

# Line-search second-order methods for optimization in noisy environments

**Marco Viola**

Department of Mathematics and Physics  
University of Campania “L. Vanvitelli”  
`marco.viola@unicampania.it`

Joint work with

[Daniela di Serafino](#) – University of Naples Federico II  
[Nataša Krejić](#), [Nataša Krklec Jerinkić](#) – University of Novi Sad

**BOS/SOR2020 Conference**  
Palais Staszic, Warsaw  
December 15, 2020



●  
●  
Università  
degli Studi  
della Campania  
*Luigi Vanvitelli*

# Outline

- 1 Problem, motivations and contribution
- 2 The LSOS framework
- 3 Numerical experiments with LSOS
- 4 Specializing LSOS for finite sums
- 5 Numerical experiments with LSOS-BFGS
- 6 Conclusions and future work

# Outline

- 1 Problem, motivations and contribution
- 2 The LSOS framework
- 3 Numerical experiments with LSOS
- 4 Specializing LSOS for finite sums
- 5 Numerical experiments with LSOS-BFGS
- 6 Conclusions and future work

# The problem

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \phi(\mathbf{x})$$

$\phi(\mathbf{x})$  twice continuously differentiable function in a **noisy environment**, i.e.  $\phi(\mathbf{x})$ ,  $\nabla\phi(\mathbf{x})$  and  $\nabla^2\phi(\mathbf{x})$  are only accessible with some level of noise:

$$f(\mathbf{x}) = \phi(\mathbf{x}) + \varepsilon_f(\mathbf{x})$$

$$\mathbf{g}(\mathbf{x}) = \nabla\phi(\mathbf{x}) + \varepsilon_g(\mathbf{x})$$

$$B(\mathbf{x}) = \nabla^2\phi(\mathbf{x}) + \varepsilon_B(\mathbf{x})$$

$\varepsilon_f(\mathbf{x})$  random number,  $\varepsilon_g(\mathbf{x})$  random vector,  $\varepsilon_B(\mathbf{x})$  symmetric random matrix

# The problem (cont'd)

The error may derive from:

- uncertainty on data;
- measurement errors;
- communication errors;
- computational inaccuracy (data come from a simulation);
- ...

# The problem (cont'd)

The error may derive from:

- uncertainty on data;
- measurement errors;
- communication errors;
- computational inaccuracy (data come from a simulation);
- ...

Special cases:

- mathematical expectation:

$$\phi(\mathbf{x}) = E_{\xi \sim \mathcal{D}} [v(\mathbf{x}, \xi)], \quad \text{and} \quad f(\mathbf{x}) = v(\mathbf{x}, \bar{\xi}), \quad \text{with} \quad \bar{\xi} \sim \mathcal{D}$$

# The problem (cont'd)

The error may derive from:

- uncertainty on data;
- measurement errors;
- communication errors;
- computational inaccuracy (data come from a simulation);
- ...

Special cases:

- mathematical expectation:

$$\phi(\mathbf{x}) = E_{\xi \sim \mathcal{D}} [v(\mathbf{x}, \xi)], \quad \text{and} \quad f(\mathbf{x}) = v(\mathbf{x}, \bar{\xi}), \quad \text{with} \quad \bar{\xi} \sim \mathcal{D}$$

- (large) finite sum of functions:

$$\phi(\mathbf{x}) = \sum_{i=1}^N \phi_i(\mathbf{x}), \quad \text{and} \quad f(\mathbf{x}) = \sum_{i \in \mathcal{S}} \phi_i(\mathbf{x}), \quad \text{with} \quad \mathcal{S} \subseteq \{1, \dots, N\}$$

# Stochastic optimization methods

## First-order methods (NON-exhaustive list)

- Stochastic Approximation - SA (Stochastic Gradient - SG)  
[Robbins & Monro, Ann. Math. Statistics 1951] (convergence in probability with harmonic-type step length, also almost sure (a.s.) convergence with SA variants)
- In the “realm” of machine learning:
  - ▶ minibatch gradient methods, see e.g. [Bottou, Curtis & Nocedal, SIREV 2018] (convergence in expectation of obj fun error with constant or harmonic-type step length)
  - ▶ variance-reduction gradient methods, e.g. SVRG [Johnson & Zhang, NIPS 2013], SAGA [Defazio, Bach & Lacoste-Julien, NIPS 2014], JacSketch [Gower, Richtárik & Bach, Math Prog 2020] (linear convergence in expectation with constant step length)



# Stochastic optimization methods (cont'd)

## Methods using second-order info (NON-exhaustive list)

- Stochastic versions of Newton-type methods
  - ▶ Ruppert, Ann Statist 1985
  - ▶ Spall, Proc various IEEE Conferences 1994, 1995, 1005
  - ▶ Byrd, Chin, Neveitt & Nocedal, SIOPT 2011
  - ▶ Byrd, Chin, Nocedal & Wu, Math Program 2012
  - ▶ Bellavia, Krejić & Krklec Jerinkić, IMA JNA 2019
  - ▶ Bollapragada, Byrd & Nocedal, IMA JNA 2019
- Stochastic BFGS
  - ▶ Byrd, Chin, Neveitt & Nocedal, SIOPT 2011
  - ▶ Moktari & Ribeiro, IEEE TSP 2014
  - ▶ Byrd, Hansen, Nocedal & Singer, SIOPT 2016
  - ▶ Gower, Goldfarb & Richtárik, Proc ICML 2016
  - ▶ Moritz, Nishihara & Jordan, Proc MLR 2016

# Our family of methods: LSOS

- **L**ine-search **S**econd-**O**rders **S**tochastic algorithmic framework, where Newton-type and quasi-Newton directions are used
- Almost sure convergence of the sequence of iterates generated by the methods fitting into the LSOS framework and effectiveness in practice
- For finite-sum objective functions (e.g. in machine learning)
  - ▶ stochastic L-BFGS for Hessian estimates + SAGA-type for gradient estimates + line search
  - ▶ almost sure convergence of the sequence of iterates (for state-of-the-art stochastic L-BFGS convergence in expectation of the obj function error)
  - ▶ linear convergence rate and worst-case  $\mathcal{O}(\log(\varepsilon^{-1}))$  complexity
  - ▶ practical efficiency (comparison with state-of-the-art stochastic optimization methods)

# Outline

- 1 Problem, motivations and contribution
- 2 The LSOS framework
- 3 Numerical experiments with LSOS
- 4 Specializing LSOS for finite sums
- 5 Numerical experiments with LSOS-BFGS
- 6 Conclusions and future work

# SOS: Second-Order Stochastic method

---

## Sketch of SOS method

---

- 1: given  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $\{\alpha_k\} \subset \mathbb{R}_+$
  - 2: **for**  $k = 0, 1, 2, \dots$  **do**
  - 3:   compute  $\mathbf{d}_k \in \mathbb{R}^n$
  - 4:   set  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$
  - 5: **end for**
- 

$\mathbf{d}_k$  specified later

# SOS: Second-Order Stochastic method

---

## Sketch of SOS method

---

- 1: given  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $\{\alpha_k\} \subset \mathbb{R}_+$
  - 2: **for**  $k = 0, 1, 2, \dots$  **do**
  - 3:   compute  $\mathbf{d}_k \in \mathbb{R}^n$
  - 4:   set  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$
  - 5: **end for**
- 

$\mathbf{d}_k$  specified later

## Basic assumptions

- 1  $\phi$  strongly convex with Lipschitz-continuous gradient:
  - ▶  $\mathbf{x}_*$  unique solution
  - ▶  $\mu I \preceq \nabla^2 \phi(\mathbf{x}) \preceq LI$
- 2 Harmonic step-length sequence:  
 $\alpha_k > 0, \sum_k \alpha_k = \infty, \sum_k \alpha_k^2 < \infty$
- 3 Unbiased gradient estimator and bounded variance of gradient errors:  
 $\mathbb{E}(\boldsymbol{\varepsilon}_g(\mathbf{x}) | \mathcal{F}_k) = 0$  and  $\mathbb{E}(\|\boldsymbol{\varepsilon}_g(\mathbf{x})\|^2 | \mathcal{F}_k) \leq M$   
( $\mathcal{F}_k = \sigma$ -algebra generated by  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k$ )

# Basic assumptions on the search directions

## Deterministic case:

$c_i > 0$  constants

- ③ “Sufficient” descent direction:

$$\nabla\phi(\mathbf{x}_k)^\top \mathbf{d}_k \leq -c_2 \|\nabla\phi(\mathbf{x}_k)\|^2$$

- ④ Direction norm bounded by gradient:

$$\|\mathbf{d}_k\| \leq c_3 \|\nabla\phi(\mathbf{x}_k)\|$$

# Basic assumptions on the search directions

## Stochastic case:

$c_i > 0$  constants

- ③ Deviation from descent direction allowed:

$$\nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k | \mathcal{F}_k) \leq c_1 \delta_k - c_2 \|\nabla\phi(\mathbf{x}_k)\|^2, \quad \delta_k > 0, \quad \sum_k \alpha_k \delta_k < \infty$$

- ④ Direction norm bounded by **noisy** gradient:

$$\|\mathbf{d}_k\| \leq c_3 \|\mathbf{g}(\mathbf{x}_k)\| \quad \text{a.s.}$$

# Basic assumptions on the search directions

## Stochastic case:

$c_i > 0$  constants

- ③ Deviation from descent direction allowed:

$$\nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k | \mathcal{F}_k) \leq c_1 \delta_k - c_2 \|\nabla\phi(\mathbf{x}_k)\|^2, \quad \delta_k > 0, \quad \sum_k \alpha_k \delta_k < \infty$$

- ④ Direction norm bounded by **noisy** gradient:

$$\|\mathbf{d}_k\| \leq c_3 \|\mathbf{g}(\mathbf{x}_k)\| \quad \text{a.s.}$$

## Theorem

*Under the previous assumptions, the sequence  $\{\mathbf{x}_k\}$  converges to  $\mathbf{x}_*$  a.s.*



# Search directions using second-order information

## Further (reasonable) assumptions

- 6 Positive definite and bounded approximate Hessians:  $\mu I \preceq B(\mathbf{x}) \preceq LI$
- 7 Mutually independent noise terms  $\varepsilon_f(\mathbf{x})$ ,  $\varepsilon_g(\mathbf{x})$  and  $\varepsilon_B(\mathbf{x})$  (to be relaxed for finite-sum problems)

# Search directions using second-order information

## Further (reasonable) assumptions

- 6 Positive definite and bounded approximate Hessians:  $\mu I \preceq B(\mathbf{x}) \preceq LI$
- 7 Mutually independent noise terms  $\varepsilon_f(\mathbf{x})$ ,  $\varepsilon_g(\mathbf{x})$  and  $\varepsilon_B(\mathbf{x})$  (to be relaxed for finite-sum problems)

## Possible directions guaranteeing convergence:

- Newton directions:

$$B(\mathbf{x}_k)\mathbf{d}_k = -\mathbf{g}(\mathbf{x}_k)$$

- “Inexact” Newton directions:

$$\|B(\mathbf{x}_k)\mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)\| \leq \delta_k \gamma_k$$

$\gamma_k$  random variable with bounded variance

# Search directions using second-order information

## Further (reasonable) assumptions

- 6 Positive definite and bounded approximate Hessians:  $\mu I \preceq B(\mathbf{x}) \preceq LI$
- 7 Mutually independent noise terms  $\varepsilon_f(\mathbf{x})$ ,  $\varepsilon_g(\mathbf{x})$  and  $\varepsilon_B(\mathbf{x})$  (to be relaxed for finite-sum problems)

## Possible directions guaranteeing convergence:

- Newton directions:

$$B(\mathbf{x}_k)\mathbf{d}_k = -\mathbf{g}(\mathbf{x}_k)$$

- “Inexact” Newton directions:

$$\|B(\mathbf{x}_k)\mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)\| \leq \delta_k(\omega_1\eta_k + \omega_2\|\mathbf{g}(\mathbf{x}_k)\|)$$

$\omega_1, \omega_2 \geq 0$  constant,  $\eta_k$  random variable with bounded variance

# LSOS: Line-search SOS

- A harmonic step-length sequence ( $\sum_k \alpha_k = \infty$ ,  $\sum_k \alpha_k^2 < \infty$ ) may make the algorithm slow (the steplength becomes too small soon)
- Tuning is necessary to ensure reasonable results; if the steplengths are not small enough the algorithm may diverge

# LSOS: Line-search SOS

- A harmonic step-length sequence ( $\sum_k \alpha_k = \infty$ ,  $\sum_k \alpha_k^2 < \infty$ ) may make the algorithm slow (the steplength becomes too small soon)
- Tuning is necessary to ensure reasonable results; if the steplengths are not small enough the algorithm may diverge

**IDEA:** start with line search and move to harmonic step lengths only if the line search produces small step lengths

# LSOS: Line-search SOS

- A harmonic step-length sequence ( $\sum_k \alpha_k = \infty$ ,  $\sum_k \alpha_k^2 < \infty$ ) may make the algorithm slow (the steplength becomes too small soon)
- Tuning is necessary to ensure reasonable results; if the steplengths are not small enough the algorithm may diverge

**IDEA:** start with line search and move to harmonic step lengths only if the line search produces small step lengths

- At each step the direction is not guaranteed to be a descent direction for  $\phi(\mathbf{x})$

**IDEA:** use nonmonotone line search

# LSOS: Line-search SOS (cont'd)

---

## LSOS algorithm

---

- 1: given  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\eta \in (0, 1)$ ,  $t_{\min} > 0$  and  $\{\alpha_k\}, \{\delta_k\}, \{\zeta_k\} \subset \mathbb{R}_+$
- 2: set LSPHase = *active*
- 3: **for**  $k = 0, 1, 2, \dots$  **do**
- 4:   compute a search direction  $\mathbf{d}_k$  such that

$$\|B(\mathbf{x}_k)\mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)\| \leq \delta_k \|\mathbf{g}(\mathbf{x}_k)\|$$

10: **end for**

---

# LSOS: Line-search SOS (cont'd)

---

## LSOS algorithm

---

- 1: given  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\eta \in (0, 1)$ ,  $t_{\min} > 0$  and  $\{\alpha_k\}, \{\delta_k\}, \{\zeta_k\} \subset \mathbb{R}_+$
- 2: set LSphase = *active*
- 3: **for**  $k = 0, 1, 2, \dots$  **do**
- 4:   compute a search direction  $\mathbf{d}_k$  such that

$$\|B(\mathbf{x}_k)\mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)\| \leq \delta_k \|\mathbf{g}(\mathbf{x}_k)\|$$

- 5:   find a step length  $t_k$  as follows:
  - 6:     **if** LSphase = *active* **then** find  $t_k$  that satisfies
$$f(\mathbf{x}_k + t_k \mathbf{d}_k) \leq f(\mathbf{x}_k) + \eta t_k \mathbf{g}(\mathbf{x}_k)^\top \mathbf{d}_k + \zeta_k$$
  - 7:     **if**  $t_k < t_{\min}$  **then** set LSphase = *inactive*
  - 8:     **if** LSphase = *inactive* **then** set  $t_k = \alpha_k$
  - 9:   set  $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$
  - 10: **end for**
-



## Theorem

Assume that  $\{\zeta_k\}$  is summable and the objective function estimator  $f$  is unbiased, i.e.

$$\mathbb{E}(\varepsilon_f(\mathbf{x})|\mathcal{F}_k) = 0.$$

If the sequence  $\{\mathbf{x}_k\}$  generated by LSOS is bounded, then  $\mathbf{x}_k \rightarrow \mathbf{x}_*$  a.s..

# Outline

- 1 Problem, motivations and contribution
- 2 The LSOS framework
- 3 Numerical experiments with LSOS**
- 4 Specializing LSOS for finite sums
- 5 Numerical experiments with LSOS-BFGS
- 6 Conclusions and future work

# Convex random problems (type 1)

$$\phi(\mathbf{x}) = \sum_{i=1}^n \lambda_i (e^{x_i} - x_i) + (\mathbf{x} - \mathbf{1})^\top A(\mathbf{x} - \mathbf{1})$$

# Convex random problems (type 1)

$$\phi(\mathbf{x}) = \sum_{i=1}^n \lambda_i (e^{x_i} - x_i) + (\mathbf{x} - \mathbf{1})^\top A (\mathbf{x} - \mathbf{1})$$

- $\lambda_i$ 's logarithmically spaced between 1 and  $\kappa$
- $A \in \mathbb{R}^{n \times n}$  spd with eigenvalues  $\lambda_i$  (generated by sprandsym)
- $n = 10^3$ ,  $\kappa = 10^2, 10^3, 10^4$
- $\varepsilon_f(\mathbf{x}) \sim \mathcal{N}(0, \sigma)$ ,  $(\varepsilon_g(\mathbf{x}))_i \sim \mathcal{N}(0, \sigma)$  and  
 $\varepsilon_B(\mathbf{x}) = \text{diag}(\mu_1, \dots, \mu_n)$ ,  $\mu_i \sim \mathcal{N}(0, \sigma)$
- $\sigma = 0.1\% \kappa, 0.5\% \kappa, 1\% \kappa$
- $x_*$  computed with high accuracy using deterministic L-BFGS  
(M. Schmidt, <https://www.cs.ubc.ca/~schmidtm/Software/minFunc.html>)

# Convex random problems (type 1)

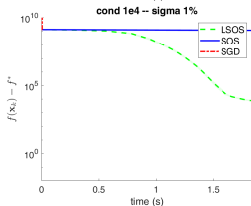
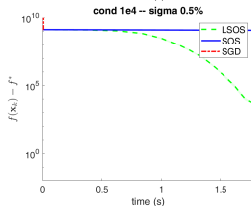
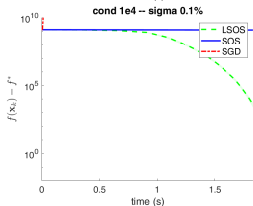
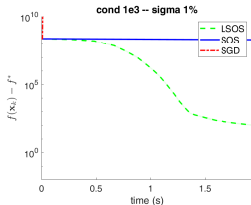
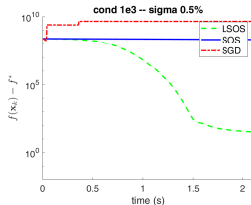
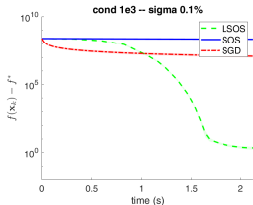
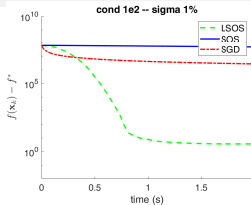
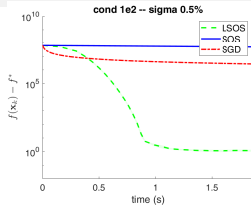
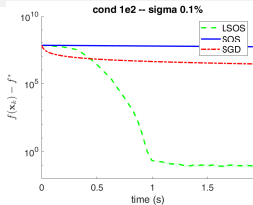
$$\phi(\mathbf{x}) = \sum_{i=1}^n \lambda_i (e^{x_i} - x_i) + (\mathbf{x} - \mathbf{1})^\top A (\mathbf{x} - \mathbf{1})$$

- $\lambda_i$ 's logarithmically spaced between 1 and  $\kappa$
- $A \in \mathbb{R}^{n \times n}$  spd with eigenvalues  $\lambda_i$  (generated by sprandsym)
- $n = 10^3$ ,  $\kappa = 10^2, 10^3, 10^4$
- $\varepsilon_f(\mathbf{x}) \sim \mathcal{N}(0, \sigma)$ ,  $(\varepsilon_g(\mathbf{x}))_i \sim \mathcal{N}(0, \sigma)$  and  
 $\varepsilon_B(\mathbf{x}) = \text{diag}(\mu_1, \dots, \mu_n)$ ,  $\mu_i \sim \mathcal{N}(0, \sigma)$
- $\sigma = 0.1\% \kappa, 0.5\% \kappa, 1\% \kappa$
- $x_*$  computed with high accuracy using deterministic L-BFGS  
(M. Schmidt, <https://www.cs.ubc.ca/~schmidtm/Software/minFunc.html>)

## Comparison of

- LSOS with exact solution of noisy Newton systems
- SOS with pre-defined step length  $\alpha_k = \frac{1}{\|\mathbf{d}_0\|} \frac{T}{T+k}$ ,  $T = 10^6$
- Stochastic Gradient Descent (SGD) with step length  $\alpha_k$

# Convex random problems (type 1): obj fun error vs time



## Convex random problems (type 2)

$$\phi(\mathbf{x}) = \sum_{i=1}^n \lambda_i (e^{x_i} - x_i) + (\mathbf{x} - \mathbf{1})^\top A (\mathbf{x} - \mathbf{1})$$

$$A = V D V^\top, \quad D = \text{diag}(\lambda_1, \dots, \lambda_n), \quad V = (I - 2 \mathbf{v}_3 \mathbf{v}_3^\top)(I - 2 \mathbf{v}_2 \mathbf{v}_2^\top)(I - 2 \mathbf{v}_1 \mathbf{v}_1^\top),$$

$$\mathbf{v}_j \text{ random, } \|\mathbf{v}_j\| = 1$$

- $n = 2 \cdot 10^4$ ,  $\kappa = 10^2, 10^3, 10^4$
- $\sigma = 0.1\% \kappa, 0.5\% \kappa, 1\% \kappa$
- Hessian in factorized form  $\implies$  (noisy) Newton system must be solved inexactly (e.g., by CG)

## Convex random problems (type 2)

$$\phi(\mathbf{x}) = \sum_{i=1}^n \lambda_i (e^{x_i} - x_i) + (\mathbf{x} - \mathbf{1})^\top A (\mathbf{x} - \mathbf{1})$$

$$A = V D V^\top, \quad D = \text{diag}(\lambda_1, \dots, \lambda_n), \quad V = (I - 2\mathbf{v}_3\mathbf{v}_3^\top)(I - 2\mathbf{v}_2\mathbf{v}_2^\top)(I - 2\mathbf{v}_1\mathbf{v}_1^\top),$$

$$\mathbf{v}_j \text{ random, } \|\mathbf{v}_j\| = 1$$

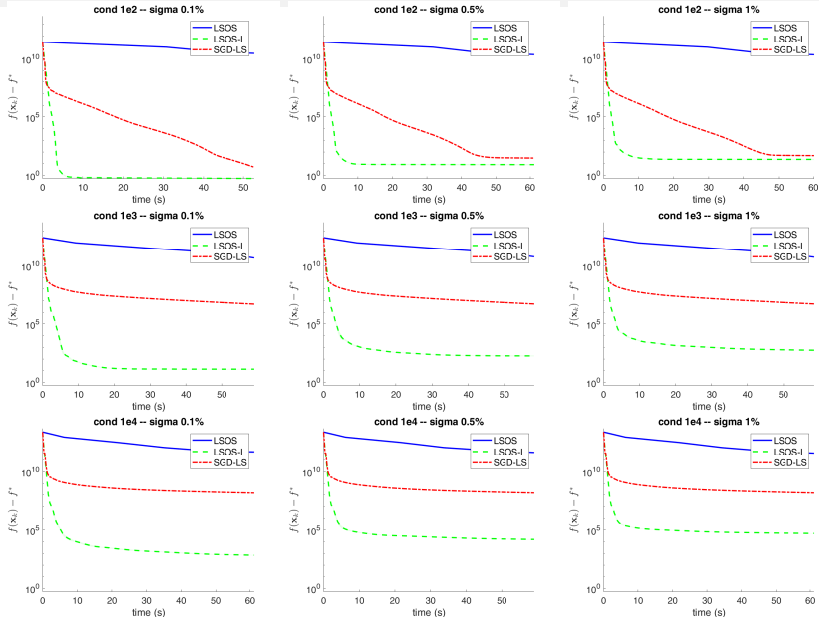
- $n = 2 \cdot 10^4$ ,  $\kappa = 10^2, 10^3, 10^4$
- $\sigma = 0.1\% \kappa, 0.5\% \kappa, 1\% \kappa$
- Hessian in factorized form  $\implies$  (noisy) Newton system must be solved inexactly (e.g., by CG)

### Comparison of

- **LSOS** (“exact” solution of noisy Newton systems - CG tolerance  $1e-6$ )
- **LSOS-I** (inexact solution of noisy Newton systems - decreasing tolerance sequence)
- **SGD-LS** (SGD with line search)



# Convex random problems (type 2): obj fun error vs time



# Outline

- 1 Problem, motivations and contribution
- 2 The LSOS framework
- 3 Numerical experiments with LSOS
- 4 Specializing LSOS for finite sums**
- 5 Numerical experiments with LSOS-BFGS
- 6 Conclusions and future work

# The finite sum case

$$\phi(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \phi_i(\mathbf{x})$$

$\phi_i(\mathbf{x}) \in \mathcal{C}^2$   $\bar{\mu}$ -strongly convex, with Lipschitz-continuous gradient with constant  $\bar{L}$

# The finite sum case

$$\phi(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \phi_i(\mathbf{x})$$

$\phi_i(\mathbf{x}) \in \mathcal{C}^2$   $\bar{\mu}$ -strongly convex, with Lipschitz-continuous gradient with constant  $\bar{L}$

**Subsampling:** at each iter  $k$ , a sample  $\mathcal{N}_k$  of size  $N_k \ll N$  is chosen randomly and uniformly from  $\mathcal{N} = \{1, \dots, N\}$ :

$$f_{\mathcal{N}_k}(\mathbf{x}) = \frac{1}{N_k} \sum_{i \in \mathcal{N}_k} \phi_i(\mathbf{x}), \quad \mathbf{g}_{\mathcal{N}_k}(\mathbf{x}) = \frac{1}{N_k} \sum_{i \in \mathcal{N}_k} \nabla \phi_i(\mathbf{x}),$$

$$B_{\mathcal{N}_k}(\mathbf{x}) = \frac{1}{N_k} \sum_{i \in \mathcal{N}_k} \nabla^2 \phi_i(\mathbf{x})$$

(unbiased estimators of  $\phi(\mathbf{x})$ ,  $\nabla \phi(\mathbf{x})$  and  $\nabla^2 \phi(\mathbf{x})$ )

# Stochastic variant of L-BFGS

Hessian approximation from stochastic variant of Limited-memory BFGS (L-BFGS)  
[Byrd, Hansen, Nocedal & Singer, SIOPT 2016]

$H_k$  defined by applying  $m$  BFGS updates to an initial matrix, using the  $m$  most recent correction pairs  $(s_j, y_j)$  obtained averaging iterates over  $r$  steps ( $j = k/r$ ):

$$H_k = H_k^{(m)}, \quad \text{where} \quad H_k^{(0)} = \frac{s_m^\top y_m}{\|y_m\|^2} I$$
$$H_k^{(j)} = \left( I - \frac{s_j y_j^\top}{s_j^\top y_j} \right)^\top H_k^{(j-1)} \left( I - \frac{y_j s_j^\top}{s_j^\top y_j} \right) + \frac{s_j s_j^\top}{s_j^\top y_j}, \quad j = 1, \dots, m$$

$$s_j = w_j - w_{j-1}, \quad y_j = B_{\mathcal{T}_j}(w_j) s_j, \quad \mathcal{T}_j \subset \{1, \dots, N\}$$

$$w_j = \frac{1}{r} \sum_{i=k-r+1}^k x_i, \quad w_{j-1} = \frac{1}{r} \sum_{i=k-2r+1}^{k-r} x_i$$

# Mini-batch SAGA

Subsampled gradient estimate by a mini-batch variant of SAGA

[Defazio, Bach & Lacoste-Julien, NIPS 2014; Gower, Richtárik & Bach, Math Prog 2020]

$$\mathbf{g}_{\mathcal{N}_k}^{\text{SAGA}}(\mathbf{x}_k) = \frac{1}{N_k} \sum_{i \in \mathcal{N}_k} \left( \nabla \phi_i(\mathbf{x}_k) - \mathbf{J}_k^{(i)} \right) + \frac{1}{N} \sum_{r=1}^N \mathbf{J}_k^{(r)}$$

$$\mathbf{J}_{k+1}^{(i)} = \begin{cases} \mathbf{J}_k^{(i)} & \text{if } i \notin \mathcal{N}_k \\ \nabla \phi_i(\mathbf{x}_{k+1}) & \text{if } i \in \mathcal{N}_k \end{cases}, \quad \mathbf{J}_0^{(i)} = \nabla \phi_i(\mathbf{x}_0)$$

$\{1, \dots, N\}$  partitioned into a fixed number  $n_b$  of random mini-batches, which are used in order

**Advantage of SAGA over SVRG:** full gradient computation only at the beginning of the algorithm (SVRG: full gradient computation each  $n_b$  iterations)

# LSOS-BFGS: Finite-Sum LSOS with L-BFGS

---

## LSOS-BFGS

---

- 1: given  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $m, r \in \mathbb{N}$ ,  $\eta, \vartheta \in (0, 1)$
- 2: **for**  $k = 0, 1, 2, \dots$  **do**
- 3:   compute a partition  $\{\mathcal{K}_0, \mathcal{K}_1, \dots, \mathcal{K}_{n_b-1}\}$  of  $\{1, \dots, N\}$

13: **end for**

---





# LSOS-BFGS: Finite-Sum LSOS with L-BFGS

---

## LSOS-BFGS

---

- 1: given  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $m, r \in \mathbb{N}$ ,  $\eta, \vartheta \in (0, 1)$
  - 2: **for**  $k = 0, 1, 2, \dots$  **do**
  - 3:   compute a partition  $\{\mathcal{K}_0, \mathcal{K}_1, \dots, \mathcal{K}_{n_b-1}\}$  of  $\{1, \dots, N\}$
  - 4:   **for**  $s = 0, \dots, n_b - 1$  **do**
  - 5:     choose  $\mathcal{N}_k = \mathcal{K}_s$  and compute  $\mathbf{g}(\mathbf{x}_k) = \mathbf{g}_{\mathcal{N}_k}^{\text{SAGA}}(\mathbf{x}_k)$
  - 6:     compute  $\mathbf{d}_k = -H_k \mathbf{g}(\mathbf{x}_k)$  with  $H_k$  defined by stochastic L-BFGS
  - 7:     find a step length  $t_k$  such that
$$f_{\mathcal{N}_k}(\mathbf{x}_k + t_k \mathbf{d}_k) \leq f_{\mathcal{N}_k}(\mathbf{x}_k) + \eta t_k \mathbf{g}(\mathbf{x}_k)^\top \mathbf{d}_k + \vartheta^k$$
  - 8:     set  $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$ ;
  - 9:     **if**  $\text{mod}(k, r) = 0$  and  $k \geq 2r$  **then**
  - 10:       update the L-BFGS correction pairs
  - 11:     **end if**
  - 12:   **end for**
  - 13: **end for**
-

# FS-LSOS: convergence

## Theorem (convergence)

*Assume  $\{t_k\}$  is bounded away from zero. Then  $\{x_k\}$  converges a.s. to the unique minimizer of  $\phi$ .*

# FS-LSOS: convergence

## Theorem (convergence)

Assume  $\{t_k\}$  is bounded away from zero. Then  $\{\mathbf{x}_k\}$  converges a.s. to the unique minimizer of  $\phi$ .

## Theorem (convergence rate)

Let  $\{t_k\}$  be bounded away from zero. Then there exist  $\rho \in (0, 1)$  and  $C > 0$  such that

$$\mathbb{E}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \leq C\rho^k.$$

## Theorem (complexity bound)

In order to achieve  $\mathbb{E}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \leq \varepsilon$  for some  $\varepsilon \in (0, e^{-1})$ , LSOS-FS takes at most

$$k_{\max} = \left\lceil \frac{|\log(C)| + 1}{|\log(\rho)|} \log(\varepsilon^{-1}) \right\rceil = \mathcal{O}(\log(\varepsilon^{-1}))$$

with  $\rho \in (0, 1)$  and  $C > 0$ .

# Outline

- 1 Problem, motivations and contribution
- 2 The LSOS framework
- 3 Numerical experiments with LSOS
- 4 Specializing LSOS for finite sums
- 5 Numerical experiments with LSOS-BFGS**
- 6 Conclusions and future work

# Linear classification problems

Training a linear classifier by minimizing the  $\ell_2$ -regularized logistic regression

Given  $N$  pairs  $(\mathbf{a}_i, b_i)$ ,  $\mathbf{a}_i \in \mathbb{R}^n$  training point,  $b_i \in \{-1, 1\}$  corresponding label, a hyperplane approximately separating the two classes can be found by minimizing

$$\phi(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \phi_i(\mathbf{x}), \quad \text{with } \phi_i(\mathbf{x}) = \log(1 + e^{-b_i \mathbf{a}_i^\top \mathbf{x}}) + \frac{\mu}{2} \|\mathbf{x}\|^2, \quad \mu > 0$$

# Linear classification problems

Training a linear classifier by minimizing the  $\ell_2$ -regularized logistic regression

Given  $N$  pairs  $(\mathbf{a}_i, b_i)$ ,  $\mathbf{a}_i \in \mathbb{R}^n$  training point,  $b_i \in \{-1, 1\}$  corresponding label, a hyperplane approximately separating the two classes can be found by minimizing

$$\phi(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \phi_i(\mathbf{x}), \quad \text{with } \phi_i(\mathbf{x}) = \log\left(1 + e^{-b_i \mathbf{a}_i^\top \mathbf{x}}\right) + \frac{\mu}{2} \|\mathbf{x}\|^2, \quad \mu > 0$$

Note that

$$\nabla \phi_i(\mathbf{x}) = \frac{1 - z_i(\mathbf{x})}{z_i(\mathbf{x})} b_i \mathbf{a}_i + \mu \mathbf{x}, \quad \nabla^2 \phi_i(\mathbf{x}) = \frac{z_i(\mathbf{x}) - 1}{z_i^2(\mathbf{x})} \mathbf{a}_i \mathbf{a}_i^\top + \mu I, \quad z_i(\mathbf{x}) = 1 + e^{-b_i \mathbf{a}_i^\top \mathbf{x}}$$

$\Downarrow$

$$\phi_i \text{ } \mu\text{-strongly convex,} \quad \mu I \preceq \nabla^2 \phi_i(\mathbf{x}) \preceq LI, \quad L = \mu + \max_{i=1, \dots, N} \|\mathbf{a}_i\|^2$$

# Linear classification problems (cont'd)

LIBSVM datasets (<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/>)

name	$N$	$n$
covtype	406709	54
w8a	49749	300
epsilon	400000	2000
gisette	6000	5000
real-sim	50617	20958
rcv1	20242	47236

NOTE:  $\mu = 1/N$ , sample size =  $\lceil \sqrt{N} \rceil$

# Linear classification problems (cont'd)

LIBSVM datasets (<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/>)

name	$N$	$n$
covtype	406709	54
w8a	49749	300
epsilon	400000	2000
gisette	6000	5000
real-sim	50617	20958
rcv1	20242	47236

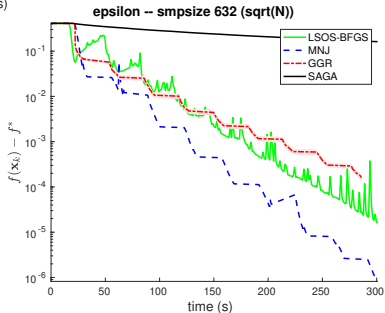
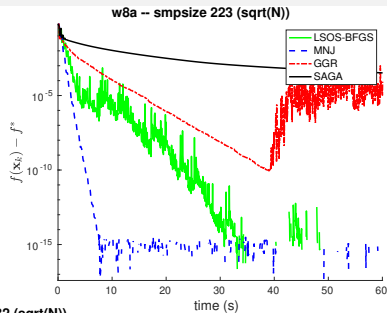
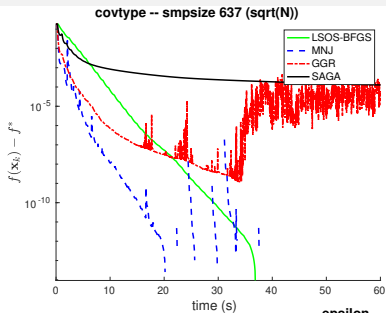
NOTE:  $\mu = 1/N$ , sample size =  $\lceil \sqrt{N} \rceil$

Comparison between

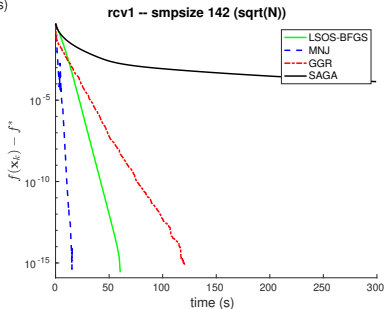
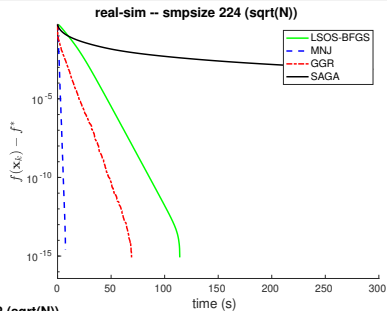
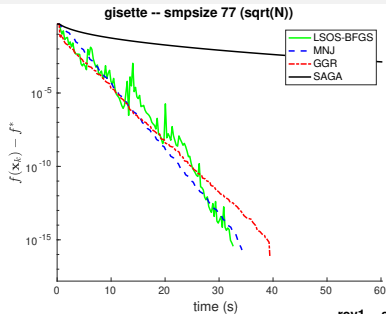
- **LSOS-BFGS**, with  $m = 10$  and  $r = 5$
- **GGR** [Gower, Goldfarb & Richtárik, Proc ICML 2016]
- **MNJ** [Moritz, Nishihara & Jordan, Proc MLR 2016]
- **Mini-batch variant of SAGA**, with the same line search as LSOS-BFGS



# Classification problems: obj fun error vs time



# Classification problems: obj fun error vs time



# Outline

- 1 Problem, motivations and contribution
- 2 The LSOS framework
- 3 Numerical experiments with LSOS
- 4 Specializing LSOS for finite sums
- 5 Numerical experiments with LSOS-BFGS
- 6 Conclusions and future work**

# Conclusions and future work

- We introduced LSOS a flexible second-order framework for optimization in noisy environments
- Almost sure convergence holds for the sequences generated by all the LSOS variants
- For finite-sum problems, we proved linear convergence rate on the obj. fun. error and worst-case complexity bound  $\mathcal{O}(\log(\varepsilon^{-1}))$  for LSOS with stochastic L-BFGS Hessian and any Lipschitz-continuous unbiased gradient estimates are used

# Conclusions and future work

- We introduced LSOS a flexible second-order framework for optimization in noisy environments
- Almost sure convergence holds for the sequences generated by all the LSOS variants
- For finite-sum problems, we proved linear convergence rate on the obj. fun. error and worst-case complexity bound  $\mathcal{O}(\log(\varepsilon^{-1}))$  for LSOS with stochastic L-BFGS Hessian and any Lipschitz-continuous unbiased gradient estimates are used
- Numerical experiments confirm that line-search techniques in second-order stochastic methods yield a significant improvement over predefined step-length sequences
- For finite sum problems LSOS-BFGS highly competitive with state-of-the art second-order stochastic optimization methods

# Conclusions and future work

- We introduced LSOS a flexible second-order framework for optimization in noisy environments
- Almost sure convergence holds for the sequences generated by all the LSOS variants
- For finite-sum problems, we proved linear convergence rate on the obj. fun. error and worst-case complexity bound  $\mathcal{O}(\log(\varepsilon^{-1}))$  for LSOS with stochastic L-BFGS Hessian and any Lipschitz-continuous unbiased gradient estimates are used
- Numerical experiments confirm that line-search techniques in second-order stochastic methods yield a significant improvement over predefined step-length sequences
- For finite sum problems LSOS-BFGS highly competitive with state-of-the art second-order stochastic optimization methods
- **What's next?** Possible extension to problems not satisfying the strong convexity assumption and to constrained problems

# Thanks for the attention!

## Any questions?

Do you want to know more?

D. di Serafino, N. Krejić, N. Krklec Jerinkić, M. Viola, *LSOS: Line-search Second-Order Stochastic optimization methods*, submitted (also available on ArXiv and Optimization Online)