

A Stochastic Composite Gradient Method with Incremental Variance Reduction

Junyu Zhang (University of Minnesota)
and
Lin Xiao (Microsoft Research)

Neural Information Processing Systems (NeurIPS)

Vancouver, Canada
December 8-14 , 2019

1

Composite stochastic optimization

- composition with expectation

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f(\mathbf{E}_g[g_\zeta(x)]) + r(x)$$

- $f : \mathbb{R}^p \rightarrow \mathbb{R}$ smooth and possibly nonconvex
- $g_\zeta : \mathbb{R}^d \rightarrow \mathbb{R}^p$ smooth vector mapping for every ζ
- $r : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ convex but possibly nonsmooth

- composition with finite sum

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f\left(\frac{1}{n} \sum_{i=1}^n g_i(x)\right) + r(x)$$

- applications beyond ERM

- policy evaluation in reinforcement learning
- risk-averse optimization, financial mathematics
- ...

2

Examples

- policy evaluation with linear function approximation

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad \|\mathbf{E}[A]x - \mathbf{E}[b]\|^2$$

A, b random, generated by MDP under fixed policy

- risk-averse optimization

$$\underset{x \in \mathbb{R}^d}{\text{maximize}} \quad \underbrace{\frac{1}{n} \sum_{i=1}^n h_i(x)}_{\text{average reward}} - \lambda \underbrace{\frac{1}{n} \sum_{i=1}^n (h_i(x) - \frac{1}{n} \sum_{i=1}^n h_i(x))^2}_{\text{variance of rewards (risk)}}$$

- often treated as two-level composite finite-sum optimization

simple transformation using $\text{Var}(a) = \mathbf{E}[a^2] - \|\mathbf{E}[a]\|^2$

$$\underset{x \in \mathbb{R}^d}{\text{maximize}} \quad \frac{1}{n} \sum_{j=1}^n h_j(x) - \lambda \left(\frac{1}{n} \sum_{j=1}^n h_j^2(x) - \left(\frac{1}{n} \sum_{j=1}^n h_j(x) \right)^2 \right)$$

actually a one-level composite finite-sum problem

3

Technical challenge and related work

- challenge: biased gradient estimator

denote $F(x) := f(g(x))$ where $g(x) := \mathbf{E}_g[g_\zeta(x)]$

$$F'(x) = [g'(x)]^T f'(g(x))$$

- subsampled estimators

$$y = \frac{1}{|S|} \sum_{\zeta \in S} g_\zeta(x), \quad z = \frac{1}{|S|} \sum_{\zeta \in S} g'_\zeta(x)$$

$$\mathbf{E}[y] = g(x) \text{ and } \mathbf{E}[z] = g'(x), \text{ but } \mathbf{E}[z^T f'(y)] \neq F'(x)$$

- related work

more general composite stochastic optimization
(Wang, Fang & Liu 2017; Wang, Liu & Fang 2017; ...)

two-level composite finite-sum: extending SVRG and SAGA
(Lian, Wang & Liu 2017; Huo, Gu, & Huang 2018; Lin, Fan, Wang & Jordan 2018; Zhang & Xiao 2019, ...)

4

Convergence analysis

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad \frac{f(\mathbf{E}_g[g_\zeta(x)]) + r(x)}{F(x)}$$

- assumptions

f is L_f -Lipschitz and f' is L_f -Lipschitz

r convex but can be non-smooth

g_ζ and g'_ζ are mean-square Lipschitz with constants ℓ_g and L_g

$$\mathbf{E}\|g_\zeta(x) - g_\zeta(y)\|^2 \leq \ell_g^2 \|x - y\|^2$$

$$\mathbf{E}\|g'_\zeta(x) - g'_\zeta(y)\|^2 \leq L_g^2 \|x - y\|^2$$

g_ζ and g'_ζ have bounded variance

$$\mathbf{E}\|g_\zeta(x) - g(y)\|^2 \leq \sigma_g^2, \quad \mathbf{E}\|g'_\zeta(x) - g'(y)\|^2 \leq \sigma_{g'}^2$$

- sample complexity for $\mathbf{E}[\|\mathcal{G}(x)\|^2] \leq \epsilon$, where

$$G(x) = \frac{1}{\eta} (x - \text{prox}_r^y(x - \eta F'(x))) = F'(x) \text{ if } r \equiv 0$$

5

Main results

$$\underset{x}{\text{minimize}} \quad \Psi(x) + r(x) \triangleq f(\mathbf{E}_g[g_\zeta(x)]) + r(x)$$

idea: use SARAH/SPIDER estimator for both $g(x)$ and $g'(x)$

sample complexities

assumptions (common: F and g_ζ Lipschitz and smooth, thus F smooth)

F nonconvex	F μ -gradient dominant	F convex, r convex	Φ optimally strongly convex
r convex	$r = 0$	Φ	μ -optimally strongly convex
$\mathbf{E}[\mathcal{O}(\epsilon^{-3/2})]$	$\mathcal{O}((\nu\epsilon^{-1}) \log \epsilon^{-1})$	$\mathcal{O}((\mu^{-1}\epsilon^{-1}) \log \epsilon^{-1})$	$\mathcal{O}((\mu^{-1}\epsilon^{-1}) \log \epsilon^{-1})$
$\sum \mathcal{O}(\min\{\epsilon^{-3/2}, \mu^2/\nu\epsilon^{-1}\})$	$\mathcal{O}((n + \nu\epsilon^{-1}) \log \epsilon^{-1})$	$\mathcal{O}((n + \mu^{-1}\epsilon^{-1}) \log \epsilon^{-1})$	$\mathcal{O}((n + \mu^{-1}\epsilon^{-1}) \log \epsilon^{-1})$

same as best complexities for problems without composition

lower bound $\mathcal{O}(\min\{\epsilon^{-3/2}, \mu^2/\nu\epsilon^{-1}\})$ (Fang, Li, Lin & Zhang 2018)

composition (biased estimator) does not incur higher complexity

6

Composite Incremental Variance Reduction (CIVR)

input: $x_0^1, \eta > 0, T \geq 1, \{\tau_t, B_t, S_t\}$ for $t = 1, \dots, T$

for $t = 0, \dots, T-1$

- sample set B_t with size B_t and construct the estimates

$$y_t^t = \frac{1}{B_t} \sum_{\zeta \in B_t} g_\zeta(x_0^t), \quad z_t^t = \frac{1}{B_t} \sum_{\zeta \in B_t} g'_\zeta(x_0^t)$$

- compute $\tilde{V}F(x_0^t) = (z_t^t)^T F'(y_t^t)$ and let $x_t^t = \text{prox}_r^y(x_0^t - \eta \tilde{V}F(x_0^t))$

- for $i = 0, \dots, t-1$
 - sample a set S_i^t with size S_i^t and construct the estimates

$$y_i^t = y_{i-1}^t + \frac{1}{S_i^t} \sum_{\zeta \in S_i^t} (g_\zeta(x_i^t) - g_\zeta(x_{i-1}^t))$$

$$z_i^t = z_{i-1}^t + \frac{1}{S_i^t} \sum_{\zeta \in S_i^t} (g'_\zeta(x_i^t) - g'_\zeta(x_{i-1}^t))$$

- compute $\tilde{V}F(x_i^t) = (z_i^t)^T F'(y_i^t)$ and $x_{i+1}^t = \text{prox}_r^y(x_i^t - \eta \tilde{V}F(x_i^t))$

- set $x_0^{t+1} = x_t^t$.

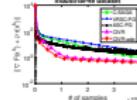
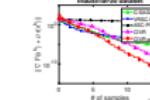
output: \bar{x} randomly chosen from $\{x_i^t\}_{i=0, \dots, t-1}^{t=1, \dots, T}$

6

Experiments on risk-averse optimization

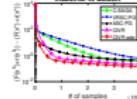
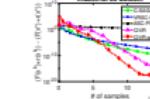
- reduction of gradient norm

Industrial-49 dataset



- reduction of objective value

Industrial-49 dataset



9