



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

The Limitation and Practical Acceleration of Stochastic Gradient Algorithms in Inverse Problems.

Citation for published version:

Tang, BJ, Egiazarian, K & Davies, M 2019, The Limitation and Practical Acceleration of Stochastic Gradient Algorithms in Inverse Problems. in *2019 IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2019 - Proceedings.*, 18778660, ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings, vol. 2019-May, Institute of Electrical and Electronics Engineers Inc., 44th IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2019, Brighton, United Kingdom, 12/05/19.
<https://doi.org/10.1109/ICASSP.2019.8683368>

Digital Object Identifier (DOI):

[10.1109/ICASSP.2019.8683368](https://doi.org/10.1109/ICASSP.2019.8683368)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

2019 IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2019 - Proceedings

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



THE LIMITATION AND PRACTICAL ACCELERATION OF STOCHASTIC GRADIENT ALGORITHMS IN INVERSE PROBLEMS

Junqi Tang[†], Karen Egiazarian*, Mike Davies[†]

[†] School of Engineering, University of Edinburgh, UK
*Noiseless Imaging Ltd, Finland.

ABSTRACT

In this work we investigate the practicability of stochastic gradient descent and recently introduced variants with variance-reduction techniques in imaging inverse problems, such as non-uniform image deblurring. Such algorithms have been shown in machine learning literature to have optimal complexities in theory, and provide great improvement empirically over the full gradient methods. Surprisingly, in some tasks such as image deblurring, many of such methods fails to converge faster than the accelerated full gradient method (FISTA), even in terms of epoch counts. We investigate this phenomenon and propose a theory-inspired mechanism to characterize whether a given inverse problem should be preferred to be solved by stochastic optimization technique with a known sampling pattern. Furthermore, to overcome another key bottleneck of stochastic optimization which is the heavy computation of proximal operators while maintaining fast convergence, we propose an accelerated primal-dual SGD algorithm and demonstrate the effectiveness of our approach in image deblurring experiments.

Index Terms— Stochastic Optimization, Inverse Problems, Image Processing

1. INTRODUCTION

The stochastic gradient methods [1, 2] and recently introduced variants with variance-reduction [3, 4, 5] have been widely used to solve large-scale convex optimization problem in machine learning applications. Such tasks can be formulated as the following:

$$x^* \in \arg \min_{x \in \mathcal{X}} \left\{ F(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) + \lambda g(x) \right\}, \quad (1)$$

where $\mathcal{X} \in \mathbb{R}^d$ is a close convex set and we denote $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ the data fidelity term. Each $f_i(x)$ is assumed to be convex and L -smooth, while the regularization term $g(x)$ is a simple convex function and is possibly non-smooth. With Nesterov’s acceleration [6, 7], researchers [8, 9, 10] have developed several “optimal” algorithms which can provably achieve the worse-case optimal convergence rate for (1).

While having been a proven success both in theory and in machine learning applications, there is no convincing result

so far in the literature which reports the performance of the stochastic gradient methods in image processing applications (except for tomography reconstruction [11, 12, 13]), which also involve large-scale optimization tasks in the same form of (1). In this work we investigate the practical performance of such methods, using non-uniform deblurring as a running example.

We make the following contributions:

(Evaluating the limitation of stochastic gradient algorithms.) We investigate the fundamental limit of possible acceleration of a stochastic gradient method over its full gradient counterpart by measuring the *Stochastic Acceleration* (SA) factor which is based on the ratio of the Lipschitz constants of the minibatched stochastic gradient and the full gradient. We discover that the SA factor is indeed able to characterize the potential of a certain optimization task being speedily solved by applying randomization techniques.

(Breaking the computational bottleneck of expensive/multiple proximal operators for Nesterov-type momentum SGD.) Another factor in image processing practice which significantly affects the SGD-type methods’ actual performance is the frequent calculation of the costly proximal operator for the regularization terms which have a linear operator, such as the TV semi-norm – SGD methods needs to calculate it much more frequently than full gradient methods. Moreover most of the fast SGD methods can not cope with more than one non-smooth regularization terms. To overcome these we propose an accelerated primal-dual SGD algorithm which can efficiently handle (1) regularization with a linear operator, (2) multiple regularization terms, while (3) maintaining Nesterov-type accelerated convergence speed in practice.

2. FAILURES OF STOCHASTIC OPTIMIZATION

We start by a simple non-uniform deblurring [14] example where the central part (sized 128 by 128) of the “Kodim05” image from *Kodak Lossless True Color Image Suite* [15] is blurred with an non-uniform blur kernel which imposes less blurring at the center but increasingly severe blurring towards the edge. We also add a small amount of noise to the blurred image.

We test the effectiveness of several algorithms by solving the same TV-regularized least-squares problem, to get an estimation of the ground truth image. The algorithms we test in the experiments include the accelerated full gradient method

FISTA [16], SGD with momentum [17], the proximal SVRG [18] and its accelerated variant, Katyusha [9].

Perhaps surprisingly, on this experiment we report a negative result for the randomized algorithms. The most efficient solver in this task is the full gradient method FISTA both in terms of wall clock time and number of datapasses. The state-of-the-art stochastic gradient method Katyusha even cannot beat FISTA in terms of epoch counts. For all the randomized algorithms we use a minibatch size which is 10 percent of the total data size. For stochastic gradient methods, smaller minibatch size in this case does not provide better performance in datapasses and will significantly slow down running time due to the multiple calls of proximal operator.

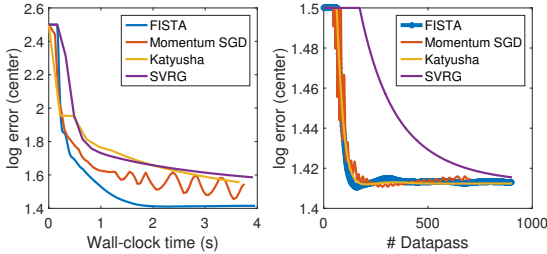


Fig. 1: The estimation error plot for the deblurring experiment. The plots correspond to the estimation error of the central part (100 by 100) of the image.

3. LIMITATIONS OF STOCHASTIC OPTIMIZATION

These results appear to be contrary to the popular belief among the stochastic optimization community, that stochastic gradient methods are much faster in terms of iteration complexity than deterministic gradient methods in solving large scale problems: to be specific – to achieve an objective gap suboptimality of $F(x) - F(x^*) \leq \varepsilon$, optimal stochastic gradient methods needs only $\Theta\left(n + \sqrt{nL/\varepsilon}\right)$ evaluations of ∇f_i , while $\Theta\left(n\sqrt{L/\varepsilon}\right)$ for optimal full gradient methods. Where is the loophole?

It is often easily ignored that the complexity results above are derived under different smoothness assumption. For the convergence bound of full gradient, the full smooth part of the cost function $f(\cdot)$ is assumed to be L -smooth, while for the case of stochastic gradient, every individual function $f_i(\cdot)$ is assumed to be L -smooth. Now we can clearly see the subtlety: to compare these complexity results and make meaningful conclusions, one has to assume that these two Lipschitz constants are roughly the same. While this is true for many problems but there are exceptions – image deblurring is one of them.

Given a minibatch index $[S_0, S_1, S_2, \dots, S_K]$ such that:

$$f(x) = \frac{1}{K} \sum_{k=1}^K f_{S_k}(x), \quad f_{S_k}(x) := \frac{K}{n} \sum_{i \in S_k} f_i(x), \quad (2)$$

In order to identify the potential of a certain optimization problem to be more efficiently solved by using stochastic gradient methods, we start by comparing the single iteration convergence of one instant of the proximal stochastic gradient descent with momentum named Katyusha [9] with the proximal accelerated full gradient descent (AFG), which read:

AFG (x_0, K, L) :

$$\begin{aligned} &\text{For } s = 0, 1, 2, \dots, S \\ &\left[\begin{array}{l} x^s = \mathcal{T}_f(y^s, L) := \text{prox}_{\lambda g}^{\frac{1}{L}}(y^s - \frac{1}{L} \nabla f(y^s)); \\ \quad \quad \quad \rightarrow \text{Proximal gradient descent} \\ a_{s+1} = (1 + \sqrt{1 + 4a_s^2})/2; \\ y^{s+1} = x^s + \frac{a_s - 1}{a_{s+1}}(x^s - x^{s-1}); \rightarrow \text{Momentum} \end{array} \right. \end{aligned}$$

Katyusha (x_0, S, m, L) :

$$\begin{aligned} &\text{For } s = 0, 1, 2, \dots, S \\ &\left[\begin{array}{l} \theta = \frac{2}{s+4}; \\ (\hat{x}^{s+1}, y^{s+1}, z^{s+1}) = \mathcal{A}(x^s, y^s, z^s, L, m, \theta, \nabla f(\hat{x}^s)); \end{array} \right. \\ &\mathcal{A}(x^s, y^s, z^s, L, m, \theta, \nabla f(\hat{x}^s)): \end{aligned}$$

For $k = 0, 1, 2, \dots, m$

$$\left[\begin{array}{l} x_{k+1} = \theta z_k + \frac{1}{2} \hat{x}^s + (\frac{1}{2} - \theta) y_k; \quad \rightarrow \text{Momentum} \\ \text{Pick } i \in [1, 2, \dots, K] \text{ uniformly at random} \\ \nabla_{k+1} = \nabla f(\hat{x}^s) + \nabla f_{S_i}(x_{k+1}) - \nabla f_{S_i}(\hat{x}^s); \\ \quad \rightarrow \text{Compute a variance reduced stochastic gradient} \\ z_{k+1} = \text{prox}_{\lambda g}^{\frac{1}{3\theta L}}(z_k - \frac{1}{3\theta L} \nabla_{k+1}); \\ \quad \quad \quad \rightarrow \text{Proximal mirror descent} \\ y_{k+1} = \text{prox}_{\lambda g}^{\frac{1}{3L}}(x_{k+1} - \frac{1}{3L} \nabla_{k+1}); \\ \quad \quad \quad \rightarrow \text{Proximal gradient descent} \end{array} \right.$$

where we define the proximal operator as:

$$\text{prox}_{\lambda g}^{\eta}(\cdot) = \arg \min_{x \in \mathcal{X}} \frac{1}{2\eta} \|x - \cdot\|_2^2 + \lambda g(x). \quad (3)$$

3.1. Analysis

We start with the standard smoothness assumption [19]:

A. 1 (Smoothness of the Full-Batch and the Mini-Batches.) $f(\cdot)$ is L_f -smooth and each f_{S_k} is L_m -smooth, that is:

$$f(x) - f(y) - \nabla f(y)^T(x - y) \leq \frac{L_f}{2} \|x - y\|_2^2, \quad \forall a, b \in \mathcal{X}, \quad (4)$$

and

$$f_{S_k}(x) - f_{S_k}(y) - \nabla f_{S_k}(y)^T(x - y) \leq \frac{L_b}{2} \|x - y\|_2^2, \quad (5)$$

$\forall x, y \in \mathcal{X}$.

Now we are ready to present the main theorem, which follows from simply combining the existing convergence results of Katyusha and AFG, as well as the lower bounds for the stochastic and deterministic first-order optimization [19, 20].

Theorem 3.1 Under A.1, when $\mathcal{X} = \{x \in \mathbb{R}^d : \|x\|_2^2 \leq 1\}$, $g(\cdot) = 0$, let $x_f^s = \mathcal{T}_f(y_f^{s-1}, L_f)$, $x_{\mathcal{A}}^s = \mathcal{A}(x_{\mathcal{A}}^{s-1}, L_b, m)$, $m = 2K$, we have:

$$\frac{\mathbb{E}F(x_{\mathcal{A}}^s) - F^*}{F(x_f^{3s}) - F^*} \leq 192 \cdot \frac{L_b}{KL_f} + \frac{512(F(x^0) - F^*)}{L_f \|x^0 - x^*\|_2^2} \quad (6)$$

and moreover,

$$\frac{\mathbb{E}F(x_{\mathcal{A}}^s) - F^*}{F(x_f^{3s}) - F^*} \geq \frac{1}{51840} \cdot \frac{L_b}{KL_f}. \quad (7)$$

for a sufficiently large dimension $d = O(\frac{L^2}{\epsilon^2} \log \frac{L}{\epsilon} + K \log K)$ where $\mathbb{E}F(x_{\mathcal{A}}^s) - F^* \leq \epsilon$.

From this theorem we can see that with the same epoch count (the iteration complexity of Katyusha's 1 epoch is equivalent to 3 iterations of AFG), the ratio of objective gap achieved by each algorithm can be upper and lower bounded by $\Theta(\frac{L_b}{KL_f})$. Although the constants seem pessimistic, it is within our expectation since the lower bounds on the convergence speed of both algorithms are derived on the worst possible function which satisfies A.1. Motivated by the theory, we propose to evaluate the potential of stochastic acceleration simply by the ratio $\frac{L_b}{KL_f}$ which dominates our upper and lower bounds in Theorem 3.1.

3.2. Evaluating the Limitation of SGD-type Algorithms

We introduced a metric called *Stochastic Acceleration* (SA) factor. The curve of SA factor as a function of the minibatch number K (for a given minibatch pattern) is able to provide a way of evaluating and characterizing inherently whether for a given inverse problem and a certain minibatch sampling scheme, randomized gradient methods should be preferred over the deterministic full gradient methods or not.

Definition 3.2 For a given data-partitioning index $\bar{S} = [S_1, \dots, S_K]$, the *Stochastic Acceleration* (SA) factor is defined as:

$$\Upsilon(\bar{S}) = \frac{KL_f}{L_b} \quad (8)$$

We test several least-squares loss $f(x) = \|Ax - b\|_2^2$ with different types of forward operator. In this case we have

$$f(x) = \|Ax - b\|_2^2 = \frac{1}{K} \sum_{k=1}^K f_{S_k}(x), \quad (9)$$

$$f_{S_k}(x) := K \|A_{S_k} x - b_{S_k}\|_2^2, \quad (10)$$

The examples of forward operator A we consider include the non-uniform deblurring (262144 by 262144), a random compressed sensing matrix with i.i.d Gaussian random entries (500 by 2000), a fan beam X-ray CT operator (91240 by 65536), and two machine learning datasets: RCV1 dataset (20242 by 47236), and Magic04 (19000 by 50, with random features). The data-partition we choose is the interleaving sampling. From the result show by the Fig 2 we find that indeed the stochastic methods have a limitation on some optimization

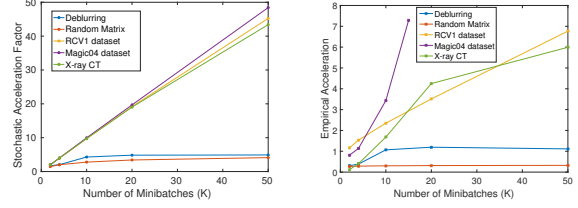


Fig. 2: Left: Stochastic Acceleration (SA) factor of inverse problems with different forward operators, Right: Empirical observation comparing the objective gap convergence of Katyusha and FISTA algorithm in 15 epochs.

problems like deblurring and inverse problems with random matrices, where we see that the curve of SA factor of such problems stays low and flat even when we increase the number of minibatches. For the machine learning datasets and X-ray CT imaging, the SA factor increases rapidly and almost linearly as we increase the number of minibatches, which is in line with observations in machine learning on the superiority of SGD and also the observation in CT image reconstruction of the benefits of the use of the ordered-subset methods [21]. The curve of SA factor on the left figure qualitatively predict the empirical comparison result of Katyusha and FISTA algorithm shown on the right, where we observe that Katyusha offers no acceleration over the FISTA on deblurring and Gaussian random inverse problem, but significantly outperforms FISTA on the other cases. Indeed, positive results for applying SGD-type algorithms on these problems are well-known already [2, 18, 21], hence we have shown that the SA factor we propose is useful in characterizing whether an inverse problem is inherently a suitable application for stochastic gradient methods.

4. PRACTICAL ACCELERATION FOR SGD

The previous section suggests that stochastic gradient methods do not always offer an intrinsic advantage for some problems. There are also several other causes for this failure. The most obvious one is that stochastic gradient methods in the primal need to calculate the proximal operator many more times than full gradient methods and hence slow down dramatically the run time. Moreover, in image processing practice often more than one non-smooth regularization term is used, where most of the existing fast stochastic methods such as Katyusha are inapplicable.

To avoid the frequent oracle call on the TV proximal operator, we can first reformulate the original optimization problem as a convex-concave saddle-point form. To be specific, the given problem:

$$x^* \in \min_{x \in \mathbb{R}^d} \{f(x) + \lambda g(Dx) + \gamma h(x)\}, \quad (11)$$

where $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(a_i, x)$ is the data-fidelity term, $g(Dx)$ is a regularization term with a linear operator $-$ for example the TV regularization ($g(\cdot) = \|\cdot\|_1$, $D \in \mathbb{R}^{r \times d}$ is the

7. REFERENCES

- [1] Herbert Robbins and Sutton Monro, “A stochastic approximation method,” in *Herbert Robbins Selected Papers*, pp. 102–109. Springer, 1985.
- [2] Léon Bottou, “Large-scale machine learning with stochastic gradient descent,” in *Proceedings of COMPSTAT’2010*, pp. 177–186. Springer, 2010.
- [3] M. Schmidt, N. Le Roux, and F. Bach, “Minimizing finite sums with the stochastic average gradient,” *Mathematical Programming*, pp. 1–30, 2013.
- [4] R. Johnson and T. Zhang, “Accelerating stochastic gradient descent using predictive variance reduction,” in *Advances in Neural Information Processing Systems 26*, pp. 315–323. Curran Associates, Inc., 2013.
- [5] A. Defazio, F. Bach, and S. Lacoste-Julien, “Saga: A fast incremental gradient method with support for non-strongly convex composite objectives,” in *Advances in Neural Information Processing Systems*, 2014, pp. 1646–1654.
- [6] Yurii Nesterov, “A method of solving a convex programming problem with convergence rate $O(1/k^2)$,” in *Soviet Mathematics Doklady*, 1983, vol. 27, pp. 372–376.
- [7] Y. Nesterov, “Gradient methods for minimizing composite objective function,” Tech. Rep., UCL, 2007.
- [8] Guanghui Lan and Yi Zhou, “An optimal randomized incremental gradient method,” *arXiv preprint arXiv:1507.02000*, 2015.
- [9] Zeyuan Allen-Zhu, “Katyusha: The first direct acceleration of stochastic gradient methods,” *Journal of Machine Learning Research*, vol. 18, no. 221, pp. 1–51, 2018.
- [10] Tomoya Murata and Taiji Suzuki, “Doubly accelerated stochastic variance reduced dual averaging method for regularized empirical risk minimization,” in *Advances in Neural Information Processing Systems*, 2017, pp. 608–617.
- [11] Davood Karimi and Rabab K Ward, “A hybrid stochastic-deterministic gradient descent algorithm for image reconstruction in cone-beam computed tomography,” *Biomedical Physics & Engineering Express*, vol. 2, no. 1, pp. 015008, 2016.
- [12] Davood Karimi and Rabab K Ward, “Sparse-view image reconstruction in cone-beam computed tomography with variance-reduced stochastic gradient descent and locally-adaptive proximal operation,” *Journal of Medical and Biological Engineering*, vol. 37, no. 3, pp. 420–440, 2017.
- [13] Antonin Chambolle, Matthias J Ehrhardt, Peter Richtárik, and Carola-Bibiane Schonlieb, “Stochastic primal-dual hybrid gradient algorithm with arbitrary sampling and imaging applications,” *SIAM Journal on Optimization*, vol. 28, no. 4, pp. 2783–2808, 2018.
- [14] Oliver Whyte, Josef Sivic, Andrew Zisserman, and Jean Ponce, “Non-uniform deblurring for shaken images,” *International journal of computer vision*, vol. 98, no. 2, pp. 168–186, 2012.
- [15] Kodak-Lossless-True-Color-Image-Suite, “<http://r0k.us/graphics/kodak/>,” .
- [16] A. Beck and M. Teboulle, “Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems,” *IEEE Transactions on Image Processing*, vol. 18, no. 11, pp. 2419–2434, 2009.
- [17] Guanghui Lan, “An optimal method for stochastic composite optimization,” *Mathematical Programming*, vol. 133, no. 1-2, pp. 365–397, 2012.
- [18] Lin Xiao and Tong Zhang, “A proximal stochastic gradient method with progressive variance reduction,” *SIAM Journal on Optimization*, vol. 24, no. 4, pp. 2057–2075, 2014.
- [19] Yurii Nesterov, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer Science & Business Media, 2013.
- [20] Blake E Woodworth and Nati Srebro, “Tight complexity bounds for optimizing composite objectives,” in *Advances in neural information processing systems*, 2016, pp. 3639–3647.
- [21] Hakan Erdogan and Jeffrey A Fessler, “Ordered subsets algorithms for transmission tomography,” *Physics in Medicine & Biology*, vol. 44, no. 11, pp. 2835, 1999.
- [22] Antonin Chambolle and Thomas Pock, “A first-order primal-dual algorithm for convex problems with applications to imaging,” *Journal of mathematical imaging and vision*, vol. 40, no. 1, pp. 120–145, 2011.
- [23] Antonin Chambolle and Thomas Pock, “On the ergodic convergence rates of a first-order primal–dual algorithm,” *Mathematical Programming*, vol. 159, no. 1-2, pp. 253–287, 2016.
- [24] Renbo Zhao and Volkan Cevher, “Stochastic three-composite convex minimization with a linear operator,” in *International Conference on Artificial Intelligence and Statistics*, 2018, pp. 765–774.
- [25] Hua Ouyang, Niao He, Long Tran, and Alexander Gray, “Stochastic alternating direction method of multipliers,” in *International Conference on Machine Learning*, 2013, pp. 80–88.
- [26] Wenliang Zhong and James Kwok, “Accelerated stochastic gradient method for composite regularization,” in *Artificial Intelligence and Statistics*, 2014, pp. 1086–1094.
- [27] Zeyuan Allen-Zhu, “Katyusha x: Practical momentum method for stochastic sum-of-nonconvex optimization,” *arXiv preprint arXiv:1802.03866*, 2018.