
Natasha: Faster Non-Convex Stochastic Optimization via Strongly Non-Convex Parameter

Zeyuan Allen-Zhu¹

Abstract

Given a non-convex function $f(x)$ that is an average of n smooth functions, we design stochastic first-order methods to find its approximate stationary points. The performance of our new methods depend on the smallest (negative) eigenvalue $-\sigma$ of the Hessian. This parameter σ captures how strongly non-convex $f(x)$ is, and is analogous to the strong convexity parameter for convex optimization. At least in theory, our methods outperform known results for a range of parameter σ , and can also be used to find approximate local minima. Our result implies an interesting dichotomy: there exists a threshold σ_0 so that the (currently) fastest methods for $\sigma > \sigma_0$ and for $\sigma < \sigma_0$ have different behaviors: the former scales with $n^{2/3}$ and the latter scales with $n^{3/4}$.

1 Introduction

We study the problem of composite *non-convex* minimization:

$$\min_{x \in \mathbb{R}^d} \left\{ F(x) := \psi(x) + f(x) := \psi(x) + \frac{1}{n} \sum_{i=1}^n f_i(x) \right\} \quad (1.1)$$

where each $f_i(x)$ is *nonconvex but smooth*, and $\psi(\cdot)$ is proper convex, possibly nonsmooth, but relatively simple. We are interested in finding a point x that is an *approximate local minimum* of $F(x)$.

- The finite-sum structure $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ arises prominently in large-scale machine learning tasks. In particular, when minimizing loss over a training set, each example i corresponds to one loss function $f_i(\cdot)$ in the summation. This finite-sum structure allows one to perform stochastic gradient descent with respect to a

Future version of this paper shall be found at <http://arxiv.org/abs/1702.00763>. ¹Microsoft Research. Correspondence to: Zeyuan Allen-Zhu <zeyuan@csail.mit.edu>.

random $\nabla f_i(x)$.

- The so-called *proximal* term $\psi(x)$ adds more generality to the model. For instance, if $\psi(x)$ is the indicator function of a convex set, then problem (1.1) becomes constraint minimization; if $\psi(x) = \|x\|_1$, then we can allow problem (1.1) to perform feature selection. In general, $\psi(x)$ has to be a simple function where the projection operation $\arg \min_x \{ \psi(x) + \frac{1}{2\eta} \|x - x_0\|^2 \}$ is efficiently computable. At a first reading of this paper, one can assume $\psi(x) \equiv 0$ for simplicity.

Many non-convex machine learning problems fall into problem (1.1). Most notably, training *deep neural networks* and classifications with *sigmoid loss* correspond to (1.1) where neither $f_i(x)$ or $f(x)$ is convex. However, our understanding to this challenging non-convex problem is very limited.

1.1 Strongly Non-Convex Optimization

Let L be the smoothness parameter for each $f_i(x)$, meaning all the eigenvalues of $\nabla^2 f_i(x)$ lie in $[-L, L]$.¹

We denote by $\sigma \in [0, L]$ the *strong-nonconvexity* parameter of $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$, meaning that

all the eigenvalues of $\nabla^2 f(x)$ lie in $[-\sigma, L]$.

We emphasize that parameter σ is analogous to the *strong-convexity* parameter μ for convex optimization, where all the eigenvalues of $\nabla^2 f(x)$ lie in $[\mu, L]$ for some $\mu > 0$.

We wish to find an ε -approximate stationary point (a.k.a. critical point) of $F(x)$, that is

a point x satisfying $\|\mathcal{G}(x)\| \leq \varepsilon$

where $\mathcal{G}(x)$ is the so-called gradient mapping of $F(x)$ (see Section 2 for a formal definition). In the special case of $\psi(\cdot) \equiv 0$, gradient mapping $\mathcal{G}(x)$ is the same as gradient $\nabla f(x)$, so x satisfies $\|\nabla f(x)\| \leq \varepsilon$.

Since $f(\cdot)$ is σ -strongly nonconvex, any ε -approximate stationary point is automatically also an (ε, σ) -approximate local minimum — meaning that the Hessian of the output point $\nabla^2 f(x) \succeq -\sigma \mathbf{I}$ is approximately positive semidefinite (PSD).

¹This definition also applies to functions $f(x)$ that are not twice differentiable, see Section 2 for details.

1.2 Motivations and Remarks

- We focus on strongly non-convex optimization because introducing this parameter σ allows us to perform a *more refined study* of non-convex optimization. If σ equals L then L -strongly nonconvex optimization is equivalent to the general non-convex optimization.
- We focus only on finding stationary points as opposed to local minima, because in a recent study — see Appendix A— researchers have shown that finding (ε, δ) -approximate local minima *reduces* to finding ε -approximate stationary points in an $O(\delta)$ -strongly non-convex function.
- Parameter σ is often not constant and can be much smaller than L . For instance, second-order methods often find $(\varepsilon, \sqrt{\varepsilon})$ -approximate local minima (Nesterov, 2008) and this corresponds to $\sigma = \sqrt{\varepsilon}$.

1.3 Known Results

Despite the widespread use of nonconvex models in machine learning and related fields, our understanding to non-convex optimization is still very limited. Until recently, nearly all research papers have been mostly focusing on either $\sigma = 0$ or $\sigma = L$:

- If $\sigma = 0$, the accelerated SVRG method (Shalev-Shwartz, 2016; Allen-Zhu & Yuan, 2016) finds x satisfying $F(x) - F(x^*) \leq \varepsilon$, in gradient complexity $\tilde{O}(n + n^{3/4}\sqrt{L/\varepsilon})$.² This result is irrelevant to this paper because $f(x)$ is simply convex.
- If $\sigma = L$, the SVRG method (Allen-Zhu & Hazan, 2016) finds an ε -approximate stationary point of $F(x)$ in gradient complexity $O(n + n^{2/3}L/\varepsilon^2)$.
- If $\sigma = L$, gradient descent finds an ε -approximate stationary point in gradient complexity $O(nL/\varepsilon^2)$.
- If $\sigma = L$, stochastic gradient descent finds an ε -approx. stationary point in gradient complexity $O(L^2/\varepsilon^4)$.

Throughout this paper, we refer to *gradient complexity* as the total number of stochastic gradient computations $\nabla f_i(x)$ and proximal computations $y \leftarrow \text{Prox}_{\psi, \eta}(x) := \arg \min_y \{\psi(y) + \frac{1}{2\eta}\|y - x\|^2\}$.³

Very recently, it was observed by two independent groups (Agarwal et al., 2017; Carmon et al., 2016) — although implicitly, see Section 2.1— that for solving the σ -strongly nonconvex problem, one can repeatedly regularize $F(x)$ to make it σ -strongly convex, and then apply the accelerated SVRG method to minimize this regularized

²We use \tilde{O} to hide poly-logarithmic factors in $n, L, 1/\varepsilon$.

³Some authors also refer to them as incremental first-order oracle (IFO) and proximal oracle (PO) calls. In most machine learning applications, each IFO and PO call can be implemented to run in time $O(d)$ where d is the dimension of the model, or even in time $O(s)$ if s is the average sparsity of the data vectors.

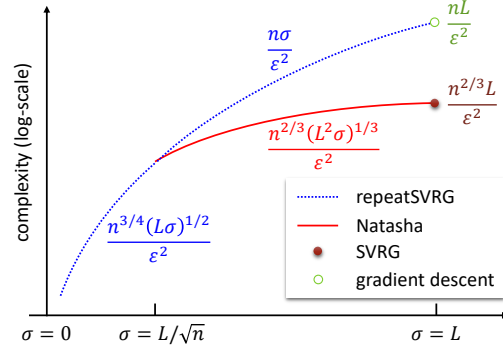


Figure 1: Comparison to prior works

function. Under mild assumption $\sigma \geq \varepsilon^2$, this approach

- finds an ε -approximate stationary point in gradient complexity $\tilde{O}\left(\frac{n\sigma + n^{3/4}\sqrt{L\sigma}}{\varepsilon^2}\right)$.

We call this method repeatSVRG in this paper. Unfortunately, repeatSVRG is even slower than the vanilla SVRG for $\sigma = L$ by a factor $n^{1/3}$, see Figure 1.

Remark on SGD. Stochastic gradient descent (SGD) has a *slower* convergence rate (i.e., in terms of $1/\varepsilon^4$) than other cited first-order methods (i.e., in terms of $1/\varepsilon^2$), see for instance (Ghadimi & Lan, 2015). However, the complexity of SGD does not depend on n and thus is incomparable to gradient descent, SVRG, or repeatSVRG.⁴ This is one of the main motivations to study how to reduce the complexity of non-SGD methods, especially in terms of n .

1.4 Our New Results

In this paper, we identify an interesting dichotomy with respect to the spectrum of the nonconvexity parameter $\sigma \in [0, L]$. In particular, we showed that if $\sigma \geq L/\sqrt{n}$, then our new method Natasha finds an ε -approximate stationary point of $F(x)$ in gradient complexity

$$O\left(n \log \frac{1}{\varepsilon} + \frac{n^{2/3}(L^2\sigma)^{1/3}}{\varepsilon^2}\right).$$

In other words, together with repeatSVRG, we have improved the gradient complexity for σ -strongly nonconvex optimization to⁵

$$\tilde{O}\left(\min\left\{\frac{n^{3/4}\sqrt{L\sigma}}{\varepsilon^2}, \frac{n^{2/3}(L^2\sigma)^{1/3}}{\varepsilon^2}\right\}\right)$$

and the first term in the min is smaller if $\sigma < L/\sqrt{n}$ and the second term is smaller if $\sigma > L/\sqrt{n}$. We illustrate our

⁴In practice, there are examples in non-convex empirical risk minimization (Allen-Zhu & Hazan, 2016) and in training neural networks (Allen-Zhu & Hazan, 2016; Reddi et al., 2016) where SVRG can outperform SGD. Of course, for deep learning tasks, SGD remains to be the best practical method of choice.

⁵We remark here that this is under mild assumptions for ε being sufficiently small. For instance, the result of (Agarwal et al., 2017; Carmon et al., 2016) requires $\varepsilon^2 \leq \sigma$. In our result, the term $n \log \frac{1}{\varepsilon}$ disappears when $\varepsilon^6 \leq L^2\sigma/n$.

performance improvement in Figure 1. Our result matches that of SVRG for $\sigma = L$, and has a much simpler analysis.

Additional Results. One can take a step further and ask what if each function $f_i(x)$ is (ℓ_1, ℓ_2) -smooth for parameters $\ell_1, \ell_2 \geq \sigma$, meaning that all the eigenvalues of $\nabla^2 f_i(x)$ lie in $[-\ell_2, \ell_1]$.

We show that a variant of our method, which we call $\text{Natasha}^{\text{full}}$, solves this more refined problem of (1.1) with total gradient complexity $O(n \log \frac{1}{\epsilon} + \frac{n^{2/3}(\ell_1 \ell_2 \sigma)^{1/3}}{\epsilon^2})$ as long as $\frac{\ell_1 \ell_2}{\sigma^2} \leq n^2$.

Remark 1.1. In applications, ℓ_1 and ℓ_2 can be of very different magnitudes. The most influential example is finding the leading eigenvector of a symmetric matrix. Using the so-called shift-and-invert reduction (Garber et al., 2016), computing the leading eigenvector reduces to the convex version of problem (1.1), where each $f_i(x)$ is $(\lambda, 1)$ -smooth for $\lambda \ll 1$. Other examples include all the applications that are built on shift-and-invert, including high rank SVD/PCA (Allen-Zhu & Li, 2016), canonical component analysis (Allen-Zhu & Li, 2017a), online matrix learning (Allen-Zhu & Li, 2017b), and approximate local minima algorithms (Agarwal et al., 2017; Carmon et al., 2016).

Mini-Batch. Our result generalizes trivially to the mini-batch stochastic setting, where in each iteration one computes $\nabla f_i(x)$ for b random choices of index $i \in [n]$ and average them. The stated gradient complexities of Natasha and $\text{Natasha}^{\text{full}}$ can be adjusted so that the factor $n^{2/3}$ is replaced with $n^{2/3}b^{1/3}$.

1.5 Our Techniques

Let us first recall the main idea behind stochastic variance-reduced methods, such as SVRG (Johnson & Zhang, 2013).

The SVRG method divides iterations into epochs, each of length n . It maintains a snapshot point \tilde{x} for each epoch, and computes the full gradient $\nabla f(\tilde{x})$ only for snapshots. Then, in each iteration t at point x_t , SVRG defines gradient estimator $\tilde{\nabla} = \nabla f_i(x_t) - \nabla f_i(\tilde{x}) + \nabla f(\tilde{x})$ which satisfies $\mathbb{E}_i[\tilde{\nabla}] = \nabla f(x_t)$, and performs proximal update $x_{t+1} \leftarrow \text{Prox}_{\psi, \alpha}(x_t - \alpha \tilde{\nabla})$ for some learning rate α . (Recall that if $\psi(\cdot) \equiv 0$ then we would have $x_{t+1} \leftarrow x_t - \alpha \tilde{\nabla}$.)

In nearly all the aforementioned results for nonconvex optimization, researchers have either directly applied SVRG (Allen-Zhu & Hazan, 2016) (for the case $\sigma = L$), or repeatedly applied SVRG (Agarwal et al., 2017; Carmon et al., 2016) (for general $\sigma \in [0, L]$). This puts some limitation in the algorithmic design, because SVRG requires each epoch to be of length exactly n .⁶

⁶The epoch length of SVRG is always n (or a constant multiple of n in practice), because this ensures the computation of $\tilde{\nabla}$ is of amortized gradient complexity $O(1)$. The per-iteration complexity of SVRG is thus the same as the traditional stochastic

Our New Idea. In this paper, we propose Natasha and $\text{Natasha}^{\text{full}}$, two methods that are no longer black-box reductions to SVRG. Both of them still divide iterations into epochs of length n , and compute gradient estimators $\tilde{\nabla}$ the same way as SVRG. However, we do not apply compute $x_t - \alpha \tilde{\nabla}$ directly.

- In our base algorithm Natasha , we divide each epoch into p sub-epochs, each with a starting vector \hat{x} . Our theory suggests the choice $p \approx (\frac{\sigma^2}{L^2}n)^{1/3}$. Then, we replace the use of $\tilde{\nabla}$ with $\tilde{\nabla} + 2\sigma(x_t - \hat{x})$. This is equivalent to replacing $f(x)$ with its regularized version $f(x) + \sigma\|x - \hat{x}\|^2$, where the center \hat{x} varies across sub-epochs. We provide pseudocode in Algorithm 1 and illustrate it in Figure 2.

We view this additional term $2\sigma(x_t - \hat{x})$ as a type of **retraction**, which stabilizes the algorithm by moving the vector a bit in the backward direction towards \hat{x} .

- In our full algorithm $\text{Natasha}^{\text{full}}$, we add one more ingredient on top of Natasha . That is, we perform updates $z_{t+1} \leftarrow \text{Prox}_{\psi, \alpha}(z_t - \alpha \tilde{\nabla})$ with respect to a different sequence $\{z_t\}$, and then define $x_t = \frac{1}{2}z_t + \frac{1}{2}\hat{x}$ and compute gradient estimators $\tilde{\nabla}$ at points x_t . We provide pseudocode in Algorithm 2 in the appendix.

We view this averaging $x_t = \frac{1}{2}z_t + \frac{1}{2}\hat{x}$ as another type of **retraction**, which stabilizes the algorithm by moving towards \hat{x} . The technique of computing gradients at points x_t but moving a different sequence of points z_t is related to the *Katyusha momentum* recently developed for convex optimization (Allen-Zhu, 2017).

1.6 Other Related Work

Methods based on variance-reduced stochastic gradients were first introduced for convex optimization. The first such method is SAG by Schmidt et al (Schmidt et al., 2013). The two most popular choices for gradient estimators are the SVRG-like one we adopted in this paper (independently introduced by (Johnson & Zhang, 2013; Zhang et al., 2013), and the SAGA-like one introduced by (Defazio et al., 2014). In nearly all applications, the results proven for SVRG-like estimators and SAGA-like estimators are simply exchangeable (therefore, the results of this paper naturally generalize to SAGA-like estimators).

The first “non-convex use” of variance reduction is by Shalev-Shwartz (Shalev-Shwartz, 2016) who assumes that each $f_i(x)$ is non-convex but their average $f(x)$ is still convex. This result has been slightly improved to several more refined settings (Allen-Zhu & Yuan, 2016). The first truly non-convex use of variance reduction (i.e., for $f(x)$ being also non-convex) is independently by (Allen-Zhu & Hazan, 2016) and (Reddi et al., 2016). First-order methods only gradient descent (SGD).

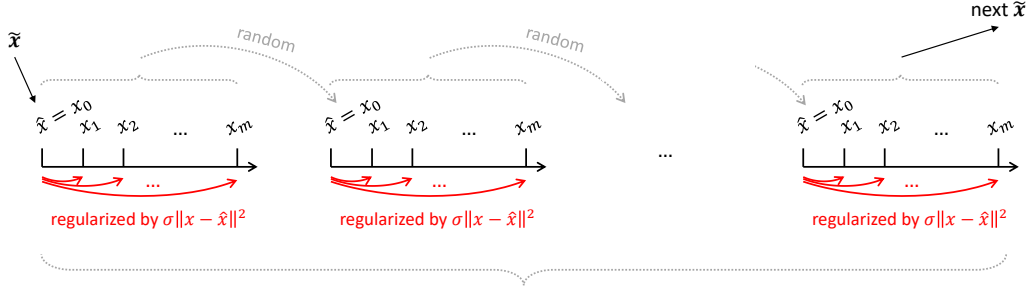


Figure 2: One full epoch of Natasha. The n iterations are divided into p sub-epochs, each consisting of $m = n/p$ steps.

find stationary points (unless there is extra assumption on the randomness of the data), and converge no faster than $1/\varepsilon^2$.

When the second-order Hessian information is used, one can (1) find local minima instead of stationary points, and (2) improve the $1/\varepsilon^2$ rate to $1/\varepsilon^{1.5}$. The first such result is by cubic-regularized Newton’s method (Nesterov, 2008); however, its per-iteration complexity is very high. Very recently, two independent groups of authors tackled this problem from a somewhat similar viewpoint (Carmon et al., 2016; Agarwal et al., 2017): if the computation of Hessian-vector multiplications (i.e., $(\nabla^2 f_i(x))v$) is on the same order of the computation of gradients $\nabla f_i(x)$,⁷ then one can obtain a $(\varepsilon, \sqrt{\varepsilon})$ -approximate local minimum in gradient complexity $\tilde{O}(\frac{n}{\varepsilon^{1.5}} + \frac{n^{3/4}}{\varepsilon^{1.75}})$, if we use big- O to also hide dependencies on the smoothness parameters.

Other related papers include Ge et al. (Ge et al., 2015) where the authors showed that a noise-injected version of SGD converges to local minima instead of critical points, as long as the underlying function is “strict-saddle.” Their theoretical running time is a large polynomial in the dimension. Lee et al. (Lee et al., 2016) showed that gradient descent, starting from a random point, almost surely converges to a local minimum if the function is “strict-saddle”. The rate of convergence required is somewhat unknown.

2 Preliminaries

Throughout this paper, we denote by $\|\cdot\|$ the Euclidean norm. We use $i \in_R [n]$ to denote that i is generated from $[n] = \{1, 2, \dots, n\}$ uniformly at random. We denote by $\nabla f(x)$ the full gradient of function f if it is differentiable, and $\partial f(x)$ any subgradient if f is only Lipschitz continuous at point x . We let x^* be any minimizer of $F(x)$.

Recall some definitions on strong convexity (SC), strongly nonconvexity, and smoothness.

Definition 2.1. For a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$,

⁷A lot of interesting problems satisfy this property, including training neural nets.

- f is σ -strongly convex if $\forall x, y \in \mathbb{R}^d$, it satisfies

$$f(y) \geq f(x) + \langle \partial f(x), y - x \rangle + \frac{\sigma}{2} \|x - y\|^2 .$$
- f is σ -strongly nonconvex if $\forall x, y \in \mathbb{R}^d$, it satisfies

$$f(y) \geq f(x) + \langle \partial f(x), y - x \rangle - \frac{\sigma}{2} \|x - y\|^2 .$$
- f is (ℓ_1, ℓ_2) -smooth if $\forall x, y \in \mathbb{R}^d$, it satisfies

$$\begin{aligned} f(x) + \langle \nabla f(x), y - x \rangle + \frac{\ell_1}{2} \|x - y\|^2 &\geq f(y) \\ &\geq f(x) + \langle \nabla f(x), y - x \rangle - \frac{\ell_2}{2} \|x - y\|^2 . \end{aligned}$$
- f is L -smooth if it is (L, L) -smooth.

The (ℓ_1, ℓ_2) -smoothness parameters were introduced in (Allen-Zhu & Yuan, 2016) to tackle the convex setting of problem (1.1). The notion of strong nonconvexity is also known as “almost convexity (Carmon et al., 2016)” or “lower smoothness (Allen-Zhu & Yuan, 2016).” We refrain from using the name “almost convexity” because it coincides with several other non-equivalent definitions in optimization literatures.

Definition 2.2. Given a parameter $\eta > 0$, the gradient mapping of $F(\cdot)$ in (1.1) at point x is

$$\mathcal{G}_\eta(x) := \frac{1}{\eta} (x - x')$$

where $x' = \arg \min_y \{ \psi(y) + \langle \nabla f(x), y \rangle + \frac{1}{2\eta} \|y - x\|^2 \}$. In particular, if $\psi(\cdot) \equiv 0$, then $\mathcal{G}_\eta(x) \equiv \nabla f(x)$.

The following theorem for the SVRG method can be found for instance in (Allen-Zhu & Yuan, 2016), which is built on top of the results (Shalev-Shwartz, 2016; Lin et al., 2015; Frostig et al., 2015):

Theorem 2.3 (SVRG). Let $G(y) := \psi(y) + \frac{1}{n} \sum_{i=1}^n g_i(y)$ be σ -strongly convex, then the SVRG method finds a point y satisfying $G(y) - G(y^*) \leq \varepsilon$

- with gradient complexity $O((n + \frac{L^2}{\sigma^2}) \log \frac{1}{\varepsilon})$, if each $g_i(\cdot)$ is L -smooth (for $L \geq \sigma$); or
- with gradient complexity $O((n + \frac{\ell_1 \ell_2}{\sigma^2}) \log \frac{1}{\varepsilon})$, if each $g_i(\cdot)$ is (ℓ_1, ℓ_2) -smooth (for $\ell_1, \ell_2 \geq \sigma$).

If one performs acceleration, the running times become $\tilde{O}(n + n^{3/4}\sqrt{L/\sigma})$ and $\tilde{O}(n + n^{3/4}(\ell_1\ell_2\sigma^2)^{1/4})$.

2.1 RepeatSVRG

We recall the idea behind a simple algorithm—that we call repeatSVRG—which finds the ε -approximate stationary points for problem (1.1) when $f(x)$ is σ -strongly nonconvex. The algorithm is divided into stages. In each stage t , consider a modified function $F_t(x) := F(x) + \sigma\|x - x_t\|^2$. It is easy to see that $F_t(x)$ is σ -strongly convex, so one can apply the accelerated SVRG method to minimize $F_t(x)$. Let x_{t+1} be any sufficiently accurate approximate minimizer of $F_t(x)$.⁸

Now, one can prove (c.f. Section 4) that x_{t+1} is an $O(\sigma\|x_t - x_{t+1}\|)$ -approximate stationary point for $F(x)$. Therefore, if $\sigma\|x_t - x_{t+1}\| \leq \varepsilon$ we can stop the algorithm because we have already found an $O(\varepsilon)$ -approximate stationary point. If $\sigma\|x_t - x_{t+1}\| > \varepsilon$, then it must satisfy that $F(x_t) - F(x_{t+1}) \geq \sigma\|x_t - x_{t+1}\|^2 \geq \Omega(\varepsilon^2/\sigma)$, but this cannot happen for more than $T = O(\frac{\sigma}{\varepsilon^2}(F(x_0) - F^*))$ stages. Therefore, the total gradient complexity is T multiplied with the complexity of accelerated SVRG in each stage (which is $\tilde{O}(n + n^{3/4}\sqrt{L/\sigma})$ according to Theorem 2.3).

Remark 2.4. The complexity of repeatSVRG can be inferred from (Agarwal et al., 2017; Carmon et al., 2016), but is not explicitly stated. For instance, the paper (Carmon et al., 2016) does not allow $F(x)$ to have a non-smooth proximal term $\psi(x)$, and applies accelerated gradient descent instead of accelerated SVRG.

3 Our Algorithms

We introduce two variants of our algorithms: (1) the base method *Natasha* targets on the simple regime when $f(x)$ and each $f_i(x)$ are both L -smooth, and (2) the full method *Natasha*^{full} targets on the more refined regime when $f(x)$ is L -smooth but each $f_i(x)$ is (ℓ_1, ℓ_2) -smooth.

Both methods follow the general idea of variance-reduced stochastic gradient descent: in each inner-most iteration, they compute a *gradient estimator* $\tilde{\nabla}$ that is of the form $\tilde{\nabla} = \nabla f(\tilde{x}) - \nabla f_i(\tilde{x}) + \nabla f_i(x)$ and satisfies $\mathbb{E}_{i \in_R[n]}[\tilde{\nabla}] = \nabla f(x)$. Here, \tilde{x} is a *snapshot* point that is changed once every n iterations (i.e., for each different $k = 1, 2, \dots, T'$ in the pseudocode), and we call it a *full epoch* for every distinct k . Notice that the amortized gradient complexity for computing $\tilde{\nabla}$ is $O(1)$ per-iteration.

Base Method. In *Natasha* (see Algorithm 1), as illustrated by Figure 2, we divide each full epoch into p sub-epochs $s = 0, 1, \dots, p - 1$, each of length $m = n/p$. In

⁸Since the accelerated SVRG method has a linear convergence rate for strongly convex functions, the complexity to find such x_{t+1} only depends logarithmically on this accuracy.

each sub-epoch s , we start with a point $x_0 = \hat{x}$, and replace $f(x)$ with its regularized version $f^s(x) := f(x) + \sigma\|x - \hat{x}\|^2$. Then, in each iteration t of the sub-epoch s , we

- compute gradient estimator $\tilde{\nabla}$ with respect to $f^s(x_t)$,
- perform update $x_{t+1} = \arg \min_y \{ \psi(y) + \langle \tilde{\nabla}, y \rangle + \frac{1}{2\alpha}\|y - x_t\|^2 \}$ with learning rate α .

Effectively, the introduction of the regularizer $\sigma\|x - \hat{x}\|^2$ makes sure that when performing update $x_t \leftarrow x_{t+1}$, we also move a bit towards point \hat{x} (i.e., retraction by regularization). Finally, when the sub-epoch is done, we define \hat{x} to be a random one from $\{x_0, \dots, x_{m-1}\}$.

Full Method. In *Natasha*^{full} (see full version), we also divide each full epoch into p sub-epochs. In each sub-epoch s , we start with a point $x_0 = z_0 = \hat{x}$ and define $f^s(x) := f(x) + \sigma\|x - \hat{x}\|^2$. However, this time in each iteration t , we

- compute gradient estimator $\tilde{\nabla}$ with respect to $f^s(x_t)$,
- perform update $z_{t+1} = \arg \min_y \{ \psi(y) + \langle \tilde{\nabla}, y \rangle + \frac{1}{2\alpha}\|y - z_t\|^2 \}$ with learning rate α , and
- choose $x_{t+1} = \frac{1}{2}z_{t+1} + \frac{1}{2}\hat{x}$.

Effectively, the regularizer $\sigma\|x - \hat{x}\|^2$ makes sure that when performing updates, we move a bit towards point \hat{x} (i.e., retraction by regularization); at the same time, the choice $x_{t+1} = \frac{1}{2}z_{t+1} + \frac{1}{2}\hat{x}$ also helps us move towards point \hat{x} (i.e., retraction by the so-called “Katyusha momentum”⁹). Finally, when the sub-epoch is over, we define \hat{x} to be a random one from the set $\{x_0, \dots, x_{m-1}\}$, and move to the next sub-epoch.

4 A Sufficient Stopping Criterion

In this section, we present a sufficient condition for finding approximate stationary points in a σ -strongly nonconvex function. Lemma 4.1 below states that, if we regularize the original function and define $G(x) := F(x) + \sigma\|x - \hat{x}\|^2$ for an arbitrary point \hat{x} , then the minimizer of $G(x)$ is an approximate saddle-point for $F(x)$.

Lemma 4.1. *Suppose $G(y) = F(y) + \sigma\|y - \hat{x}\|^2$ for some given point \hat{x} , and let x^* be the minimizer of $G(y)$. If we minimize $G(y)$ and obtain a point x satisfying*

$$G(x) - G(x^*) \leq \delta^2 \sigma,$$

then for every $\eta \in (0, \frac{1}{\max\{L, 4\sigma\}}]$ we have the gradient mapping

$$\|\mathcal{G}_\eta(x)\|^2 \leq 12\sigma^2\|x^* - \hat{x}\|^2 + O(\delta^2).$$

Notice that when $\psi(x) \equiv 0$ this lemma is trivial, and can be found for instance in (Carmon et al., 2016). The main

⁹The idea for this second kind of retraction, and the idea of having the updates on a sequence z_t but computing gradients at points x_t , is largely motivated by our recent work on the Katyusha momentum and the Katyusha acceleration (Allen-Zhu, 2017).

Algorithm 1 $\text{Natasha}(x^\varnothing, p, T', \alpha)$

Input: starting vector x^\varnothing , sub-epoch count $p \in [n]$, epoch count T' , learning rate $\alpha > 0$.

Output: vector x^{out} .

- 1: $\hat{x} \leftarrow x^\varnothing; m \leftarrow n/p; X \leftarrow [];$
- 2: **for** $k \leftarrow 1$ **to** T' **do** $\diamond T'$ full epochs
- 3: $\tilde{x} \leftarrow \hat{x}; \mu \leftarrow \nabla f(\tilde{x});$
- 4: **for** $s \leftarrow 0$ **to** $p - 1$ **do** $\diamond p$ sub-epochs in each epoch
- 5: $x_0 \leftarrow \tilde{x}; X \leftarrow [X, \tilde{x}];$
- 6: **for** $t \leftarrow 0$ **to** $m - 1$ **do** $\diamond m$ iterations in each sub-epoch
- 7: $i \leftarrow$ a random choice from $\{1, \dots, n\}$.
- 8: $\tilde{\nabla} \leftarrow \nabla f_i(x_t) - \nabla f_i(\tilde{x}) + \mu + 2\sigma(x_t - \tilde{x})$ $\diamond \mathbb{E}_i[\tilde{\nabla}] = \nabla(f(x) + \sigma\|x - \tilde{x}\|^2)|_{x_t}$
- 9: $x_{t+1} = \arg \min_{y \in \mathbb{R}^d} \left\{ \psi(y) + \frac{1}{2\alpha} \|y - x_t\|^2 + \langle \tilde{\nabla}, y \rangle \right\}$
- 10: **end for**
- 11: $\hat{x} \leftarrow$ a random choice from $\{x_0, x_1, \dots, x_{m-1}\};$ \diamond for practitioners, choose the average
- 12: **end for**
- 13: **end for**
- 14: $\hat{x} \leftarrow$ a random vector in $X;$ \diamond for practitioners, choose the last
- 15: $x^{\text{out}} \leftarrow$ an approximate minimizer of $G(y) := F(y) + \sigma\|y - \hat{x}\|^2$ using SVRG.
- 16: **return** $x^{\text{out}}.$ \diamond it suffices to run SVRG for $O(n \log \frac{1}{\epsilon})$ iterations.

technical difficulty arises in order to deal with $\psi(x) \neq 0$. The proof is included in the full version.

5 Base Method: Analysis for One Full Epoch

In this section, we consider problem (1.1) where each $f_i(x)$ is L -smooth and $F(x)$ is σ -strongly nonconvex. We use our base method *Natasha* to minimize $F(x)$, and analyze its behavior for one full epoch in this section. We assume $\sigma \leq L$ without loss of generality, because any L -smooth function is also L -strongly nonconvex.

Notations. We introduce the following notations for analysis purpose only.

- Let \hat{x}^s be the vector \hat{x} at the beginning of sub-epoch s .
- Let x_t^s be the vector x_t in sub-epoch s .
- Let i_t^s be the index $i \in [n]$ in sub-epoch s at iteration t .
- Let $f^s(x) := f(x) + \sigma\|x - \hat{x}^s\|^2$, $F^s(x) := F(x) + \sigma\|x - \hat{x}^s\|^2$, and $x_*^s := \arg \min_x \{F^s(x)\}$.
- Let $\tilde{\nabla} f^s(x_t^s) := \nabla f_i(x_t^s) - \nabla f_i(\tilde{x}) + \nabla f(\tilde{x}) + 2\sigma(x_t - \tilde{x})$ where $i = i_t^s$.
- Let $\tilde{\nabla} f(x_t^s) := \nabla f_i(x_t^s) - \nabla f_i(\tilde{x}) + \nabla f(\tilde{x})$ where $i = i_t^s$.

We obviously have that $f^s(x)$ and $F^s(x)$ are σ -strongly convex, and $f^s(x)$ is $(L + 2\sigma)$ -smooth.

5.1 Variance Upper Bound

The following lemma gives an upper bound on the variance of the gradient estimator $\tilde{\nabla} f^s(x_t^s)$:

Lemma 5.1. We have $\mathbb{E}_{i_t^s} [\|\tilde{\nabla} f^s(x_t^s) - \nabla f^s(x_t^s)\|^2] \leq pL^2\|x_t^s - \hat{x}^s\|^2 + pL^2 \sum_{k=0}^{s-1} \|\hat{x}^k - \hat{x}^{k+1}\|^2$.

Proof. We have

$$\begin{aligned} \mathbb{E}_{i_t^s} [\|\tilde{\nabla} f^s(x_t^s) - \nabla f^s(x_t^s)\|^2] &= \mathbb{E}_{i_t^s} [\|\tilde{\nabla} f(x_t^s) - \nabla f(x_t^s)\|^2] \\ &= \mathbb{E}_{i \in R[n]} [\|(\nabla f_i(x_t^s) - \nabla f_i(\tilde{x})) - (\nabla f(x_t^s) - \nabla f(\tilde{x}))\|^2] \\ &\stackrel{\textcircled{1}}{\leq} \mathbb{E}_{i \in R[n]} [\|\nabla f_i(x_t^s) - \nabla f_i(\tilde{x})\|^2] \\ &\stackrel{\textcircled{2}}{\leq} p \mathbb{E}_{i \in R[n]} [\|\nabla f_i(x_t^s) - \nabla f_i(\hat{x}^s)\|^2] \\ &\quad + p \sum_{k=0}^{s-1} \mathbb{E}_{i \in R[n]} [\|\nabla f_i(\hat{x}^k) - \nabla f_i(\hat{x}^{k+1})\|^2] \\ &\stackrel{\textcircled{3}}{\leq} pL^2\|x_t^s - \hat{x}^s\|^2 + pL^2 \sum_{k=0}^{s-1} \|\hat{x}^k - \hat{x}^{k+1}\|^2. \end{aligned}$$

Above, inequality $\textcircled{1}$ is because for any random vector $\zeta \in \mathbb{R}^d$, it holds that $\mathbb{E}\|\zeta - \mathbb{E}\zeta\|^2 = \mathbb{E}\|\zeta\|^2 - \|\mathbb{E}\zeta\|^2$; inequality $\textcircled{2}$ is because $\hat{x}^0 = \tilde{x}$ and for any p vectors $a_1, a_2, \dots, a_p \in \mathbb{R}^d$, it holds that $\|a_1 + \dots + a_p\|^2 \leq p\|a_1\|^2 + \dots + p\|a_p\|^2$; and inequality $\textcircled{3}$ is because each $f_i(\cdot)$ is L -smooth. \square

5.2 Analysis for One Sub-Epoch

The following inequality is classically known as the ‘‘regret inequality’’ for mirror descent (Allen-Zhu & Orecchia, 2017), and its proof is classical (see full version):

Fact 5.2. $\langle \tilde{\nabla} f^s(x_t^s), x_{t+1}^s - u \rangle + \psi(x_{t+1}^s) - \psi(u) \leq \frac{\|x_t^s - u\|^2}{2\alpha} - \frac{\|x_{t+1}^s - u\|^2}{2\alpha} - \frac{\|x_{t+1}^s - x_t^s\|^2}{2\alpha}$ for every $u \in \mathbb{R}^d$.

The following lemma is our main contribution for the base method *Natasha*.

Lemma 5.3. As long as $\alpha \leq \frac{1}{2L+4\sigma}$, we have

$$\begin{aligned} &\mathbb{E} \left[(F^s(\hat{x}^{s+1}) - F^s(x_*^s)) \right] \\ &\leq \mathbb{E} \left[\frac{F^s(\hat{x}^s) - F^s(x_*^s)}{\sigma\alpha m/2} + \alpha p L^2 \left(\sum_{k=0}^s \|\hat{x}^k - \hat{x}^{k+1}\|^2 \right) \right]. \end{aligned}$$

Proof. We first compute that

$$\begin{aligned}
 F^s(x_{t+1}^s) - F^s(u) &= f^s(x_{t+1}^s) - f^s(u) + \psi(x_{t+1}^s) - \psi(u) \\
 &\stackrel{\textcircled{1}}{\leq} f^s(x_t^s) + \langle \nabla f^s(x_t^s), x_{t+1}^s - x_t^s \rangle + \frac{L+2\sigma}{2} \|x_t^s - x_{t+1}^s\|^2 \\
 &\quad - f^s(u) + \psi(x_{t+1}^s) - \psi(u) \\
 &\stackrel{\textcircled{2}}{\leq} \langle \nabla f^s(x_t^s), x_{t+1}^s - x_t^s \rangle + \frac{L+2\sigma}{2} \|x_t^s - x_{t+1}^s\|^2 \\
 &\quad + \langle \nabla f^s(x_t^s), x_t^s - u \rangle + \psi(x_{t+1}^s) - \psi(u). \tag{5.1}
 \end{aligned}$$

Above, inequality $\textcircled{1}$ uses the fact that $f^s(\cdot)$ is $(L+2\sigma)$ -smooth; and inequality $\textcircled{2}$ uses the convexity of $f^s(\cdot)$. Now, we take expectation with respect to i_t^s on both sides of (5.1), and derive that:

$$\begin{aligned}
 &\mathbb{E}_{i_t^s} [F^s(x_{t+1}^s)] - F^s(u) \\
 &\stackrel{\textcircled{1}}{\leq} \mathbb{E}_{i_t^s} \left[\langle \tilde{\nabla} f^s(x_t^s) - \nabla f^s(x_t^s), x_t^s - x_{t+1}^s \rangle + \langle \tilde{\nabla} f^s(x_t^s), x_{t+1}^s - u \rangle \right. \\
 &\quad \left. + \frac{L+2\sigma}{2} \|x_t^s - x_{t+1}^s\|^2 + \psi(x_{t+1}^s) - \psi(u) \right] \\
 &\stackrel{\textcircled{2}}{\leq} \mathbb{E}_{i_t^s} \left[\langle \tilde{\nabla} f^s(x_t^s) - \nabla f^s(x_t^s), x_t^s - x_{t+1}^s \rangle + \frac{\|x_t^s - u\|^2}{2\alpha} \right. \\
 &\quad \left. - \frac{\|x_{t+1}^s - u\|^2}{2\alpha} - \left(\frac{1}{2\alpha} - \frac{L+2\sigma}{2} \right) \|x_{t+1}^s - x_t^s\|^2 \right] \\
 &\stackrel{\textcircled{3}}{\leq} \mathbb{E}_{i_t^s} \left[\alpha \|\tilde{\nabla} f^s(x_t^s) - \nabla f^s(x_t^s)\|^2 + \frac{\|x_t^s - u\|^2}{2\alpha} - \frac{\|x_{t+1}^s - u\|^2}{2\alpha} \right] \\
 &\stackrel{\textcircled{4}}{\leq} \mathbb{E}_{i_t^s} \left[\alpha p L^2 \|x_t^s - \hat{x}^s\|^2 + \alpha p L^2 \sum_{k=0}^{s-1} \|\hat{x}^k - \hat{x}^{k+1}\|^2 \right. \\
 &\quad \left. + \frac{\|x_t^s - u\|^2}{2\alpha} - \frac{\|x_{t+1}^s - u\|^2}{2\alpha} \right]. \tag{5.2}
 \end{aligned}$$

Above, inequality $\textcircled{1}$ is follows from (5.1) together with the fact that $\mathbb{E}_{i_t^s} [\tilde{\nabla} f^s(x_t^s)] = \nabla f^s(x_t^s)$ implies

$$\begin{aligned}
 &\mathbb{E}_{i_t^s} [\langle \nabla f^s(x_t^s), x_{t+1}^s - x_t^s \rangle + \langle \nabla f^s(x_t^s), x_t^s - u \rangle] \\
 &= \mathbb{E}_{i_t^s} [\langle \tilde{\nabla} f^s(x_t^s) - \nabla f^s(x_t^s), x_{t+1}^s - x_t^s \rangle + \langle \tilde{\nabla} f^s(x_t^s), x_{t+1}^s - u \rangle];
 \end{aligned}$$

inequality $\textcircled{2}$ uses Fact 5.2; inequality $\textcircled{3}$ uses $\alpha \leq \frac{1}{2L+4\sigma}$ together with Young's inequality $\langle a, b \rangle \leq \frac{1}{2} \|a\|^2 + \frac{1}{2} \|b\|^2$; and inequality $\textcircled{4}$ uses Lemma 5.1.

Finally, choosing $u = x_*^s$ to be the (unique) minimizer of $F^s(\cdot) = f^s(\cdot) + \psi(\cdot)$, and telescoping inequality (5.2) for $t = 0, 1, \dots, m-1$, we have

$$\begin{aligned}
 &\mathbb{E} \left[\sum_{t=1}^{m-1} (F^s(x_t^s) - F^s(x_*^s)) \right] \\
 &\leq \mathbb{E} \left[\frac{\|x_0^s - x_*^s\|^2}{2\alpha} + \sum_{t=0}^{m-1} \left(\alpha p L^2 \|x_t^s - \hat{x}^s\|^2 \right. \right. \\
 &\quad \left. \left. + \alpha p L^2 \sum_{k=0}^{s-1} \|\hat{x}^k - \hat{x}^{k+1}\|^2 \right) \right] \\
 &\leq \mathbb{E} \left[\frac{F^s(\hat{x}^s) - F^s(x_*^s)}{\sigma\alpha} + \alpha p m L^2 \left(\sum_{k=0}^s \|\hat{x}^k - \hat{x}^{k+1}\|^2 \right) \right].
 \end{aligned}$$

Above, the second inequality uses the fact that \hat{x}^{s+1} is chosen from $\{x_0^s, \dots, x_{m-1}^s\}$ uniformly at random, as well as the σ -strong convexity of $F^s(\cdot)$.

Dividing both sides by m and rearranging the terms (using $\frac{1}{2\sigma\alpha} \geq 1$), we have

$$\begin{aligned}
 &\mathbb{E} \left[(F^s(\hat{x}^{s+1}) - F^s(x_*^s)) \right] \\
 &\leq \mathbb{E} \left[\frac{F^s(\hat{x}^s) - F^s(x_*^s)}{\sigma\alpha m/2} + \alpha p L^2 \left(\sum_{k=0}^s \|\hat{x}^k - \hat{x}^{k+1}\|^2 \right) \right]. \quad \square
 \end{aligned}$$

5.3 Analysis for One Full Epoch

One can telescope Lemma 5.3 for an entire epoch and arrive at the following lemma (see full version):

Lemma 5.4. *If $\alpha \leq \frac{1}{2L+4\sigma}$, $\alpha \geq \frac{4}{\sigma m}$ and $\alpha \leq \frac{\sigma}{p^2 L^2}$, we have*

$$\sum_{s=0}^{p-1} \mathbb{E} \left[(F^s(\hat{x}^s) - F^s(x_*^s)) \right] \leq 2\mathbb{E} \left[F(\hat{x}^0) - F(\hat{x}^p) \right].$$

6 Base Method: Final Theorem

We are now ready to state and prove our main convergence theorem for Natasha:

Theorem 1. *Suppose in (1.1), each $f_i(x)$ is L -smooth and $F(x)$ is σ -strongly nonconvex for $\sigma \leq L$. Then, if $\frac{L^2}{\sigma^2} \leq n$, $p = \Theta\left(\left(\frac{\sigma^2}{L^2}n\right)^{1/3}\right)$ and $\alpha = \Theta\left(\frac{\sigma}{p^2 L^2}\right)$, our base method Natasha outputs a point x^{out} satisfying*

$$\mathbb{E}[\|\mathcal{G}_\eta(x^{\text{out}})\|^2] \leq O\left(\frac{(L^2\sigma)^{1/3}n^{2/3}}{T'n}\right) \cdot (F(x^\varnothing) - F^*).$$

for every $\eta \in \left(0, \frac{1}{\max\{L, 4\sigma\}}\right]$. In other words, to obtain $\mathbb{E}[\|\mathcal{G}_\eta(x^{\text{out}})\|^2] \leq \varepsilon^2$, we need gradient complexity

$$O\left(n \log \frac{1}{\varepsilon} + \frac{(L^2\sigma)^{1/3}n^{2/3}}{\varepsilon^2} \cdot (F(x^\varnothing) - F^*)\right).$$

In the above theorem, we have assumed $\sigma \leq L$ without loss of generality because any L -smooth function is also L -strongly nonconvex. Also, we have assumed $\frac{L^2}{\sigma^2} \leq n$ and if this inequality does not hold, then one should apply repeatSVRG for a faster running time (see Figure 1).

Proof of Theorem 1. We choose $p = \left(\frac{\sigma^2}{24L^2}n\right)^{1/3}$, $m = n/p$, and $\alpha = \frac{4}{\sigma m} = \frac{\sigma}{6p^2 L^2} \leq \frac{1}{2L+4\sigma}$, so we can apply Lemma 5.4. If we telescope Lemma 5.4 for the entire algorithm (which has T' full epochs), and use the fact that \hat{x}^p of the previous epoch equals \hat{x}^0 of the next epoch, we conclude that if we choose a random epoch and a random subepoch s , we will have

$$\mathbb{E}[F^s(\hat{x}^s) - F^s(x_*^s)] \leq \frac{2}{pT'} (F(x^\varnothing) - F^*).$$

By the σ -strong convexity of $F^s(\cdot)$, we have $\mathbb{E}[\sigma \|\hat{x}^s - x_*^s\|^2] \leq \frac{4}{pT'} (F(x^\varnothing) - F^*)$.

Now, $F^s(x) = F(x) + \sigma \|x - \hat{x}^s\|^2$ satisfies the assumption of $G(x)$ in Lemma 4.1. If we use the SVRG method (see Theorem 2.3) to minimize the convex function $F^s(x)$, we

get an output x^{out} satisfying $F^s(x^{\text{out}}) - F^s(x_*^s) \leq \varepsilon^2 \sigma$ in gradient complexity $O\left((n + \frac{L^2}{\sigma^2}) \log \frac{1}{\varepsilon}\right) \leq O(n \log \frac{1}{\varepsilon})$.

We can therefore apply Lemma 4.1 and conclude that this output x^{out} satisfies

$$\begin{aligned} \mathbb{E}[\|\mathcal{G}_\eta(x^{\text{out}})\|^2] &\leq O\left(\frac{\sigma}{pT'}\right) \cdot (F(x^\varnothing) - F^*) \\ &= O\left(\frac{(L^2\sigma)^{1/3}n^{2/3}}{T'n}\right) \cdot (F(x^\varnothing) - F^*) . \end{aligned}$$

In other words, we obtain $\mathbb{E}[\|\mathcal{G}_\eta(x^{\text{out}})\|^2] \leq \varepsilon^2$ using

$$T'n = O\left(n + \frac{(L^2\sigma)^{1/3}n^{2/3}}{\varepsilon^2} \cdot (F(x^\varnothing) - F^*)\right)$$

computations of the stochastic gradients. Here, the additive term n is because $T' \geq 1$.

Finally, adding this with $O(n \log \frac{1}{\varepsilon})$, the gradient complexity for the application of SVRG in the last line of `Natasha`, we finish the proof of the total gradient complexity. \square

7 Full Method: Final Theorem

We analyze and state the main theorems for our full method `Natashafull` in the full version of this paper.

8 Conclusion

Stochastic gradient descent and gradient descent (including alternating minimization) have become the canonical methods for solving non-convex machine learning tasks. However, can we design new non-convex methods to run even faster than SGD or GD?

This present paper tries to tackle this general question, by providing a new `Natasha` method which is intrinsically different from GD or SGD. It runs faster than GD and SVRG-based methods at least in theory. We hope that this could be a non-negligible step towards our better understanding of non-convex optimization.

Finally, our results give rise to an interesting dichotomy in the best-known complexity of first-order non-convex optimization: the complexity scales with $n^{3/4}$ for $\sigma < L/\sqrt{n}$ and with $n^{2/3}$ for $\sigma > L/\sqrt{n}$. It remains open to investigate whether this dichotomy is intrinsic, or we can design a more efficient algorithm that outperforms both.

References

Agarwal, Naman, Allen-Zhu, Zeyuan, Bullins, Brian, Hazan, Elad, and Ma, Tengyu. Finding Approximate Local Minima for Nonconvex Optimization in Linear Time. In *STOC*, 2017.

Allen-Zhu, Zeyuan. Katyusha: The First Direct Acceleration of Stochastic Gradient Methods. In *STOC*, 2017.

Allen-Zhu, Zeyuan and Hazan, Elad. Variance Reduction for Faster Non-Convex Optimization. In *NIPS*, 2016.

Allen-Zhu, Zeyuan and Li, Yuanzhi. LazySVD: Even Faster SVD Decomposition Yet Without Agonizing Pain. In *NIPS*, 2016.

Allen-Zhu, Zeyuan and Li, Yuanzhi. Doubly Accelerated Methods for Faster CCA and Generalized Eigendecomposition. In *Proceedings of the 34th International Conference on Machine Learning*, ICML '17, 2017a.

Allen-Zhu, Zeyuan and Li, Yuanzhi. Follow the Compressed Leader: Faster Online Learning of Eigenvectors and Faster MMWU. In *Proceedings of the 34th International Conference on Machine Learning*, ICML '17, 2017b.

Allen-Zhu, Zeyuan and Orecchia, Lorenzo. Linear Coupling: An Ultimate Unification of Gradient and Mirror Descent. In *Proceedings of the 8th Innovations in Theoretical Computer Science*, ITCS '17, 2017.

Allen-Zhu, Zeyuan and Yuan, Yang. Improved SVRG for Non-Strongly-Convex or Sum-of-Non-Convex Objectives. In *ICML*, 2016.

Carmon, Yair, Duchi, John C., Hinder, Oliver, and Sidford, Aaron. Accelerated Methods for Non-Convex Optimization. *ArXiv e-prints*, abs/1611.00756, November 2016.

Defazio, Aaron, Bach, Francis, and Lacoste-Julien, Simon. SAGA: A Fast Incremental Gradient Method With Support for Non-Strongly Convex Composite Objectives. In *NIPS*, 2014. URL <http://arxiv.org/abs/1407.0202>.

Frostig, Roy, Ge, Rong, Kakade, Sham M., and Sidford, Aaron. Un-regularizing: approximate proximal point and faster stochastic algorithms for empirical risk minimization. In *ICML*, volume 37, pp. 1–28, 2015. URL <http://arxiv.org/abs/1506.07512>.

Garber, Dan, Hazan, Elad, Jin, Chi, Kakade, Sham M., Musco, Cameron, Netrapalli, Praneeth, and Sidford, Aaron. Robust shift-and-invert preconditioning: Faster and more sample efficient algorithms for eigenvector computation. In *ICML*, 2016.

Ge, Rong, Huang, Furong, Jin, Chi, and Yuan, Yang. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Proceedings of the 28th Annual Conference on Learning Theory*, COLT 2015, 2015.

Ghadimi, Saeed and Lan, Guanghui. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, pp. 1–26, feb 2015. ISSN 0025-5610. doi: 10.1007/s10107-015-0871-8. URL <http://arxiv.org/abs/1310.3787><http://link.springer.com/10.1007/s10107-015-0871-8>.

- Johnson, Rie and Zhang, Tong. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, NIPS 2013, pp. 315–323, 2013.
- Lee, Jason D., Simchowitz, Max, Jordan, Michael I., and Recht, Benjamin. Gradient descent only converges to minimizers. In *Proceedings of the 29th Conference on Learning Theory, COLT 2016, New York, USA, June 23-26, 2016*, pp. 1246–1257, 2016.
- Lin, Hongzhou, Mairal, Julien, and Harchaoui, Zaid. A Universal Catalyst for First-Order Optimization. In *NIPS*, 2015. URL <http://arxiv.org/pdf/1506.02186v1.pdf>.
- Nesterov, Yurii. Accelerating the cubic regularization of newton’s method on convex problems. *Mathematical Programming*, 112(1):159–181, 2008.
- Reddi, Sashank J., Hefny, Ahmed, Sra, Suvrit, Póczos, Barnabas, and Smola, Alex. Stochastic variance reduction for nonconvex optimization. *ArXiv e-prints*, abs/1603.06160, March 2016.
- Schmidt, Mark, Le Roux, Nicolas, and Bach, Francis. Minimizing finite sums with the stochastic average gradient. *arXiv preprint arXiv:1309.2388*, pp. 1–45, 2013. URL <http://arxiv.org/abs/1309.2388>. Preliminary version appeared in NIPS 2012.
- Shalev-Shwartz, Shai. SDCA without Duality, Regularization, and Individual Convexity. In *ICML*, 2016.
- Zhang, Lijun, Mahdavi, Mehrdad, and Jin, Rong. Linear convergence with condition number independent access of full gradients. In *Advances in Neural Information Processing Systems*, pp. 980–988, 2013.