

OPTIMIZATION METHODS USING RANDOM MODELS AND EXAMPLES FROM MACHINE LEARNING

Stefania Bellavia
Dipartimento di Ingegneria Industriale
Università di Firenze



UNIVERSITÀ
DEGLI STUDI
FIRENZE

Modern Techniques of Very Large Scale Optimization

Edinburgh, 19th-20th May 2022

Acknowledgements

- Gianmarco Gurioli, Benedetta Morini, Simone Rebegoldi
University of Florence, Italy
- Nataša Krejić
University of Novi Sad, Serbia
- Philippe Toint
University of Namur, Belgium

- Introduction: **random models** and motivating applications.
- Trust-region procedures with random models: **adaptive** choice of sample size and learning rate.
- Complexity results in **expectation**.
- Finite sum: trust region & **Inexact restoration**.
- Conclusions.

Unconstrained Optimization Problems

$$\min_{x \in \mathbb{R}^n} f(x),$$

with $f : \mathbb{R}^n \rightarrow \mathbb{R}$ sufficiently smooth ($f \in C^2$ for second-order methods), bounded below, possibly nonconvex.

- $f(x)$, $\nabla f(x)$ and $\nabla^2 f(x)$ evaluations are subject to **random noise** and we can only compute random estimates

$$\bar{f}(x) = \bar{f}(x, \xi), \quad \bar{\nabla} f(x) = \bar{\nabla} f(x, \xi) \quad \bar{\nabla}^2 f(x) = \bar{\nabla}^2 f(x, \xi)$$

where ξ is a random variable.

First-order

- first-order random model: $m(p) = f(x) + \overline{\nabla f(x)}^T p$;
- first-order regularized random model: $m(p) = f(x) + \overline{\nabla f(x)}^T p + \frac{\sigma}{2} \|x\|^2$, $\sigma > 0$
- ϵ - approximate first-order critical point:

$$\|\nabla f(\hat{x})\|_2 \leq \epsilon.$$

Second-order

- second-order random model: $m(p) = f(x) + \overline{\nabla f(x)}^T p + \frac{1}{2} p^T \overline{\nabla^2 f(x)} p$
- second-order regularized random model: $m(p) = f(x) + \overline{\nabla f(x)}^T p + \frac{1}{2} p^T \overline{\nabla^2 f(x)} p + \frac{\sigma}{3} \|x\|^3$
- ϵ approximate first and second-order critical point:

$$\begin{cases} \|\nabla f(\hat{x})\|_2 \leq \epsilon \\ \lambda_{\min}(\nabla^2 f(\hat{x})) \geq -\epsilon. \end{cases}$$

Motivating applications

Finite-sum minimization problem:

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{N} \sum_{i=1}^N \phi_i(x),$$

where $\phi_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, \dots, N$.

- Several problems can be cast in the previous form: classification, data fitting, sample average approximation ...
- Supervised machine learning: given a family of prediction function $h(\cdot; x)$, $x \in \mathbb{R}^n$, a loss function ℓ and a set of examples $\{(a_i, b_i)\}_{i=1}^N$ (training set), $a_i \in \mathbb{R}^d$ (feature), $b_i \in \mathbb{R}$ (label),

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{N} \sum_{i=1}^N \underbrace{\ell(h(a_i; x), b_i)}_{\phi_i(x)} \quad \text{Empirical Risk}$$

- The function f is often nonconvex, e.g. in the case of neural networks
- Big data applications $\Rightarrow N$ very large $\Rightarrow f$ and derivatives are very expensive!

Subsampled functions, gradients and Hessians

N is large

- M : sample size
- I_M : a randomly selected nonempty subset of $\{1, \dots, N\}$ of cardinality M

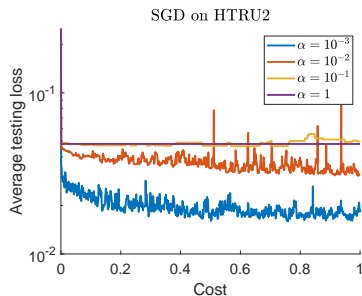
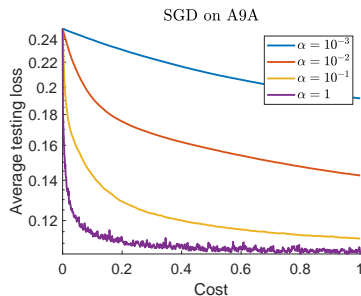
$$I_M \subseteq \{1, \dots, N\}, \quad |I_M| = M, \quad M \geq 1,$$

Use:

$$\begin{aligned}\bar{f}(x) &= \frac{1}{M} \sum_{i \in I_M} \phi_i(x) \\ \overline{\nabla f}(x) &= \frac{1}{M} \sum_{i \in I_M} \nabla \phi_i(x) \\ \overline{\nabla^2 f}(x) &= \frac{1}{M} \sum_{i \in I_M} \nabla^2 \phi_i(x)\end{aligned}$$

- A training set shows redundancy in the data \Rightarrow using all the sample data in every optimization iteration is inefficient
- Overall less expensive when N is large
- Computational evidence that they are more robust than fully deterministic approaches.

Stochastic gradient methods



$$x_{k+1} = x_k - \alpha_k \overline{\nabla f}(x_k), \quad k = 0, 1, \dots$$

- ✓ The expected value of the average norm of the gradients can be made small by picking a sufficiently small α
- ✗ ... but the smaller α , the slower the convergence rate!
- ✗ The optimal α (and the mini-batch size) are problem-dependent!
- ✗ For large-scale, real-world systems, expensive parameter tuning efforts is required!

Adaptive stochastic optimization methods

- SGD and its variants employ stochastic (possibly and occasionally full) gradient estimates and **do not rely on any machinery from standard globally convergent optimization procedures**, such as linesearch or trust-region.
- Strategies for selecting the steplength that mimic traditional step acceptance rules using stochastic estimates of functions and gradients:
 - Some criterion to **accept/reject the step** is tested
 - **Stochastic estimates** of functions and derivatives are computed.



random models are employed.

Bandeira, Vicente Scheinberg, SIOPT, 2014, [Trust-region](#)
Chen, Menickelly, Scheinberg, Math. Prog., 2018, [Trust-region](#)
Bollapragada, Byrd, and Nocedal, IMA JNA, [Inexact Newton](#)
Blanchet, Cartis, Menickelly, Scheinberg, INFORMS J. on Opt. 2019, [Trust-region](#)
B., Gurioli, Morini, Toint, SIOPT 2019, & J. of Complexity 2021, [Adaptive regularized](#)
B., Krejić, N Krklec Jerinkić, [Inexact-Newton](#), [Line-search](#) Paquette, Scheinberg, SIOPT 2020 [Line-search](#)
Xu, Roosta, Mahoney, Math. Prog. 2020 [Newton](#), [Trust-region](#) and [Adaptive regularized](#)
Berahas, Cao, Scheinberg, SIOPT 2021 [Line-search](#)
B., Gurioli, Morini, Toint, ArXiv, 2021 [Trust-region](#).
B., Krejić, Morini, Rebegoldi, ArXiv, 2021, [Trust-region](#)
di Serafino, Krejić, Krklec Jerinkić, Viola, ArXiv 2021, [Quasi-Newton](#), [Line-search](#)
Bergou, Diouane, Kunc, Kungursteu, Royer, INFORMS J. Optim., 2022, [Quasi-Newton](#), [Line-search](#)
Wang, Yuan, JCAM, 2022, [Trust-region](#)

Deterministic Trust-Region method

k th iteration

0. Given $x_k \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\gamma > 1$, and the trust-region radius $\delta_k > 0$.

1. Compute a trial step

Compute the model $m_k(p)$ and an (approximate) solution of the trust-region problem

$$\min_p m_k(p) \quad \text{s.t. } \|p\| \leq \delta_k$$

2. Check decrease

$$\rho_k(p_k) = \frac{f(x_k) - f(x_k + p_k)}{m_k(0) - m_k(p_k)}$$

3. Successful iteration

If $\rho_k \geq \eta$ then set $\delta_{k+1} = \gamma\delta_k$ and $x_{k+1} = x_k + p_k$.

4. Unsuccessful iteration

If $\rho_k < \eta$ then $\delta_{k+1} = \gamma^{-1}\delta_k$ and $x_{k+1} = x_k$

Trust-Region method with random models

k th iteration

0. Given $x_k \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\gamma > 1$, and the trust-region radius $\delta_k > 0$.
1. **Compute a trial step**
Compute a **random model** $\bar{m}_k(p)$ and an (approximate) solution of the trust-region problem

$$\min_p \bar{m}_k(p) \quad \text{s.t. } \|p\| \leq \delta_k$$

2. **Check decrease**

$$\rho_k(p_k) = \frac{f(x_k) - f(x_k + p_k)}{\bar{m}_k(0) - \bar{m}_k(p_k)}$$

3. **Successful iteration**
If $\rho_k \geq \eta$ then set $\delta_{k+1} = \gamma\delta_k$ and $x_{k+1} = x_k + p_k$.
4. **Unsuccessful iteration**
If $\rho_k < \eta$ then set $\delta_{k+1} = \gamma^{-1}\delta_k$ and $x_{k+1} = x_k$

Stochastic Trust-Region

k th iteration

0. Given $x_k \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\gamma > 1$, and the trust-region radius $\delta_k > 0$.
1. **Compute a trial step**
Compute a **random model** $\bar{m}_k(p)$ and an (approximate) solution of the trust-region problem

$$\min_p \bar{m}_k(p) \quad \text{s.t. } \|p\| \leq \delta_k$$

2. **Guess decrease**
Compute $\bar{f}(x_k)$ and $\bar{f}(x_k + p_k)$ estimate of $f(x_k)$ and $f(x_k + p_k)$ and

$$\rho_k(p_k) = \frac{\bar{f}(x_k) - \bar{f}(x_k + p_k)}{\bar{m}_k(0) - \bar{m}_k(p_k)}$$

3. **Successful iteration**
If $\rho_k \geq \eta$ then set $\delta_{k+1} = \gamma\delta_k$ and $x_{k+1} = x_k + p_k$.
4. **Unsuccessful iteration**
If $\rho_k < \eta$ then set $\delta_{k+1} = \gamma^{-1}\delta_k$ and $x_{k+1} = x_k$

 Blanchet, Cartis, Menickelly, Scheinberg, *INFORMS J. on Opt.* (2019)

 Wang, Yuan, *JCAM*, (2022)

 B., Gurioli, Morini, Toint, *arXiv:2112.06176* (2021)

Stochastic Trust-Region -First order method

kth iteration

0. Given $x_k \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\gamma > 1$, and the trust-region radius $\delta_k > 0$,
1. **Compute a trial step**
Compute a **random estimate** $\bar{\nabla}f(x_k)$ of $\nabla f(x_k)$ and set

$$p_k = - \underbrace{\frac{\delta_k}{\|\nabla f(x_k)\|}}_{\alpha_k} \bar{\nabla}f(x_k)$$

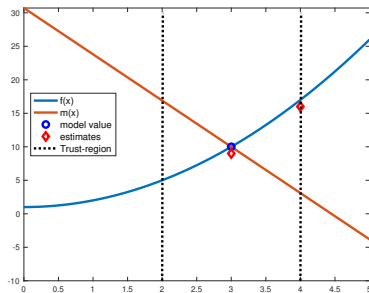
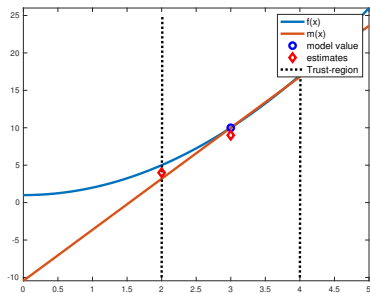
2. **Guess decrease**
Compute $\bar{f}(x_k)$ and $\bar{f}(x_k + p_k)$ estimate of $f(x_k)$ and $f(x_k + p_k)$ and

$$\rho_k(p_k) = \frac{\bar{f}(x_k) - \bar{f}(x_k + p_k)}{\|\bar{\nabla}f(x_k)\| \delta_k}$$

3. **Successful/unsuccessful iteration**
If $\rho_k \geq \eta$ then set $\delta_{k+1} = \gamma \delta_k$ and $x_{k+1} = x_k + p_k$.
If $\rho_k < \eta$ then set $\delta_{k+1} = \gamma^{-1} \delta_k$ and $x_{k+1} = x_k$

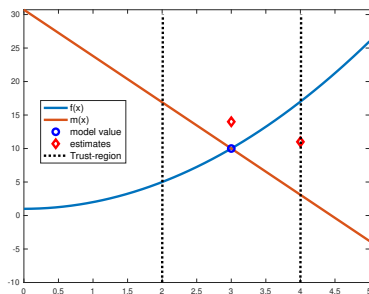
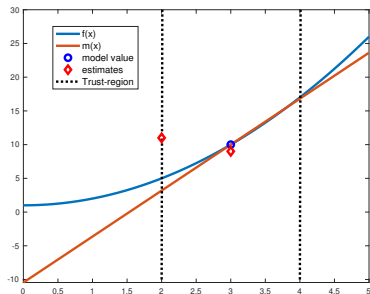
Stochastic gradient method with adaptive choice of the steplength (learning rate)! 

Possible iteration outcomes



Good model, good estimates (left), Bad model, good estimates (right)
True successful True unsuccessful

Possible iteration outcomes



Good model, bad estimates (left), Bad model, bad estimates (right)
False unsuccessful False successful

- What does it mean "good" model/estimations?
- How often can we have false successful/unsuccessful iterations?

Adaptive accuracy requirements

What does it mean "good" model/estimations?

- $m_k(p) = f(x_k) + \overline{\nabla f}(x_k)^T p$ is a "good" model if

$$\|\overline{\nabla f}(x_k) - \nabla f(x_k)\| \leq \nu \|\overline{\nabla f}(x_k)\| \quad \nu = \frac{1}{4}(1 - \eta)$$

- $\bar{f}(x_k)$ and $\bar{f}(x_k + p_k)$ are "good" function estimates if

$$\max\{|\bar{f}(x_k) - f(x_k)|, |\bar{f}(x_k + p_k) - f(x_k + p_k)|\} \leq \nu \|\overline{\nabla f}(x_k)\| \delta_k$$



B., Gurioli, Morini, Toint, [arXiv:2112.06176](https://arxiv.org/abs/2112.06176) (2021)

Similar accuracy requirements are used in other TR approaches and in linesearch and adaptive regularized methods.

The probabilistic setting

Consider the events

$$\mathcal{M}_k = \{\|G_k - \nabla f(X_k)\| \leq \nu \|G_k\|\}$$

$$\mathcal{F}_k = \{\max\{|F_k^0 - f(X_k)|, |F_k^p - f(X_k + P_k)|\} \leq \nu \|G_k\| \Delta_k \}$$

How often can we have false successful/unsuccessful iterations?

An *informal* statement of our assumptions:

We assume that

$$\text{Probability}[\mathcal{M}_k \cap \mathcal{F}_k | \text{conditioned by the past}] = p_* > \frac{1}{2}$$

the expected value of $f(X_k) - f(X_k + P_k)$ at false successful iterations, conditioned by the past, is positive.

+ f bounded below and Lipschitz continuity of $\nabla f(x)$

=====

X_k, Δ_k, P_k are the random variables corresponding to the realizations x_k, δ_k, p_k .

G_k is the random variable associate with the realization $\overline{\nabla f}(x_k)$.

F_k^0, F_k^p are the random variables associated with the realizations $\overline{f}(x_k), \overline{f}(x_k + p_k)$.

Iteration complexity

Let

$$N_\epsilon = \inf \{k \geq 0 \mid \|\nabla f(X_k)\| \leq \epsilon\}.$$

If the stochastic Trust-region algorithm is applied to the problem

$$\min f(x)$$

then, under the stated assumptions,

$$\mathbf{E}[N_\epsilon] = \mathbf{O}(\epsilon^{-2})$$

$\mathbf{O}(\epsilon^{-2})$ iteration bound is sharp for TR methods using exact function and gradient evaluations.

Probability p_* is constant along the iterations and we only require $p_* > 1/2$.
Accurate model and accurate functions “happen more often than not”



B., Gurioli, Morini, Toint [arXiv:2112.06176](https://arxiv.org/abs/2112.06176) (2021)

Ensuring the Accuracy Requirements

The Finite-Sum Minimisation Setting - Uniform Random Subsampling

Consider the *finite-sum* minimisation setting: $\min_{x \in \mathbb{R}^n} f(x)$, $f = \frac{1}{N} \sum_{i=1}^N \phi_i(x)$.

- Subsampling:

$$\bar{f}(x_k) = \frac{1}{|\mathcal{D}_k^f|} \sum_{i \in \mathcal{D}_k^f} \phi_i(x_k), \quad \nabla \bar{f}(x_k) = \frac{1}{|\mathcal{D}_k^g|} \sum_{i \in \mathcal{D}_k^g} \nabla \phi_i(x_k),$$

with $\mathcal{D}_k^f, \mathcal{D}_k^g \subseteq \{1, 2, \dots, N\}$ (randomly and uniformly taken).

- assume that $\exists \kappa_\phi(x_k) > 0$ s.t. $\kappa_\phi(x_k) \geq \max_{i \in \{1, \dots, N\}} \|\phi_i(x_k)\|$;
- Then, given the accuracy requirement ζ_k and a prefixed probability $\alpha_* \in (0, 1)$, using the Bernstein Inequality

$$|\mathcal{D}_k| \geq \min \left\{ N, \left\lceil \frac{4\kappa_\phi(x_k)}{\zeta_k} \left(\frac{2\kappa_\phi(x_k)}{\zeta_k} + \frac{1}{3} \right) \ln \left(\frac{1}{1 - \alpha_*} \right) \right\rceil \right\}$$

⇓

$$\Pr(|\bar{f}(x_k) - f(x_k)| \leq \zeta_k) \geq \alpha_*.$$



The Finite-Sum Minimisation Setting - Adaptive choice of the sample size

- Events:

$$\begin{aligned}\mathcal{M}_k &= \{\|G_k - \nabla f(X_k)\| \leq \nu \|G_k\|\}, \\ \mathcal{F}_k &= \{\max\{|F_k^0 - f(X_k)|, |F_k^p - f(X_k + P_k)|\} \leq \nu \|G_k\| \Delta_k\}\end{aligned}$$

- Given $\alpha_*, \beta_* \in [0, 1]$ such that $p_* = \alpha_* \beta_* > \frac{1}{2}$. if

$$|\mathcal{D}_k^f| = O\left(\frac{1}{\nu \|\overline{\nabla f}(x_k)\|^2 \delta_k^2} \log\left(\frac{1}{1 - \alpha_*}\right)\right)$$

$$|\mathcal{D}_k^g| = O\left(\frac{1}{\zeta_k^2} \log\left(\frac{1}{1 - \beta_*}\right)\right) \quad \zeta_k < \nu \|\overline{\nabla f}(x_k)\|$$

then

$$\text{Probability}[\mathcal{M}_k \cap \mathcal{F}_k | \text{conditioned by the past}] \geq p_*.$$

The Finite-Sum Minimisation Setting - Adaptive choice of the sample size

- Events:

$$\begin{aligned}\mathcal{M}_k &= \{ \|\mathbf{G}_k - \nabla f(\mathbf{X}_k)\| \leq \nu \|\mathbf{G}_k\| \}, \\ \mathcal{F}_k &= \{ \max\{|F_k^0 - f(\mathbf{X}_k)|, |F_k^p - f(\mathbf{X}_k + P_k)|\} \leq \nu \|\mathbf{G}_k\| \Delta_k \}\end{aligned}$$

- Given $\alpha_*, \beta_* \in [0, 1]$ such that $p_* = \alpha_* \beta_* > \frac{1}{2}$. if

$$|\mathcal{D}_k^f| = O\left(\frac{1}{\nu \|\overline{\nabla f}(\mathbf{x}_k)\|^2 \delta_k^2} \log\left(\frac{1}{1 - \alpha_*}\right)\right)$$

$$|\mathcal{D}_k^g| = O\left(\frac{1}{\zeta_k^2} \log\left(\frac{1}{1 - \beta_*}\right)\right) \quad \zeta_k < \nu \|\overline{\nabla f}(\mathbf{x}_k)\|$$

then

$$\text{Probability}[\mathcal{M}_k \cap \mathcal{F}_k | \text{conditioned by the past}] \geq p_*.$$

- The computation of \mathcal{D}_k^g requires an inner loop.
- This choice of $|\mathcal{D}_k^f|$ also provide a **positive expected value of $f(\mathbf{X}_k) - f(\mathbf{X}_k + P_k)$ at false successful iterations, conditioned by the past.**

An example: classification problems

- **Logistic loss:** given $\{(a_i, b_i)\}_{i=1}^N$

$$f(x) = \frac{1}{N} \sum_{i=1}^N \underbrace{\log(1 + e^{-b_i a_i^T x})}_{\phi_i(x)} + \frac{1}{2N} \|x\|^2,$$

- **Nonlinear least squares problems:** given $\{(a_i, b_i)\}_{i=1}^N$

$$f(x) = \frac{1}{N} \sum_{i=1}^N \underbrace{\left(b_i - \frac{1}{1 + e^{-a_i^T x}} \right)^2}_{\phi_i(x)}$$

The classifier is such that

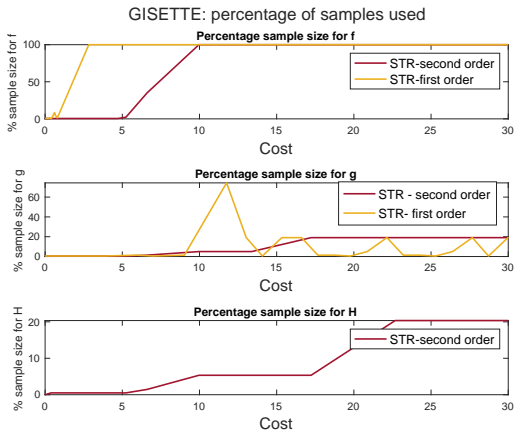
$$\begin{aligned} \frac{1}{1 + e^{-a_i^T x}} &\geq 0.5 & b_i = 1 \\ \frac{1}{1 + e^{-a_i^T x}} &< 0.5, & b_i = 0 \end{aligned}$$

- **Props:** Number of Propagations (1 full function and gradient evaluation is counted as 2 Prop). A maximum number of Props is considered as a termination criterion.
- Computing $\bar{f}(x)$ and $\bar{\nabla}f(x)$ costs $\frac{|\mathcal{D}_k^f| + |\mathcal{D}_k^g|}{N}$ props.

STR - first and second order: Adaptive sample size choice

N : 4800 $n = 5000$, Testing 1200

Average Accuracy STR- first order 87.85%, STR- second order 94.67%



Stochastic trust region & inexact restoration

k th iteration

0. Given $x_k \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\gamma > 1$, and the trust-region radius $\Delta_k > 0$,

1. Compute a trial step

Choose *randomly and uniformly* $\mathcal{D}_k^g \subseteq \{1, 2, \dots, N\}$, compute

$$\overline{\nabla f}(x_k) = \frac{1}{|\mathcal{D}_k^g|} \sum_{i \in \mathcal{D}_k^g} \nabla \phi_i(x_k) \text{ and set}$$

$$p_k = -\frac{\delta_k}{\|\overline{\nabla f}(x_k)\|} \overline{\nabla f}(x_k)$$

2. Guess decrease

Compute $\bar{f}(x_k + p_k)$ and $\bar{f}(x_k)$ by subsampling in \mathcal{D}_k^g
and $\rho_k(p_k)$ given by the inexact-restoration step acceptance rule.

3. Successful/unsuccessful iteration

If $\rho_k \geq \eta$ and $\|\overline{\nabla f}(x_k)\| \geq \eta_2 \delta_k$ then set $\delta_{k+1} = \gamma \delta_k$ and $x_{k+1} = x_k + p_k$.
Otherwise set $\delta_{k+1} = \gamma^{-1} \delta_k$ and $x_{k+1} = x_k$

Stochastic trust region & inexact restoration

k th iteration

0. Given $x_k \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\gamma > 1$, and the trust-region radius $\Delta_k > 0$,

1. Compute a trial step

Choose *randomly and uniformly* $\mathcal{D}_k^g \subseteq \{1, 2, \dots, N\}$, compute

$$\overline{\nabla f}(x_k) = \frac{1}{|\mathcal{D}_k^g|} \sum_{i \in \mathcal{D}_k^g} \nabla \phi_i(x_k) \text{ and set}$$

$$p_k = -\frac{\delta_k}{\|\overline{\nabla f}(x_k)\|} \overline{\nabla f}(x_k)$$

2. Guess decrease

Compute $\bar{f}(x_k + p_k)$ and $\bar{f}(x_k)$ by subsampling in \mathcal{D}_k^g
and $\rho_k(p_k)$ given by the inexact-restoration step acceptance rule.

3. Successful/unsuccessful iteration

If $\rho_k \geq \eta$ and $\|\overline{\nabla f}(x_k)\| \geq \eta_2 \delta_k$ then set $\delta_{k+1} = \gamma \delta_k$ and $x_{k+1} = x_k + p_k$.
Otherwise set $\delta_{k+1} = \gamma^{-1} \delta_k$ and $x_{k+1} = x_k$

The function approximation is computed averaging in the same subsample used for the gradient approximation!



B., Krejić, Morini, Rebegoldi *A stochastic first-order trust-region method with inexact restoration for finite-sum minimization*, Arxiv2107.03129, 2021

Inexact-restoration step acceptance

Given $x_k, \mathcal{D}_k^g, \mathcal{D}_{k-1}^g, \theta_k, p_k$.

- Let $\bar{f}_{k-1}(x_k) = \frac{1}{|\mathcal{D}_{k-1}^g|} \sum_{i \in \mathcal{D}_{k-1}^g} \phi_i(x_k)$ be the estimate computed at the previous iteration and

$$\rho_k = \frac{\text{Ared}_k(\theta_{k+1})}{\text{Pred}_k(\theta_{k+1})}$$

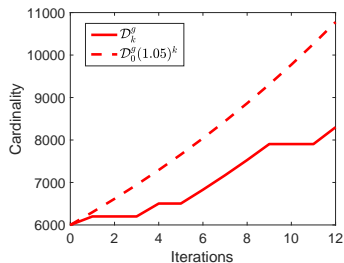
- $\text{Pred}_k(\theta_{k+1}) = \theta_{k+1}(\bar{f}_{k-1}(x_k) - \underbrace{(\bar{f}(x_k) + \nabla \bar{f}(x_k)^T p_k)}_{m_k(p_k)}) + (1 - \theta_{k+1}) \frac{|\mathcal{D}_k^g| - |\mathcal{D}_{k-1}^g|}{N}$
- $\text{Ared}_k(\theta_{k+1}) = \theta_{k+1}(\bar{f}_{k-1}(x_k) - \bar{f}(x_k + p_k)) + (1 - \theta_{k+1}) \frac{|\mathcal{D}_k^g| - |\mathcal{D}_{k-1}^g|}{N}$
- $\theta_{k+1} \in (0, 1)$ s.t.

$$\text{Pred}_k(\theta_{k+1}) \geq \eta \frac{|\mathcal{D}_k^g| - |\mathcal{D}_{k-1}^g|}{N}.$$

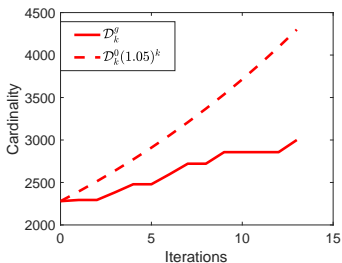
We balance the increase/decrease in the approximated objective function with the increase/decrease in the sample size.

History: sample size versus iterations

MNIST problem $N = 60000$
Average accuracy: 86,90%



A9A problem $N = 22793$.
Average accuracy: 98,32%



'-': SIRTR '- -': $\mathcal{D}_{k+1}^g = 1.05\mathcal{D}_k^g$

TRISH: trust-region without adaptive choice of the learning rate

SIRTR versus Trust-Region-ish algorithm (TRish)

TRish is a stochastic gradient method based on a trust-region methodology. Normalized steps are used in a dynamic manner whenever the norm of the stochastic gradient is within a prefixed interval. The k -th iteration of TRish is given by

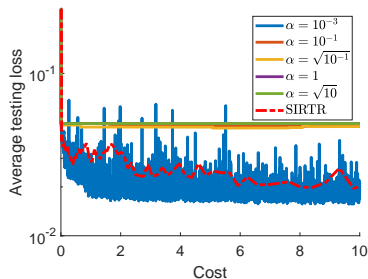
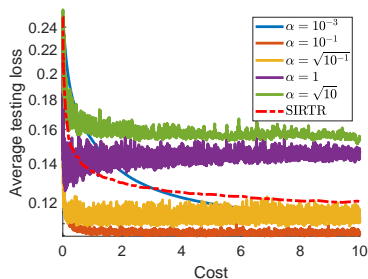
$$x_{k+1} = x_k - \begin{cases} \gamma_{1,k} \alpha_k \bar{\nabla} f(x_k), & \text{if } \|\bar{\nabla} f(x_k)\| \in \left[0, \frac{1}{\gamma_{1,k}}\right) \\ \alpha_k \frac{\bar{\nabla} f(x_k)}{\|\bar{\nabla} f(x_k)\|}, & \text{if } \|\bar{\nabla} f(x_k)\| \in \left[\frac{1}{\gamma_{1,k}}, \frac{1}{\gamma_{2,k}}\right] \\ \gamma_{2,k} \alpha_k \bar{\nabla} f(x_k), & \text{if } \|\bar{\nabla} f(x_k)\| \in \left(\frac{1}{\gamma_{2,k}}, \infty\right) \end{cases}$$

where $\alpha_k > 0$ is the steplength parameter and $0 < \gamma_{2,k} < \gamma_{1,k}$ are positive constants.



F.E. Curtis, K. Scheinberg, R. Shi, [INFORMS Journal on Optimization \(2019\)](#)

Avoiding learning rate tuning



SIRTR versus TRish algorithm for several choices of the steplength α .

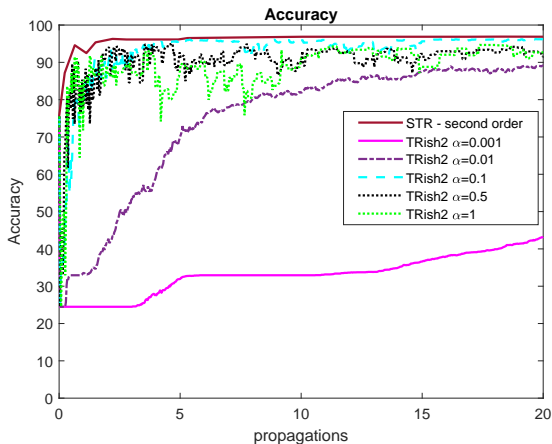
Decrease of the (average) testing loss $\bar{F}(x_k)$ w.r.t. the (average) computational time.

From left to right: a9a and htru2 datasets.

Avoiding Learning rate tuning (2)

Mushrooms dataset, Training $N = 5000$, $n = 112$, Testing 1600, batch-size=50

TRish2 $\gamma_1 = 4/G$, $\gamma_2 = 1/(2G)$ G : average norm of stochastic gradient estimates provided by SGD, $\alpha = 0.1$.



- The stochastic trust-region approach has been extended to **polynomial models of arbitrary degree**:
 - seek for first- and second-order critical points,
and also for critical points of arbitrary order
- **Adaptive accuracy**, finite sum context:
 - adaptive choice of the steplength and of the subsample sizes
- **Second order methods**:
 - Inexact steps + matrix-free implementation produce
a significant reduction of each iteration cost
- **More numerical results**: Training neural network for monitoring the electricity consumption of a healthcare facility.

- The stochastic trust-region approach has been extended to **polynomial models of arbitrary degree**:
 - seek for first- and second-order critical points,
and also for critical points of arbitrary order
- **Adaptive accuracy**, finite sum context:
 - adaptive choice of the steplength and of the subsample sizes
- **Second order methods**:
 - Inexact steps + matrix-free implementation produce
a significant reduction of each iteration cost
- **More numerical results**: Training neural network for monitoring the electricity consumption of a healthcare facility.

Thank you!

Some references



S.B., N.Krejić, B.Morini, S.Rebegoldi, *A stochastic first-order trust-region method with inexact restoration for finite-sum minimization*, Arxiv: 2107.03129, 2021



S. B., G. Gurioli, B. Morini and Ph. L. Toint, *Trust-region algorithms: probabilistic complexity and intrinsic noise with applications to subsampling techniques*, arXiv:2112.06176, 2021.



S. B., G. Gurioli, B. Morini and Ph. L. Toint, *Adaptive regularization for nonconvex optimization using inexact function values and randomly perturbed derivatives*, Journal of Complexity, 2022.



S. B., G. Gurioli, *Complexity Analysis of a Stochastic Cubic Regularisation Method under Inexact Gradient Evaluations and Dynamic Hessian Accuracy*, Optimization, 2021.



S. B., G. Gurioli, B. Morini, *Adaptive cubic regularization methods with dynamic inexact Hessian information and applications to finite-sum minimization*, IMA Journal Numerical Analysis, 2021.



S. B., N.Krejić, B. Morini *Inexact restoration with subsampled trust-region methods for finite-sum minimization*, COAP, 2020.



S. B., N.Krejić, N. Krklec Jerinkic, *Subsampled Inexact Newton methods for minimizing large sums of convex functions*, IMA J. Numer. Anal. 2020.

Grants: Indam-GNCS, UNIFI Internationalization plan, **“Second order methods for optimisation problems in machine learning”**, project for the exchange of researchers within the frame of the executive programme of Scientific and Technological cooperation between the Italian Republic and the Republic of Serbia for the years 2019-2022.

Approximated function and derivative evaluations

- $\overline{\nabla^j f}(x_k)$ of $\nabla^j f(x_k)$ (define $\nabla^0 f \stackrel{\text{def}}{=} f$) at iteration k are given by

$$\overline{\nabla^j f}(x_k) = \frac{1}{|\mathcal{D}_{k,j}|} \sum_{i=1}^N \nabla^j f_i(x_k), \quad j \in \{0, 1, 2\},$$

with $\mathcal{D}_{k,j} \subseteq \{1, 2, \dots, N\}$ (randomly and uniformly taken) such that

$$|\mathcal{D}_{k,j}| = \min \left\{ N, \left\lceil \frac{4\kappa_{f,j}}{\zeta_{k,j}} \left(\frac{2\kappa_{f,j}}{\zeta_{k,j}} + \frac{1}{3} \right) \log \left(\frac{d_j}{t} \right) \right\rceil \right\}, \quad j \in \{0, 1, 2\},$$

where $d_0 = 2$, $d_1 = n + 1$, $d_2 = 2n$, $t = 0.2$.

- $\kappa_{f,0} = 10^{-3}$, $\kappa_{f,1} = 5 \cdot 10^{-4}$, $\kappa_{f,2} = 10^{-4}$: set in order to control the growth of the sample sizes throughout the running of the algorithm