA
- 随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)
- SVRG 收敛性分析

有限和目标函数

有限和目标函数

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

考慮优化问题

$$\min_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}),$$

其中 $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ 可微.

- 梯度下降法 (Gradient Descent, GD)

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma_k \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x})$$

SGD*

SGD* 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

- 小批量梯度下降法 (Min-Batch, MBGD)

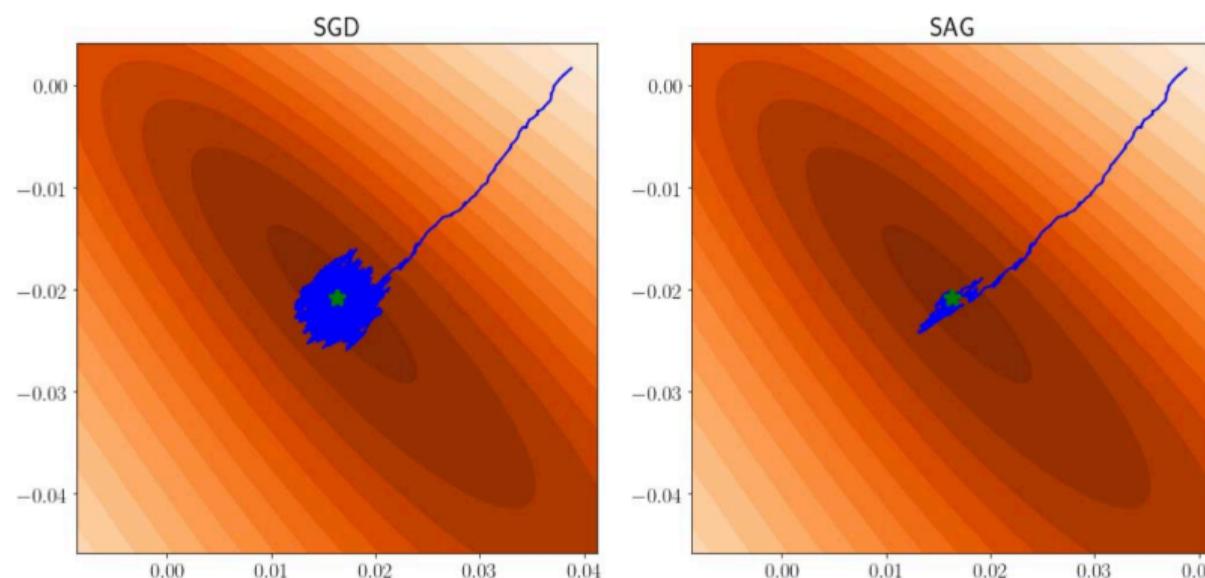
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma_k \frac{1}{|\mathcal{B}_k|} \sum_{i \in \mathcal{B}_k} \nabla f_i(\mathbf{x}) \quad \mathcal{B} \subseteq \{1, 2, \dots, n\}.$$

- 随机梯度下降法 (Stochastic Gradient Descent, SGD)

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma_k \nabla f_{i_k}(\mathbf{x}_k) \quad i_k \sim \text{Unif}(1, 2, \dots, n)$$

SGD 不收敛: $\lim_{k \rightarrow \infty} \nabla f(\mathbf{x}_k) \neq 0$

有限和目标函数



Gower, Robert M., et al. "Variance-reduced methods for machine learning."

Proceedings of the IEEE 108.11 (2020): 1968-1983.

方差缩减方法

SGD_{*}

SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

有限和目标函数

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析随机平均梯度 (Stochastic
Average Gradient, SAG)Stochastic Average Gradient
Amélioré: SAGA随机方差缩减梯度 (Stochastic
Reduced Gradient, SVRG)

SVRG 收敛性分析

引理 1

设 $X \in R^d$ 为随机向量存在有限方差, 则有

$$\mathbb{E} \left[\|X - \mathbb{E}[X]\|^2 \right] \leq \mathbb{E}[\|X\|^2]$$

有限和目标函数

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

证明.

$$\begin{aligned}\mathbb{E} [\|X - \mathbb{E} X\|^2] &= \mathbb{E} [\|X\|^2 - 2\langle X, \mathbb{E} X \rangle + \|\mathbb{E} X\|^2] \\ &= \mathbb{E}[\|X\|^2] - \|\mathbb{E} X\|^2 \leq \mathbb{E}[\|X\|^2].\end{aligned}$$

□

引理 2

设 f 为 L -smooth 凸函数, 则有

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

有限和目标函数

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

定义 3 (方差缩减性质 (Variance Reduction Property, VRP))

设随机序列 $\{\mathbf{g}_k\}_{k \geq 0} \subset \mathbb{R}^d$, 若有

$$\mathbb{E} \left[\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2 \right] \rightarrow 0 \quad \text{当 } k \rightarrow \infty$$

则称随机序列 $\{\mathbf{g}_k\}_{k \geq 0}$ 满足方差缩减性质.

有限和目标函数

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

如何构造满足方差缩减性质的随机序列 $\{\mathbf{g}_k\}_{k \geq 0}$?

- $\nabla f(\mathbf{x}_k) \approx \mathbf{g}_k$

- $\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma_k \mathbf{g}_k$

有限和目标函数

原型方法-SGD_{*}

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析随机平均梯度 (Stochastic
Average Gradient, SAG)Stochastic Average Gradient
Amélioré: SAGA随机方差缩减梯度 (Stochastic
Reduced Gradient, SVRG)

SVRG 收敛性分析

假设已知 \mathbf{x}^* 为有限和优化问题的最优解, 则定义

$$\mathbf{g}_k = \nabla f_{i_k}(\mathbf{x}_k) - \nabla f_{i_k}(\mathbf{x}^*) \quad i_k \sim \text{Unif}(1, 2, \dots, n).$$

有限和目标函数

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

 \mathbf{g}_k 是 $\nabla f(\mathbf{x}_k)$ 的条件无偏估计.

$$\begin{aligned}
 \mathbb{E}[\mathbf{g}_k | \mathbf{x}_k] &= \mathbb{E} [\nabla f_{i_k}(\mathbf{x}_k) - \nabla f_{i_k}(\mathbf{x}^*) | \mathbf{x}_k] \\
 &= \sum_{i=1}^n \frac{1}{n} (f_i(\mathbf{x}_k) - f_i(\mathbf{x}^*)) \\
 &= \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}_k) - \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}^*) \\
 &= \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*) = \nabla f(\mathbf{x}_k).
 \end{aligned}$$

有限和目标函数

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

$$\begin{aligned}\mathbb{E} \left[\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2 | \mathbf{x}_k \right] &\stackrel{(i)}{\leq} \mathbb{E} \left[\|\mathbf{g}_k\|^2 | \mathbf{x}_k \right] \\ &\leq \mathbb{E} \left[\|\nabla f_{i_k}(\mathbf{x}_k) - \nabla f_{i_k}(\mathbf{x}^*)\|^2 | \mathbf{x}_k \right]\end{aligned}$$

故

$$\mathbb{E} \left[\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2 \right] \leq \mathbb{E} \left[\|\nabla f_{i_k}(\mathbf{x}_k) - \nabla f_{i_k}(\mathbf{x}^*)\|^2 \right]$$

若 $\mathbf{x}_k \rightarrow \mathbf{x}_*$, 则 $\{\mathbf{g}_k\}$ 满足 VRP.

有限和目标函数

方差缩减方法

SGD*

SGD* 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

引理 4

设 f_i 为 L_i -smooth 的凸函数, 记 $L_{\max} = \max_{i \in [n]} L_i$. 若 $i_k \sim \text{Unif}\{1, 2, \dots, n\}$, 则

$$\mathbb{E}_{i_k} \left[\|\nabla f_{i_k}(\mathbf{x}) - \nabla f_{i_k}(\mathbf{x}^*)\|^2 \right] \leq 2L_{\max}(f(\mathbf{x}) - f(\mathbf{x}^*)) \quad \forall \mathbf{x} \in \mathbb{R}^d$$

有限目标函数

证明.

由引理 2

$$\|\nabla f_{i_k}(\mathbf{x}) - \nabla f_{i_k}(\mathbf{x}^*)\|^2 \leq 2L_{\max} \left(f_{i_k}(\mathbf{x}) - f_{i_k}(\mathbf{x}^*) - \langle \nabla f_{i_k}(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \right) \quad \forall \mathbf{x} \in \mathbb{R}^d$$

不等式两侧同时关于 i_k 求期望

$$\begin{aligned} \mathbb{E}_{i_k} \|\nabla f_{i_k}(\mathbf{x}) - \nabla f_{i_k}(\mathbf{x}^*)\|^2 &\leq \mathbb{E}_{i_k} \left[2L_{\max} \left(f_{i_k}(\mathbf{x}) - f_{i_k}(\mathbf{x}^*) - \langle \nabla f_{i_k}(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \right) \right] \\ &= 2L_{\max} \left(\frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) - \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}^*) - \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \right\rangle \right) \\ &= 2L_{\max}(f(\mathbf{x}_k) - f(\mathbf{x}^*)). \end{aligned}$$

□

有限和目标函数

方差缩减方法

SGD*

SGD* 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

定理 5

设 f_i 为 L_i -smooth 凸函数, f 为可微 μ -强凸函数, 记
 $L_{\max} = \max_{i \in [n]} L_i$, 取步长 $\gamma \in (0, 1/L_{\max})$, SGD* 迭代满足

$$\mathbb{E} \left[\|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*\|^2 \right] \leq (1 - \gamma\mu) \mathbb{E} \left[\|\boldsymbol{x}_k - \boldsymbol{x}^*\|^2 \right] \quad \forall k \geq 0$$

和 VRP.

证明.

由迭代格式

$$\begin{aligned}\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \|\mathbf{x}_k - \gamma \mathbf{g}_k - \mathbf{x}^*\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\langle \gamma \mathbf{g}_k, \mathbf{x}_k - \mathbf{x}^* \rangle + \gamma^2 \|\mathbf{g}_k\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\gamma \langle \nabla f_{i_k}(\mathbf{x}_k) - \nabla f_{i_k}(\mathbf{x}^*), \mathbf{x}_k - \mathbf{x}^* \rangle + \gamma^2 \|\nabla f_{i_k}(\mathbf{x}_k) - \nabla f_{i_k}(\mathbf{x}^*)\|^2\end{aligned}$$

不等式两端在 \mathbf{x}_k 条件下关于 i_k 求期望

$$\mathbb{E}_{i_k} \left[\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 | \mathbf{x}_k \right] = \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\gamma \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle + \gamma^2 \mathbb{E}_{i_k} \left[\|\nabla f_{i_k}(\mathbf{x}_k) - \nabla f_{i_k}(\mathbf{x}^*)\|^2 | \mathbf{x}_k \right]$$

□

有限和目标函数

方差缩减方法

SGD*

SGD* 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

证明.

f 为 μ -强凸函数

$$\langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \geq f(\mathbf{x}_k) - f(\mathbf{x}^*) + \frac{\mu}{2} \|\mathbf{x}_k - \mathbf{x}^*\|^2$$

由引理 4

$$\mathbb{E}_{i_k} \left[\|\nabla f_{i_k}(\mathbf{x}_k) - \nabla f_{i_k}(\mathbf{x}^*)\|^2 | \mathbf{x}_k \right] \leq 2L_{\max}(f(\mathbf{x}_k) - f(\mathbf{x}^*))$$

综合上述不等式

$$\mathbb{E}_{i_k} \left[\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 | \mathbf{x}_k \right] \leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 - \gamma \mu \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\gamma(1 - \gamma L_{\max})(f(\mathbf{x}_k) - f(\mathbf{x}^*))$$



证明.

注意到 $\gamma < 1/L_{\max}$ 且 $f(\mathbf{x}_k) - f(\mathbf{x}^*) \geq 0$. 故

$$\mathbb{E}_{\mathbf{x}_{k+1}} [\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 | \mathbf{x}_k] = \mathbb{E}_{i_k} [\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 | \mathbf{x}_k] \leq (1 - \mu\gamma) \|\mathbf{x}_k - \mathbf{x}^*\|^2$$

两端同时关于 \mathbf{x}_k 求期望, 由塔公式法则

$$\mathbb{E} [\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2] \leq (1 - \mu\gamma) \mathbb{E} [\|\mathbf{x}_k - \mathbf{x}^*\|^2]$$



方差缩减方法

SGD*

SGD* 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)
Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

证明.

$$\begin{aligned}\mathbb{E} \left[\|\boldsymbol{g}_k - \nabla f(\boldsymbol{x}_k)\|^2 | \boldsymbol{x}_k \right] &\leq \mathbb{E} \left[\|\nabla f_{i_k}(\boldsymbol{x}_k) - \nabla f_{i_k}(\boldsymbol{x}^*)\|^2 | \boldsymbol{x}_k \right] \\ &\leq \mathbb{E} \left[L_i^2 \|\boldsymbol{x}_k - \boldsymbol{x}^*\|^2 | \boldsymbol{x}_k \right] \\ &\leq L_{\max}^2 \|\boldsymbol{x}_k - \boldsymbol{x}^*\|^2\end{aligned}$$

故

$$\begin{aligned}\mathbb{E} \left[\|\boldsymbol{g}_k - \nabla f(\boldsymbol{x}_k)\|^2 \right] &\leq L_{\max}^2 \mathbb{E} \left[\|\boldsymbol{x}_k - \boldsymbol{x}^*\|^2 \right] \\ &\leq L_{\max}^2 (1 - \mu\gamma)^k \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2\end{aligned}$$

□

有限目标函数

定义辅助变量 $\mathbf{v}_k^{(i)}$

$$\mathbf{v}_k^{(j)} = \begin{cases} \nabla f_j(\mathbf{x}_k), & j = i_k \\ \mathbf{v}_{k-1}^{(i)}, & j \neq i_k \end{cases} \quad i_k \sim \text{Unif}(1, 2, \dots, n)$$

初始化为 $\mathbf{v}_0^{(i)} = \mathbf{0}, \forall i \in [n]$.

$$\begin{aligned} \mathbf{g}_k &= \frac{1}{n} \sum_{i=1}^n \mathbf{v}_k^{(i)} = \frac{1}{n} \sum_{j:j \neq i_k} \mathbf{v}_k^{(j)} + \frac{1}{n} \mathbf{v}_k^{(i_k)} = \frac{1}{n} \sum_{j:j \neq i_k} \mathbf{v}_{k-1}^{(j)} + \frac{1}{n} \mathbf{v}_k^{(i_k)} \\ &= \mathbf{g}_{k-1} - \frac{1}{n} \mathbf{v}_{k-1}^{(i_k)} + \frac{1}{n} \nabla f_{i_k}(\mathbf{x}_k) \end{aligned}$$

方差缩减方法

SGD*

SGD* 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

有限和目标函数

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析随机平均梯度 (Stochastic
Average Gradient, SAG)Stochastic Average Gradient
Amélioré: SAGA随机方差缩减梯度 (Stochastic
Reduced Gradient, SVRG)

SVRG 收敛性分析

\mathbf{g}_k 不是 $\nabla f(\mathbf{x}_k)$ 的一个无偏估计

$$\begin{aligned}\mathbb{E}[\mathbf{g}_k | \mathbf{x}_k] &= \mathbb{E} \left[\mathbf{g}_{k-1} - \frac{1}{n} \mathbf{v}_{k-1}^{(i_k)} + \frac{1}{n} \nabla f_{i_k}(\mathbf{x}_k) | \mathbf{x}_k \right] \\ &= \nabla f(\mathbf{x}_k) + \mathbf{g}_{k-1} - \frac{1}{n} \mathbf{g}_{k-1}\end{aligned}$$

有限目标函数

给定一组向量 $\{\mathbf{v}_i\}_{i \in [n]}$,

$$\nabla f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \left(\nabla f_i(\mathbf{x}) - \mathbf{v}_i + \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \right)$$

记

$$\begin{aligned}\nabla f_i(\mathbf{x}; \mathbf{v}_i) &= \nabla f_i(\mathbf{x}) - \mathbf{v}_i + \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \\ &= \nabla f_i(\mathbf{x}) - \mathbf{v}_i + \mathbb{E}[\mathbf{v}_i]\end{aligned}$$

定义

$$\mathbf{g}_k = \nabla f_{i_k}(\mathbf{x}_k, \mathbf{v}_{i_k}) = \nabla f_{i_k}(\mathbf{x}) - \mathbf{v}_{i_k} + \mathbb{E}[\mathbf{v}_{i_k}] \quad i_k \sim \text{Unif}(1, 2, \dots, n)$$

方差缩减方法

SGD*

SGD* 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

有限和目标函数

方差缩减性质

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

$$\begin{aligned}\mathbb{E}_{i_k} [\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2] &= \mathbb{E}_{i_k} \left[\|\nabla f_{i_k}(\mathbf{x}) - \mathbf{v}_{i_k} + \mathbb{E}[\mathbf{v}_{i_k}] - \nabla f(\mathbf{x}_k)\|^2 \right] \\ &\leq \mathbb{E}_{i_k} \left[\|\nabla f_{i_k}(\mathbf{x}_k) - \mathbf{v}_{i_k}\|^2 \right]\end{aligned}$$

若 $\nabla f_{i_k}(\mathbf{x}_k) \xrightarrow[k \rightarrow \infty]{} \mathbf{v}_{i_k}$, $\mathbb{E}_{i_k} [\|\mathbf{g}_k - \nabla f(\mathbf{x}_k)\|^2] \rightarrow 0$.

有限和目标函数

近似梯度

定义

$$\mathbf{v}_k^{(j)} = \nabla f_j(\mathbf{x}_{\setminus k})$$

其中 $\nabla f_j(\mathbf{x}_{\setminus k})$ 表示第 k 次迭代前最近一次计算的 f_j 的梯度.

迭代格式的近似梯度定义为

$$\mathbf{g}_k = \nabla f_{i_k}(\mathbf{x}_k) - \mathbf{v}_k^{(i_k)} + \frac{1}{n} \sum_{i=1}^n \mathbf{v}_k^{(i)} \quad i_k \sim \text{Unif}(1, 2, \dots, n)$$

有限和目标函数

方差缩减方法

SGD_{*}SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

高效计算

记 $\bar{\mathbf{g}}_{k-1} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_k^{(i)}$

- $i_k \in \text{Unif}(1, 2, \dots, n)$
- $\mathbf{g}_k = \nabla f_{i_k}(\mathbf{x}_k) - \mathbf{v}_k^{(i_k)} + \bar{\mathbf{g}}_{k-1}$
- $\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma_k \mathbf{g}_k$
- $\bar{\mathbf{g}}_k = \bar{\mathbf{g}}_{k-1} - \frac{1}{n} \mathbf{v}_k^{(i_k)} + \frac{1}{n} \nabla f_{i_k}(\mathbf{x}_k)$

初始化取 $\mathbf{v}_0^{(i)} = \mathbf{0}, \forall i \in [n], \bar{\mathbf{g}}_0 = \mathbf{0}$.

有限和目标函数

方差缩减方法

SGD_{*}

SGD_{*} 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

SAG 和 SAGA 需要维持向量组 $\{\mathbf{v}_i\}$, 空间复杂度 $\mathcal{O}(nd)$.

SVRG 定义近似梯度

$$\mathbf{g}_k = \nabla f_{i_k}(\mathbf{x}_k) - \nabla f_{i_k}(\bar{\mathbf{x}}_k) + \nabla f(\bar{\mathbf{x}}_k)$$

空间复杂度 $\mathcal{O}(d)$.

有限目标函数

```
1 import numpy as np
2
3 def stochastic_variance_reduced_gradient(theta, step_size, sample_num,
4     outer_max, inner_max):
5     for o in range(outer_max):
6         # compute full gradient of the finite sum objective function at theta
7         full_grad = gradient_of_obj(theta)
8         inner_theta = theta
9         for i in range(inner_max):
10            idx = np.random.randint(0, sample_num)
11            # compute individual gradient with index idx at inner_theta
12            grad1 = gradient_of_obj(inner_theta, idx)
13            # compute individual gradient with index idx at theta
14            grad2 = gradient_of_obj(theta, idx)
15            # gradient estimation
16            grad = grad1 - grad2 + full_grad
17            inner_theta += -step_size * grad
18
19            theta = inner_theta
20
21
22 return theta
```

方差缩减方法

SGD*

SGD* 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析

有限目标函数

定理 6

设 f_i 为 L_i -smooth 凸函数, f 为可微 μ -强凸函数, 记

$L_{\max} = \max_{i \in [n]} L_i$, 取步长 $\gamma \in (0, \frac{1}{2L_{\max}})$, 若 SVRG 内循环迭代次数 m 充分大使得

$$\rho \triangleq \frac{1}{\mu\gamma(1-2L_{\max}\gamma)m} + \frac{2L_{\max}\gamma}{1-2L_{\max}\gamma} < 1$$

则

$$\mathbb{E}[f(\mathbf{x}_o) - f(\mathbf{x}^*)] \leq \rho^o(f(\mathbf{x}_0) - f(\mathbf{x}^*)).$$

这里 $\mathbf{x}_o = \frac{1}{m} \sum_{i=0}^{m-1} \mathbf{x}_{o_i}$, 其中 \mathbf{x}_{o_i} 表示第 o 次外循环中的第 i 次内循环更新结果.

方差缩减方法

SGD*

SGD* 的收敛性分析

随机平均梯度 (Stochastic Average Gradient, SAG)

Stochastic Average Gradient
Amélioré: SAGA

随机方差缩减梯度 (Stochastic Reduced Gradient, SVRG)

SVRG 收敛性分析