# Minimization by Incremental Stochastic Surrogate Optimization for Large Scale Nonconvex Problems Belhal Karimi<sup>1</sup>, Hoi-To Wai<sup>2</sup>, Eric Moulines<sup>3</sup> and Ping Li<sup>1</sup> Baidu Research<sup>1</sup>, Chinese University of Hong Kong<sup>2</sup>, Ecole Polytechnique<sup>3</sup>

# Large Scale Optimization

• **Objective:** *Constrained* minimization problem of a finite sum of functions:

$$\min_{oldsymbol{ heta}\in\Theta} \mathcal{L}(oldsymbol{ heta}) \coloneqq rac{1}{n} \sum_{i=1}^n \mathcal{L}_i(oldsymbol{ heta})$$
, (1)

where  $\mathcal{L}_i : \mathbb{R}^p \to \mathbb{R}$  is bounded from below and is (possibly) nonconvex and include a nonsmooth penalty. • The gap  $\widehat{e}(\boldsymbol{\theta}; \{\overline{\boldsymbol{\theta}}_i\}_{i=1}^n)$  is L-smooth.

**Asymptotic Stationary Point Condition**:  $f(\mathbf{0} + \mathbf{1})$   $f(\mathbf{0})$ 

$$f'(\boldsymbol{\theta}, \mathbf{d}) := \lim_{t \to 0^+} \frac{\tau(\boldsymbol{\theta} + \tau \mathbf{d}) - \tau(\boldsymbol{\theta})}{t} \ge 0.$$

# The MISO Method (Mairal, 2015)

- MISO has been proposed in (Mairal, 2015)
- Initialize the majorizing surrogate functions  $\mathcal{A}_{i}^{0}(\boldsymbol{\theta}) := \widehat{\mathcal{L}}_{i}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(0)}), i \in [n].$
- For k > 0: • Pick  $i_k$  uniformly from [n].
- Update  $\mathcal{A}_{i}^{k+1}(\boldsymbol{\theta})$  as:

$$\mathcal{A}_{i}^{k+1}(\boldsymbol{ heta}) = egin{cases} \widehat{\mathcal{L}}_{i}(\boldsymbol{ heta}; \boldsymbol{ heta}^{(k)}), & ext{if } i = i_{k} \ \mathcal{A}_{i}^{k}(\boldsymbol{ heta}), & ext{otherwise.} \end{cases}$$

• Set  $\boldsymbol{\theta}^{(k+1)} \in \operatorname*{arg\,min}_{\boldsymbol{\theta}\in\Theta} \frac{1}{n} \sum_{i=1}^{n} \mathcal{A}_{i}^{k+1}(\boldsymbol{\theta}).$ • Return:  $\theta^{(k+1)}$ .

#### An Intractability for Latent Data Models

• Case when the surrogate functions computed in MISO are not tractable. • Surrogate is expressed as an integral over a set of latent variables z.

$$\widehat{\mathcal{L}}_{i}(\boldsymbol{\theta};\overline{\boldsymbol{\theta}}) := \int_{Z} r_{i}(\boldsymbol{\theta};\overline{\boldsymbol{\theta}},z_{i})p_{i}(z_{i};\overline{\boldsymbol{\theta}})\mu_{i}(dz_{i}) . \quad (2)$$

• Our scheme is based on the computation of  $\mathcal{L}_i(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, \{z_m\}_{m=1}^M)$ , a Monte Carlo approximation of the surrogate function  $\widehat{\mathcal{L}}_i(\theta; \overline{\theta})$  defined by (2) such that:

$$\widetilde{\mathcal{L}}_{i}(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, \{z_{m}\}_{m=1}^{M}) := \frac{1}{M} \sum_{m=1}^{M} r_{i}(\boldsymbol{\theta}; \overline{\boldsymbol{\theta}}, z_{m}),$$
 (3)

where  $\{z_i^m\}_{m=0}^{M-1}$  is a Monte Carlo batch.

# **MISSO:** Minimization by Incremental Stochastic Surrogate Optimization

Idea: We replace the expectation in (2) by a Monte Carlo integration and then optimizes the objective function in an incremental manner...

#### Algorithm 2 The MISSO method.

- 1: Input: initialization  $\theta^{(0)}$ ; a sequence of non-negative numbers  $\{M_{(k)}\}_{k=0}^{\infty}$ .
- 2: For all  $i \in [[1, n]]$ , draw  $M_{(0)}$  Monte Carlo samples with the stationary distribution  $p_i(\cdot; \theta^{(0)})$ .
- 3: Initialize the surrogate function as

$$\widetilde{\mathcal{A}}_{i}^{0}(\boldsymbol{ heta}) := \widetilde{\mathcal{L}}_{i}(\boldsymbol{ heta}; \boldsymbol{ heta}^{(0)}, \{z_{i,m}^{(0)}\}_{m=1}^{M_{(0)}}), \ i \in \llbracket 1, n 
rbracket .$$

- 4: for  $k = 0, 1, ..., K_{max}$  do
- Pick a function index  $i_k$  uniformly on  $[\![1, n]\!]$ . 5:
- Draw  $M_{(k)}$  Monte Carlo samples with the stationary distribution  $p_i(\cdot; \boldsymbol{\theta}^{(k)})$ . 6:
- Update the individual surrogate functions recursively as: 7:

$$\widetilde{\mathcal{A}}_{i}^{k+1}(oldsymbol{ heta}) = egin{cases} \widetilde{\mathcal{L}}_{i}(oldsymbol{ heta};oldsymbol{ heta}^{(k)},\{oldsymbol{x}\ \widetilde{\mathcal{A}}_{i}^{k}(oldsymbol{ heta}), \end{cases}$$

Set  $\boldsymbol{\theta}^{(k+1)} \in \arg\min_{\boldsymbol{\theta}\in\Theta} \widetilde{\mathcal{L}}^{(k+1)}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}$ 9: **end for** 

## How does this Monte Carlo Approximation affects the convergence rate?

## **Global Convergence Analysis**

• We bound the **Monte Carlo noise** by some constants using the notion of metric entropy (or bracketing number) see Van der Vaart (2000):  $\mathbb{E}\left|\sup_{M}\left|\frac{1}{M}\sum_{i=1}^{M}f(z_{i,m})-\mathbb{E}[f(z_{i})]\right|\right| \leq \frac{CL}{\sqrt{M}}$ . • Convergence Analysis: Assume  $\widehat{\mathcal{L}}_i(\theta; \overline{\theta}) \geq \mathcal{L}_i(\theta)$  and the error gap follows  $\|\nabla \widehat{e}(\theta; \{\overline{\theta}_i\}_{i=1}^n)\|^2 \leq 2L \widehat{e}(\theta; \{\overline{\theta}_i\}_{i=1}^n)$ 

**Theorem (MISSO Finite-Time Convergence)** For any  $K_{\max} \ge 1$ ,  $K \sim \mathcal{U}([0, K_{\max} - 1])$  independent of the  $\{i_k\}_{k=0}^{K_{\max}}$ , we have the **global rate**:

where

$$\mathbb{E}\left[\|\nabla\widehat{e}^{(\mathcal{K})}(\boldsymbol{\theta}^{(\mathcal{K})})\|^{2}\right] \leq \frac{\Delta_{(\mathcal{K}_{\max})}}{\mathcal{K}_{\max}} \quad \text{and} \quad \mathbb{E}\left[g_{-}(\boldsymbol{\theta}^{(\mathcal{K})})\right] \leq \sqrt{\frac{\Delta_{(\mathcal{K}_{\max})}}{\mathcal{K}_{\max}}} + \frac{C_{gr}}{\mathcal{K}_{\max}}\overline{\mathcal{M}}_{(\mathcal{K}_{\max})} \;.$$

$$\Delta_{(\mathcal{K}_{\max})} := 2nL\mathbb{E}\left[\widetilde{\mathcal{L}}^{(0)}(\boldsymbol{\theta}^{(0)}) - \widetilde{\mathcal{L}}^{(\mathcal{K}_{\max})}(\boldsymbol{\theta}^{(\mathcal{K}_{\max})})\right] + 4LC_{r}\overline{\mathcal{M}}_{(\mathcal{K}_{\max})} \;.$$

$$\tag{4}$$

The 33rd International Conference on Algorithmic Learning Theory, Paris, France

 $\{z_{i,m}^{(k)}\}_{m=1}^{M_{(k)}}$ , if  $i = i_k$ otherwise.

$$_{\scriptscriptstyle L}\widetilde{\mathcal{A}}_i^{k+1}(oldsymbol{ heta}).$$



#### Logistic Regression on Hemorrhage dataset

• Goal: Logistic Regression with missing values on Traumabase (severe hemorrhage):

• 16 quantitative measurements, like BMI, age, blood pressure, heart rate at different stages after the accident on 6384 patients • MISSO is an incremental MCEM in this case







 $p_i(y_i|z_i) = S(\boldsymbol{\delta}^ op ar{z}_i)^{y_i} \left(1 - S(\boldsymbol{\delta}^ op ar{z}_i)
ight)^{1-y_i}$  ,

• Goal: Train Bayesian variants of LeNet-5 and ResNet-18 on MNIST and CIFAR10:

• Variational inference and the ELBO loss to fit Bayesian Neural Networks on different datasets. • MISSO is an incremental VI in this case

#### References

Julien Mairal. Incremental majorization-minimization optimization with application to large-scale machine learning. SIAM Journal on Optimization, 25(2):829–855, 2015.

Aad W Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000.