

Optimization Problem

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{j=1}^n f_j(x) + \psi(x),$$

- $f_j : \mathbb{R}^d \rightarrow \mathbb{R}$ is M_j smooth and convex:

$$0 \preceq \nabla^2 f_j(x) \preceq M_j$$

- $\psi : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is a proper, closed and convex regularizer, admitting a cheap proximal operator
- $f \stackrel{\text{def}}{=} \frac{1}{n} \sum_j f_j$ is σ quasi strongly convex

Oracle

$G(x) \stackrel{\text{def}}{=} [\nabla f_1(x), \nabla f_2(x), \dots, \nabla f_n(x)]$: Jacobian matrix

- Oracle can be accessed via: $\mathcal{U}G(x)$, $\mathcal{S}G(x)$
- $\mathcal{U} : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ - random linear operator, identity in expectation
- $\mathcal{S} : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ - random projection operator, possibly correlated with \mathcal{U}
- \mathcal{U} , \mathcal{S} might correspond to **right matrix multiplication** (SAGA [1], JacSketch [2]), **left matrix multiplication** (SEGA [3]), their **combination** (ISAEGA) and **many more**
- Different choices of \mathcal{U} , \mathcal{S} yield different methods.

Variance reduction (unbiased)

Given sequence J^k which estimates $G(x^k)$ such that $\lim_{k \rightarrow \infty} J^k = G(x^*)$, unbiased variance reduced gradient is the following:

$$g^k = \frac{1}{n} J^k e + \frac{1}{n} \mathcal{U} (G(x^k) - J^k) e. \quad (1)$$

Jacobian Sketching

Observing $\mathcal{S}G(x^k)$ every iteration, how to design Jacobian estimator sequence J^k ? **Projecting**:

$$J^{k+1} = \underset{J}{\text{argmin}} \|J - J^k\| \quad \text{s. t. } \mathcal{S}J = \mathcal{S}G(x^k) \\ = J^k - \mathcal{S}(G(x^k) - J^k) \quad (2)$$

Algorithm

Algorithm 1 Generalized JacSketch (GJS)

- Parameters:** Step size $\alpha > 0$, random projector \mathcal{S} and unbiased sketch \mathcal{U}
- Initialization:** Choose solution estimate $x^0 \in \mathbb{R}^d$ and Jacobian estimate $J^0 \in \mathbb{R}^{d \times n}$
- for** $k = 0, 1, \dots$ **do**
- Sample realizations of \mathcal{S} and \mathcal{U} , and perform sketches $\mathcal{S}G(x^k)$ and $\mathcal{U}G(x^k)$
- $J^{k+1} = J^k - \mathcal{S}(J^k - G(x^k))$ update the Jacobian estimate via (2)
- $g^k = \frac{1}{n} J^k e + \frac{1}{n} \mathcal{U} (G(x^k) - J^k) e$ construct the gradient estimator via (1)
- $x^{k+1} = \text{prox}_{\alpha\psi}(x^k - \alpha g^k)$ perform the proximal SGD step
- end for**

Convergence rate

Single convergence theorem, **tightest known rate in every special case** (many new rates in special cases for known methods; many new methods as well).

- Let $\mathcal{M} : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ be linear operator such that $(\mathcal{M}X)_{:,j} = M_j X_{:,j}$ for any $X \in \mathbb{R}^{d \times n}$.
- Let $\mathcal{B} : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ be a linear operator (to be chosen; only for theory) such that with step size α we have:

$$(1 - \alpha\sigma) \|\mathcal{B}\mathcal{M}^{\frac{1}{2}}X\|^2 \geq \frac{2\alpha}{n^2} \mathbb{E} [\|\mathcal{U}X e\|^2] \\ + \left\| (\mathcal{I} - \mathbb{E}[\mathcal{S}])^{\frac{1}{2}} \mathcal{B}\mathcal{M}^{\frac{1}{2}}X \right\|^2 \\ \frac{1}{n} \|\mathcal{M}^{\frac{1}{2}}X\|^2 \geq \frac{2\alpha}{n^2} \mathbb{E} [\|\mathcal{U}X e\|^2] + \left\| (\mathbb{E}[\mathcal{S}])^{\frac{1}{2}} \mathcal{B}\mathcal{M}^{\frac{1}{2}}X \right\|^2$$

Theorem (simplified)

For GJS we have $\mathbb{E}[\Psi^k] \leq (1 - \alpha\sigma)^k \Psi^0$ for $\Psi^k \stackrel{\text{def}}{=} \|x^k - x^*\|^2 + \alpha \|\mathcal{B}(J^k - G(x^*))\|^2$.

- Linear convergence under minimal assumptions
- Rate depends on smoothness patterns (matrices M_j), distributions of \mathcal{S}, \mathcal{U} (controllable in practice) and quasi strong convexity σ
- Full version:** exploits possible prior knowledge about $G(x^*)$, exploits structure of ψ , extends quasi strong convexity to strong growth.

Special cases

- SAGA [1]: recovers best known results
- JacSketch [2] More general + better rate
- LSVRG: Arbitrary sampling + prox
- SEGA [3]: Better rate under arbitrary sampling
- Extensions of algorithms from [4] – arbitrary sampling and conjectured ISEAGA.
- Many more:

Choice of random operators \mathcal{S} and \mathcal{U} defining Algorithm 1	$\mathcal{U}X$	#	Name	Algorithm Comment	Section
$\mathcal{S}X$	$\mathcal{U}X$	2	SAGA	basic variant of SAGA [3]	G.1
$\mathcal{X} \sum_{j \in R} e_j e_j^\top$ w.p. p_R	$\mathcal{X} \sum_{j \in R} \frac{1}{p_j} e_j e_j^\top$ w.p. p_R	3	SAGA	SAGA with AS [26]	G.2
$\sum_{i \in L} e_i e_i^\top$ w.p. p_i	$\sum_{i \in L} \frac{1}{p_i} e_i e_i^\top$ w.p. p_i	4	SEGA	basic variant of SEGA [9]	H.1
$\sum_{i \in L} e_i e_i^\top$ w.p. p_i	$\sum_{i \in L} \frac{1}{p_i} e_i e_i^\top$ w.p. p_i	5	SEGA	SEGA [9] with AS and prox	H.2
$\begin{cases} 0 & \text{w.p. } 1 - \rho \\ \mathcal{X} & \text{w.p. } \rho \end{cases}$	$\sum_{i \in L} \frac{1}{p_i} e_i e_i^\top$ w.p. p_i	6	SVRG	NEW	H.3
$\begin{cases} 0 & \text{w.p. } 1 - \rho \\ \mathcal{X} & \text{w.p. } \rho \end{cases}$	$\mathcal{X} \sum_{j \in R} \frac{1}{p_j} e_j e_j^\top$ w.p. p_R	7	SGD-ablft	NEW	I
$\begin{cases} 0 & \text{w.p. } 1 - \rho \\ \mathcal{X} & \text{w.p. } \rho \end{cases}$	$\mathcal{X} \sum_{j \in R} \frac{1}{p_j} e_j e_j^\top$ w.p. p_R	8	LSVRG	LSVRG [14] with AS and prox	J
$\begin{cases} 0 & \text{w.p. } 1 - \rho \\ \mathcal{X} & \text{w.p. } \rho \end{cases}$	$\begin{cases} 0 & \text{w.p. } 1 - \delta \\ \frac{1}{2} \mathcal{X} & \text{w.p. } \delta \end{cases}$	9	B2	NEW	K.1
$\mathcal{X} \sum_{j \in R} e_j e_j^\top$ w.p. p_R	$\begin{cases} 0 & \text{w.p. } 1 - \delta \\ \frac{1}{2} \mathcal{X} & \text{w.p. } \delta \end{cases}$	10	LSVRG-inv	NEW	K.2
$\sum_{i \in L} e_i e_i^\top$ w.p. p_i	$\begin{cases} 0 & \text{w.p. } 1 - \delta \\ \frac{1}{2} \mathcal{X} & \text{w.p. } \delta \end{cases}$	11	SVRG-inv	NEW	K.3
$\mathcal{X} \sum_{j \in R} e_j e_j^\top$ w.p. p_R	$\sum_{i \in L} \frac{1}{p_i} e_i e_i^\top$ w.p. p_i	12	RL	NEW	L.1
$\sum_{i \in L} e_i e_i^\top$ w.p. p_i	$\mathcal{X} \sum_{j \in R} \frac{1}{p_j} e_j e_j^\top$ w.p. p_R	13	LR	NEW	L.2
$\mathcal{I}_L \mathcal{X} \mathcal{I}_R$ w.p. $p_L p_R$	$\mathcal{I}_L \left(\left(p^{-1} (p^{-1})^\top \right) \circ \mathcal{X} \right) \mathcal{I}_R$ w.p. $p_L p_R$	14	SABGA	NEW	M.1
$\begin{cases} 0 & \text{w.p. } 1 - \rho \\ \mathcal{X} & \text{w.p. } \rho \end{cases}$	$\mathcal{I}_L \left(\left(p^{-1} (p^{-1})^\top \right) \circ \mathcal{X} \right) \mathcal{I}_R$ w.p. $p_L p_R$	15	SVRG	NEW	M.2
$\sum_{i \in L} \mathcal{I}_{L_i} \mathcal{X}_{N_i} \mathcal{I}_{R_i}$	$\sum_{i \in L} \left((p^i)^{-1} (p^i)^{-1 \top} \right) \circ (\mathcal{I}_{L_i} \mathcal{X}_{N_i} \mathcal{I}_{R_i})$	16	ISABGA	NEW (reminiscent of [20])	M.3
$\sum_{i \in L} \mathcal{I}_{L_i} \mathcal{X}_{N_i}$	$\sum_{i \in L} \left((p^i)^{-1} e^{\top} \right) \circ (\mathcal{I}_{L_i} \mathcal{X}_{N_i})$	17	ISEGA	ISEGA [20] with AS	M.3
$\mathcal{X}R$	$\mathcal{X}R \mathbb{E}[R]^{-1}$	18	JS	JacSketch [8] with AS and prox.	N

Arbitrary sampling

- Tight rate under any distribution of \mathcal{S}, \mathcal{U}
- Allows to exploit data structure from smoothness (matrices M_j) and design importance samplings
- New for many well established algorithms, bridged by our analysis

Specific algorithms

Algorithm 2 SEGA with arbitrary sampling

- Require:** Step size $\alpha > 0$, starting point $x^0 \in \mathbb{R}^d$, random sampling $L \subseteq \{1, 2, \dots, d\}$
- Set $h^0 = 0$
- for** $k = 0, 1, 2, \dots$ **do**
- Sample random $L^k \subseteq \{1, 2, \dots, d\}$
- Set $h^{k+1} = h^k + \sum_{i \in L^k} (\nabla_i f(x^k) - h_i^k) e_i$
- $g^k = h^k + \sum_{i \in L^k} \frac{1}{p_i} (\nabla_i f(x^k) - h_i^k) e_i$
- $x^{k+1} = \text{prox}_{\alpha\psi}(x^k - \alpha g^k)$
- end for**

Algorithm 3 ISAEGA [NEW METHOD]

- Input:** $x^0 \in \mathbb{R}^d$, # parallel units T , each owning set of indices N_t (for $1 \leq t \leq T$), distributions \mathcal{D}_t over subsets of N_t , distributions \mathcal{D}_t over subsets coordinates $[d]$, step size α
- $J^0 = 0$
- for** $k = 0, 1, \dots$ **do**
- for** $t = 1, \dots, T$ **in parallel do**
- Sample $R_t \sim \mathcal{D}_t$, $R_t \subseteq N_t$, $L_t \sim \mathcal{D}_t$; $L_t \subseteq [d]$
- Observe $\nabla_{L_t} f_j(x^k)$ for $j \in R_t$
- Set $J_{i,j}^{k+1} = \begin{cases} \nabla_i f_j(x^k) & \text{if } j \in R_t, i \in L_t \\ J_{i,j}^k & \text{otherwise} \end{cases}$
- Send $J_{N_t}^{k+1} - J_{N_t}^k$ to master ▷ Sparse
- end for**
- $g^k = \left(J^k + \sum_{t=1}^T \left(p^{t-1} p^{-1 \top} \right) \circ \left(\sum_{i \in L_t} e_i e_i^\top \right) (J^{k+1} - J^k)_{:,N_t} \left(\sum_{j \in R_t} e_j e_j^\top \right) \right) e$
- $x^{k+1} = \text{prox}_{\alpha\psi}(x^k - \alpha g^k)$
- end for**

References

- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2014.
- Robert M Gower, Peter Richtárik, and Francis Bach. Stochastic quasi-gradient methods: Variance reduction via Jacobian sketching. *arXiv preprint arXiv:1805.02632*, 2018.
- Filip Hanzely, Konstantin Mishchenko, and Peter Richtárik. SEGA: Variance reduction via gradient sketching. In *Advances in Neural Information Processing Systems*, pages 2082–2093, 2018.
- Konstantin Mishchenko, Filip Hanzely, and Peter Richtárik. 99% of distributed optimization is a waste of time: The issue and how to fix it. *arXiv preprint arXiv:1901.09437*, 2019.