

Don't Jump Through Hoops and Remove Those Loops: SVRG and Katyusha are Better Without the Outer Loop

Dmitry Kovalev

King Abdullah University of Science and Technology, Saudi Arabia

DMITRY.KOVALEV@KAUST.EDU.SA

Samuel Horváth

King Abdullah University of Science and Technology, Saudi Arabia

SAMUEL.HORVATH@KAUST.EDU.SA

Peter Richtárik

King Abdullah University of Science and Technology, Saudi Arabia

PETER.RICHTARIK@KAUST.EDU.SA

Editors: Aryeh Kontorovich and Gergely Neu

Abstract

The stochastic variance-reduced gradient method (SVRG) and its accelerated variant (Katyusha) have attracted enormous attention in the machine learning community in the last few years due to their superior theoretical properties and empirical behaviour on training supervised machine learning models via the empirical risk minimization paradigm. A key structural element in both of these methods is the inclusion of an outer loop at the beginning of which a full pass over the training data is made in order to compute the exact gradient, which is then used in an inner loop to construct a variance-reduced estimator of the gradient using new stochastic gradient information. In this work, we design *loopless variants* of both of these methods. In particular, we remove the outer loop and replace its function by a coin flip performed in each iteration designed to trigger, with a small probability, the computation of the gradient. We prove that the new methods enjoy the same superior theoretical convergence properties as the original methods. For loopless SVRG, the same rate is obtained for a large interval of coin flip probabilities, including the probability $1/n$, where n is the number of functions. This is the first result where a variant of SVRG is shown to converge with the same rate without the need for the algorithm to know the condition number, which is often unknown or hard to estimate correctly. We demonstrate through numerical experiments that the loopless methods can have superior and more robust practical behavior.

Keywords: stochastic optimization, variance-reduced methods, SVRG

1. Introduction

Empirical risk minimization (a.k.a. finite-sum) problems form the dominant paradigm for training supervised machine learning models such as ridge regression, support vector machines, logistic regression, and neural networks. In its most general form, a finite sum problem has the form

$$\min_{x \in \mathbb{R}^d} f(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(x), \quad (1)$$

where n refers to the number of training data points (e.g., videos, images, molecules), x is the vector representation of a model using d features, and $f_i(x)$ is the loss of model x on data point i .

Variance-reduced methods. One of the most remarkable algorithmic breakthroughs in recent years was the development of *variance-reduced* stochastic gradient algorithms for solving (1). These

methods are significantly faster than SGD (Nemirovsky and Yudin, 1983; Nemirovski et al., 2009; Takáč et al., 2013) in theory and practice on convex and strongly convex problems, and faster in theory on several classes on nonconvex problems (unfortunately, these methods are not yet successful in training production-grade neural networks).

Two of the most notable and popular methods belonging to the family of variance-reduced methods are SVRG (Johnson and Zhang, 2013) and its accelerated variant known as Katyusha (Allen-Zhu, 2017). The latter method accelerates the former via the employment of a novel “negative momentum” idea. Both of these methods have a double loop design. At the beginning of the outer loop, a full pass over the training data is made to compute the gradient of f at a reference point w^k , which is chosen as the freshest iterate (SVRG) or a weighted average of recent iterates (for Katyusha). This gradient is then used in the inner loop to *adjust* the stochastic gradient $\nabla f_i(x^k)$, where i is sampled uniformly at random from $[n] \stackrel{\text{def}}{=} \{1, 2, \dots, n\}$, and x^k is the current iterate, so as to reduce its variance. In particular, both SVRG and Katyusha perform the adjustment $g^k = \nabla f_i(x^k) - \nabla f_i(w^k) + \nabla f(w^k)$. Note that, like $\nabla f_i(x^k)$, the new search direction g^k is an unbiased estimator of $\nabla f(x^k)$. Indeed,

$$\mathbb{E} [g^k] = \nabla f(x^k) - \nabla f(w^k) + \nabla f(w^k) = \nabla f(x^k). \quad (2)$$

where the expectation is taken over random choice of $i \in [n]$. However, it turns out that as the methods progress, the variance of g^k , unlike that of $\nabla f_i(x^k)$, progressively decreases to zero. The total effect of this is significantly faster convergence.

Convergence of SVRG and Katyusha for L -smooth and μ -strongly convex functions. For instance, consider the regime where f_i is L -smooth for each i , and f is μ -strongly convex:

Assumption 1 (L -smoothness) Functions $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ are L -smooth for some $L > 0$:

$$f_i(y) \leq f_i(x) + \langle \nabla f_i(x), y - x \rangle + \frac{L}{2} \|y - x\|^2, \quad \forall x, y \in \mathbb{R}^d. \quad (3)$$

Assumption 2 (μ -strong convexity) Function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex for $\mu > 0$:

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2, \quad \forall x, y \in \mathbb{R}^d. \quad (4)$$

In this regime, the iteration complexity of SVRG is $\mathcal{O}((n + L/\mu) \log 1/\epsilon)$, which is a vast improvement on the linear rate of gradient descent (GD), which is $\mathcal{O}(nL/\mu \log 1/\epsilon)$, and on the sublinear rate of SGD, which is $\mathcal{O}(L/\mu + \sigma^2/\mu^2\epsilon)$, where $\sigma^2 = 1/n \sum_i \|\nabla f_i(x^*)\|^2$ and x^* is the (necessarily unique) minimizer of f . On the other hand, Katyusha enjoys the *accelerated* rate $\mathcal{O}((n + \sqrt{nL/\mu}) \log 1/\epsilon)$, which is superior to that of SVRG in the ill-conditioned regime where $L/\mu \geq n$. This rate has been shown to be *optimal* in a certain precise sense (Nesterov, 2013).

In the past several years, an enormous effort of the machine learning and optimization communities was exerted into designing new efficient variance-reduced methods to tackle problem (1). These developments have brought about a renaissance in the field. The historically first provably variance-reduced method, the stochastic average gradient (SAG) method of Roux et al. (2012); Schmidt et al. (2017), was awarded the Lagrange prize in continuous optimization in 2018. The SAG method was later modified to an unbiased variant called SAGA (Defazio et al., 2014a), achieving the same theoretical rates. Alternative variance-reduced methods include MISO (Mairal, 2015),

FINITO (Defazio et al., 2014b), SDCA (Shalev-Shwartz, 2016), dfSDCA (Csiba and Richtárik, 2015), AdaSDCA (Csiba et al., 2015), QUARTZ (Qu et al., 2015), SBFGS (Gower et al., 2016), SDNA (Qu et al., 2016), SARAH (Nguyen et al., 2017) and S2GD (Konečný and Richtárik, 2017), mS2GD (Konečný et al., 2016), RBCN (Doikov and Richtárik, 2018), JacSketch (Gower et al., 2018) and SAGD (Bibi et al., 2018). Accelerated variance-reduced methods were developed by Shalev-Shwartz and Zhang (2014), Defazio (2016), Zhou (2018) and Zhou et al. (2018).

2. Contributions

As explained in the introduction, a trade-mark structural feature of SVRG and its accelerated variant, Katyusha, is the presence of the outer loop in which a full pass over the data is made. However, the presence of this outer loop is the source of several issues. First, the methods are harder to analyze. Second, one needs to decide at which point the inner loop is terminated and the outer loop entered. For SVRG, the theoretically optimal inner loop size depends on both L and μ . However, μ is not always known. Moreover, even when an estimate is available, as is the case in regularized problems with an explicit strongly convex regularizer, the estimate can often be very loose. Because of these issues, one often chooses the inner loop size in a suboptimal way, such as by setting it to n or $\mathcal{O}(n)$.

Two loopless methods. In this paper we address the above issues by developing *loopless* variants of both SVRG and Katyusha; we refer to them as L-SVRG and L-Katyusha, respectively. In these methods, we dispose of the outer loop and replace its role by a *biased coin-flip*, to be performed in every step of the methods, used to trigger the computation of the gradient $\nabla f(w^k)$ via a pass over the data. In particular, in each step, with (a small) probability $p > 0$ we perform a full pass over data and update the reference gradient $\nabla f(w^k)$. With probability $1 - p$ we keep the previous reference gradient. This procedure can alternatively be interpreted as *having an outer loop of a random length*. However, the resulting methods are easier to write down, comprehend and analyze.

Fast rates are preserved. We show that L-SVRG and L-Katyusha enjoy the same fast theoretical rates as their loopy forefathers. Our proofs are different and the complexity results more insightful.

For L-SVRG with fixed stepsize $\eta = 1/6L$ and probability $p = 1/n$, we show (see Theorem 5) that for the Lyapunov function

$$\Phi^k \stackrel{\text{def}}{=} \left\| x^k - x^* \right\|^2 + \frac{4\eta^2}{pn} \sum_{i=1}^n \left\| \nabla f_i(w^k) - \nabla f_i(x^*) \right\|^2. \quad (5)$$

we get $\mathbb{E}[\Phi^k] \leq \epsilon \Phi^0$ as long as $k = \mathcal{O}((n + L/\mu) \log 1/\epsilon)$. In contrast, the classical SVRG result shows convergence of the expected functional suboptimality $\mathbb{E}[f(x^k) - f(x^*)]$ to zero at the same rate. Note that the classical result follows from our theorem by utilizing the inequality $f(x^k) - f(x^*) \leq L/2 \|x^k - x^*\|^2$, which is a simple consequence of L -smoothness. However, our result provides a deeper insight into the behavior of the method. In particular, it follows that the gradients $\nabla f_i(w^k)$ at the reference points w^k converge to the gradients at the optimum. This is a key intuition behind the workings of SVRG, one not revealed by the classical analysis. Hereby we close the gap in the theoretical understanding of the the SVRG convergence mechanism. Moreover, our theory predicts that as long as p is chosen in the (possibly very large) interval

$$\min \left\{ \frac{c}{n}, \frac{c\mu}{L} \right\} \leq p \leq \max \left\{ \frac{c}{n}, \frac{c\mu}{L} \right\}, \quad (6)$$

where $c = \Theta(1)$, L-SVRG will enjoy the optimal complexity $\mathcal{O}((n + L/\mu) \log 1/\epsilon)$. In the ill-conditioned regime $L/\mu \gg n$, for instance, we roughly have $p \in [\mu/L, 1/n]$. This is in contrast with the (loopy/standard) SVRG method the outer loop of which needs to be of the size $\approx L/\mu$. To the best of our knowledge, SVRG does not enjoy this rate for an outer loop of size n (or any value independent of μ , which is often not known in practice), even though this is the setting most often used in practice. Several authors have tried to establish such a result, but without success. We thus answer an open problem since 2013, the inception of SVRG.

For L-Katyusha with stepsize $\eta = \frac{\theta_2}{(1+\theta_2)\theta_1}$ we show convergence of the Lyapunov function

$$\Psi^k = \mathcal{Z}^k + \mathcal{Y}^k + \mathcal{W}^k, \quad (7)$$

where $\mathcal{Z}^k = \frac{L(1+\eta\sigma)}{2\eta} \|z^k - x^*\|^2$, $\mathcal{Y}^k = \frac{1}{\theta_1}(f(y^k) - f(x^*))$, and $\mathcal{W}^k = \frac{\theta_2(1+\theta_1)}{p\theta_1}(f(w^k) - f(x^*))$, and where x^k, y^k and w^k are iterates produced by the method, with the parameters defined by $\sigma = \mu/L$, $\theta_1 = \min\{\sqrt{2\sigma n/3}, 1/2\}$, $\theta_2 = 1/2$, $p = 1/n$. Our main result (Theorem 11) states that $\mathbb{E}[\Psi^k] \leq \epsilon\Psi^0$ as long as $k = \mathcal{O}((n + \sqrt{nL/\mu}) \log 1/\epsilon)$.

Simplified analysis. Advantage of the loopless approach is that a *single iteration analysis is sufficient to establish convergence*. In contrast, one needs to perform elaborate aggregation across the inner loop to prove the convergence of the original loopy methods.

Superior empirical behaviour. We show through extensive numerical testing on both synthetic and real data that our loopless methods are superior to their loopy variants. We show through experiments that L-SVRG is *very robust to the choice of p from the optimal interval (6)* predicted by our theory. Moreover, *even the worst case for L-SVRG outperforms the best case for SVRG*. This shows how further randomization can significantly speed up and stabilize the algorithm.

Notation. Throughout the whole paper we use conditional expectation $\mathbb{E}[\mathcal{X} \mid x^k, w^k]$ for L-SVRG and $\mathbb{E}[\mathcal{X} \mid y^k, z^k, w^k]$ for L-Katyusha, but for simplicity we will denote these expectations as $\mathbb{E}[\mathcal{X}]$. If $\mathbb{E}[\mathcal{X}]$ refers to unconditional expectation, it is directly mentioned.

3. Loopless SVRG (L-SVRG)

In this section we describe in detail the Loopless SVRG method (L-SVRG), and its convergence.

The algorithm. The L-SVRG method, formalized as Algorithm 1, is inspired by the original SVRG (Johnson and Zhang, 2013) method. We remove the outer loop present in SVRG and instead use a probabilistic update of the full gradient.¹ This update can be also seen in a way that outer loop size is generated by geometric distribution similar to methods of Konečný and Richtárik (2017); Lei et al. (2017).

Note that the reference point w^k (at which a full gradient is computed) is updated in each iteration with probability p to the current iterate x^k , and is left unchanged with probability $1 - p$. Alternatively, the probability p can be seen as a parameter that controls the expected time before next full pass over data. To be more precise, the expected time before next full pass over data is $1/p$. Intuitively, we wish to keep p small so that full passes over data are computed rarely enough. As we shall see next, the simple choice $p = 1/n$ leads to complexity identical to that of original SVRG.

1. This idea was independently explored in Hofmann et al. (2015); we have learned about this work after a first draft of our paper was finished.

Algorithm 1 Loopless SVRG (L-SVRG)

Parameters: stepsize $\eta > 0$, probability $p \in (0, 1]$

Initialization: $x^0 = w^0 \in \mathbb{R}^d$

for $k = 0, 1, 2, \dots$ **do**

$$g^k = \nabla f_i(x^k) - \nabla f_i(w^k) + \nabla f(w^k) \quad (i \in \{1, \dots, n\} \text{ is sampled uniformly at random})$$

$$x^{k+1} = x^k - \eta g^k$$

$$w^{k+1} = \begin{cases} x^k & \text{with probability } p \\ w^k & \text{with probability } 1 - p \end{cases}$$

end for

Convergence theory. A key role in the analysis is played by the *gradient learning* quantity

$$\mathcal{D}^k \stackrel{\text{def}}{=} \frac{4\eta^2}{pn} \sum_{i=1}^n \left\| \nabla f_i(w^k) - \nabla f_i(x^*) \right\|^2 \quad (8)$$

and the Lyapunov function $\Phi^k \stackrel{\text{def}}{=} \|x^k - x^*\|^2 + \mathcal{D}^k$. The analysis involves four lemmas, followed by the main theorem. We wish to mention the lemmas as they highlight the way in which the argument works. All lemmas combined, together with the main theorem, can be proved on a single page, which underlines the simplicity of our approach.

Our first lemma upper bounds the expected squared distance of x^{k+1} from x^* in terms of the same distance but for x^k , function suboptimality, and second moment of g^k .

Lemma 1 *We have*

$$\mathbb{E} \left[\|x^{k+1} - x^*\|^2 \right] \leq (1 - \eta\mu) \|x^k - x^*\|^2 - 2\eta(f(x^k) - f(x^*)) + \eta^2 \mathbb{E} \left[\|g^k\|^2 \right]. \quad (9)$$

In our next lemma, we further bound the second moment of g^k in terms of function suboptimality and \mathcal{D}^k .

Lemma 2 *We have*

$$\mathbb{E} \left[\|g^k\|^2 \right] \leq 4L(f(x^k) - f(x^*)) + \frac{p}{2\eta^2} \mathcal{D}^k. \quad (10)$$

Finally, we bound $\mathbb{E} [\mathcal{D}^{k+1}]$ in terms of \mathcal{D}^k and function suboptimality.

Lemma 3 *We have*

$$\mathbb{E} [\mathcal{D}^{k+1}] \leq (1 - p)\mathcal{D}^k + 8L\eta^2(f(x^k) - f(x^*)). \quad (11)$$

Putting the above three lemmas together naturally leads to the following result involving Lyapunov function (5).

Lemma 4 *Let the step size $\eta \leq 1/6L$. Then for all $k \geq 0$ the following inequality holds:*

$$\mathbb{E} [\Phi^{k+1}] \leq (1 - \eta\mu) \|x^k - x^*\|^2 + \left(1 - \frac{p}{2}\right) \mathcal{D}^k. \quad (12)$$

Using this lemma we can obtain a recursion involving the Lyapunov function on the right-hand side of (12) and obtain the rate of convergence stated in the following theorem.

Theorem 5 *Let $\eta = 1/6L$, $p = 1/n$. Then $\mathbb{E} [\Phi^k] \leq \varepsilon \Phi^0$ as long as $k \geq \mathcal{O}((n + L/\mu) \log 1/\varepsilon)$.*

Proof As the corollary of Lemma 4 we have $\mathbb{E} [\Phi^k] \leq \max\{1 - \eta\mu, 1 - p/2\} \Phi^{k-1}$. Setting $\eta = 1/6L$, $p = 1/n$ and unrolling conditional probability one obtains

$$\mathbb{E} [\Phi^k] \leq \max\left\{1 - \frac{\mu}{6L}, 1 - \frac{1}{2n}\right\}^k \Phi^0,$$

which concludes the proof. ■

Note that the step size does not depend on the strong convexity parameter μ and yet the resulting complexity adapts to it.

Discussion. Examining (12), we can see that contraction of the Lyapunov function is $\max\{1 - \eta\mu, 1 - p/2\}$. Due to the limitation of $\eta \leq 1/6L$, the first term is at least $1 - \eta/6\mu$, thus the complexity cannot be better than $\mathcal{O}(L/\mu \log 1/\varepsilon)$. In terms of total complexity (number of stochastic gradient calls), L-SVRG calls the stochastic gradient oracle in expectation $\mathcal{O}(1 + pn)$ times in each iteration. Combining these two complexities together, one gets the total complexity $\mathcal{O}((1/p + n + L/\mu + Lpn/\mu) \log 1/\varepsilon)$. Note that any choice of $p \in [\min\{c/n, c\mu/L\}, \max\{c/n, c\mu/L\}]$, where $c = \Theta(1)$, leads to the optimal total complexity $\mathcal{O}((n + L/\mu) \log 1/\varepsilon)$. This fills the gap in SVRG theory, where the outer loop length (in our case $1/p$ in expectation) needs to be proportional to L/μ . Moreover, analysis for L-SVRG is much simpler and provides more insights.

4. Loopless Katyusha (L-Katyusha)

In this section we describe in detail the Loopless Katyusha method (L-Katyusha), and its convergence properties.

The algorithm. The L-Katyusha method, formalized as Algorithm 2, is inspired by the original Katyusha (Allen-Zhu, 2017) method. We use the same technique as for Algorithm 1, where we remove the outer loop present in Katyusha and instead use a probabilistic update of the full gradient.

Algorithm 2 Loopless Katyusha (L-Katyusha)

Parameters: θ_1, θ_2 , probability $p \in (0, 1]$

Initialization: Choose $y^0 = w^0 = z^0 \in \mathbb{R}^d$, stepsize $\eta = \frac{\theta_2}{(1+\theta_2)\theta_1}$ and set $\sigma = \frac{\mu}{L}$

for $k = 0, 1, 2, \dots$ **do**

$$x^k = \theta_1 z^k + \theta_2 w^k + (1 - \theta_1 - \theta_2) y^k$$

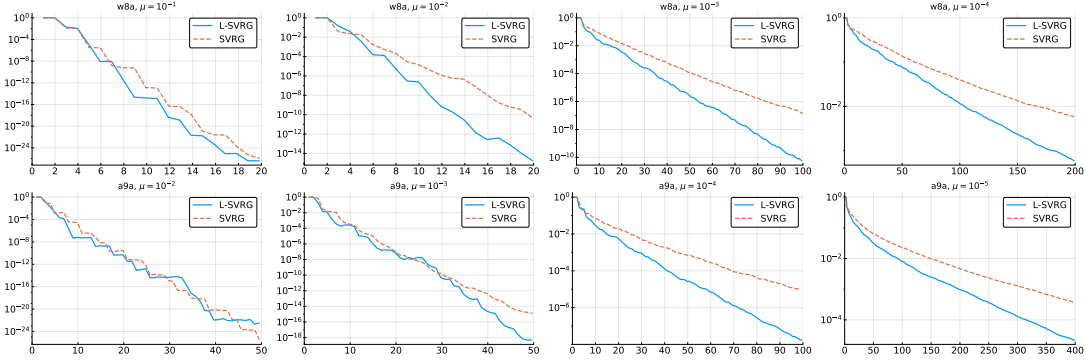
$$g^k = \nabla f_i(x^k) - \nabla f_i(w^k) + \nabla f(w^k) \quad (i \in \{1, \dots, n\} \text{ is sampled uniformly at random})$$

$$z^{k+1} = \frac{1}{1+\eta\sigma} (\eta\sigma x^k + z^k - \frac{\eta}{L} g^k)$$

$$y^{k+1} = x^k + \theta_1 (z^{k+1} - z^k)$$

$$w^{k+1} = \begin{cases} y^k & \text{with probability } p \\ w^k & \text{with probability } 1 - p \end{cases}$$

end for


 Figure 1: Comparison of SVRG and L-SVRG for different datasets and regularizer weights μ .

The exact analogy applies to the reference point w^k (at which a full gradient is computed) as for L-SVRG. Instead of updating this point in a deterministic way every m iteration, we use the probabilistic update with parameter p , when we update w^{k+1} to the current iterate y^k with this probability and is left unchanged with probability $1 - p$. As we shall see next, the same choice $p = 1/n$ as for L-SVRG leads to complexity identical to that of original Katyusha.

Convergence theory. In comparison to L-SVRG, we don't use *gradient mapping* as the key component of our analysis. Instead, we prove convergence of functional values in y^k, w^k and point-wise convergence of z^k . This is summarized in the following Lyapunov function:

$$\Psi^k = \mathcal{Z}^k + \mathcal{Y}^k + \mathcal{W}^k, \quad (13)$$

where $\mathcal{Z}^k = \frac{L(1+\eta\sigma)}{2\eta} \|z^k - x^*\|^2$, $\mathcal{Y}^k = \frac{1}{\theta_1}(f(y^k) - f(x^*))$, $\mathcal{W}^k = \frac{\theta_2(1+\theta_1)}{p\theta_1}(f(w^k) - f(x^*))$. Note that even if x^k is not in this function, its point-wise convergence is directly implied by the convergence of Ψ^k due to the definition of x^k in Algorithm 2 and L -smoothness of f .

The analysis involves five lemmas, followed by the convergence summarized in the main theorem. The lemmas highlight important steps of our analysis. The simplicity of our approach is still preserved: all lemmas and the main theorem can be proved on not more than two pages.

Our first lemma upper bounds the variance of the gradient estimator g^k , which eventually goes to zero as our algorithm progresses.

Lemma 6 *We have*

$$\mathbb{E} \left[\left\| g^k - \nabla f(x^k) \right\|^2 \right] \leq 2L \left(f(w^k) - f(x^k) - \langle \nabla f(x^k), w^k - x^k \rangle \right). \quad (14)$$

Next two lemmas are more technical, but essential for proving the convergence.

Lemma 7 *We have*

$$\langle g^k, x^* - z^{k+1} \rangle + \frac{\mu}{2} \|x^k - x^*\|^2 \geq \frac{L}{2\eta} \|z^k - z^{k+1}\|^2 + \mathcal{Z}^{k+1} - \frac{1}{1 + \eta\sigma} \mathcal{Z}^k. \quad (15)$$

Lemma 8 *We have*

$$\frac{1}{\theta_1} \left(f(y^{k+1}) - f(x^k) \right) - \frac{\theta_2}{2L\theta_1} \left\| g^k - \nabla f(x^k) \right\|^2 \leq \frac{L}{2\eta} \|z^{k+1} - z^k\|^2 + \langle g^k, z^{k+1} - z^k \rangle. \quad (16)$$

Finally, we use the update of Algorithm 2 to decompose \mathcal{W}^{k+1} in terms of \mathcal{W}^k and \mathcal{Y}^k , which is one of the main components that allow for simpler analysis than the one of original Katyusha.

Lemma 9 *We have*

$$\mathbb{E} \left[\mathcal{W}^{k+1} \right] = (1 - p)\mathcal{W}^k + \theta_2(1 + \theta_1)\mathcal{Y}^k. \quad (17)$$

Putting all lemmas together, we obtain the following contraction of the Lyapunov function (7).

Lemma 10 *Let $\theta_1, \theta_2 > 0$, $\theta_1 + \theta_2 \leq 1$, $\sigma = \frac{\mu}{L}$ and $\eta = \frac{\theta_2}{(1+\theta_2)\theta_1}$, then we have*

$$\mathbb{E} \left[\mathcal{Z}^{k+1} + \mathcal{Y}^{k+1} + \mathcal{W}^{k+1} \right] \leq \frac{1}{1 + \eta\sigma} \mathcal{Z}^k + (1 - \theta_1(1 - \theta_2))\mathcal{Y}^k + \left(1 - \frac{p\theta_1}{1 + \theta_1} \right) \mathcal{W}^k. \quad (18)$$

In order to obtain a recursion involving the Lyapunov function on the right-hand side of (18)

Theorem 11 *Let $\theta_1 = \min\{\sqrt{2\sigma n/3}, 1/2\}$, $\theta_2 = 1/2$, $p = 1/n$. Then $\mathbb{E} [\Psi^k] \leq \varepsilon\Psi^0$ after the following number of iterations: $k = \mathcal{O}((n + \sqrt{nL/\mu}) \log 1/\varepsilon)$.*

Proof From Lemma 10 we get

$$\mathbb{E} \left[\Psi^{k+1} \right] \leq \max \left\{ \frac{1}{(1 + \eta\sigma)}, 1 - \theta_1(1 - \theta_2), 1 - \frac{p\theta_1}{(1 + \theta_1)} \right\} \Psi^k.$$

Setting $p = 1/n$, $\theta_1 = \min\{\sqrt{2\sigma n/3}, 1/2\}$, $\theta_2 = 1/2$, and unrolling conditional probability one obtains $\mathbb{E} [\Psi^{k+1}] \leq (1 - \theta)\mathbb{E} [\Psi^k]$, where $\theta = \min\{\sigma/6\theta_1, \theta_1/2n\}$. Choosing $\sigma = \mu/L$ concludes the proof. ■

Discussion. One can show by analyzing (18) that for ill-conditioned problems ($n < L/\mu$), the iteration complexity is $\mathcal{O}(\sqrt{L/\mu p} \log 1/\varepsilon)$. Algorithm 2 calls stochastic gradient oracle $\mathcal{O}(1 + pn)$ times per iteration in expectation. Thus, the total complexity is $\mathcal{O}((1 + pn)\sqrt{L/\mu p} \log 1/\varepsilon)$. One can see that $p = \Theta(1/n)$ leads to optimal rate.

5. Numerical Experiments

In this section, we perform experiments with logistic regression for binary classification with L_2 regularizer, where our loss function has the form $f_i(x) = \log(1 + \exp(-b_i a_i^\top x)) + \frac{\mu}{2} \|x\|^2$, where $a_i \in \mathbb{R}^d$, $b_i \in \{-1, +1\}$, $i \in [n]$. Hence, f is smooth and μ -strongly convex. We use four LIBSVM library²: *a9a*, *w8a*, *mushrooms*, *phishing*, *cod-rna*.

We compare our methods L-SVRG and L-Katyusha with their original version. It is well-known that whenever practical, SAGA is a bit faster than SVRG. While a comparison to SAGA seems natural as it also does not have a double loop structure, we position our loopless methods for applications where the high memory requirements of SAGA prevent it to be applied. Thus, we do not compare to SAGA.

Plots are constructed in such a way that the y -axis displays $\|x^k - x^*\|^2$ for L-SVRG and $\|y^k - x^*\|^2$ for L-Katyusha, where x^* were obtained by running gradient descent for a large

2. The LIBSVM dataset collection is available at <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

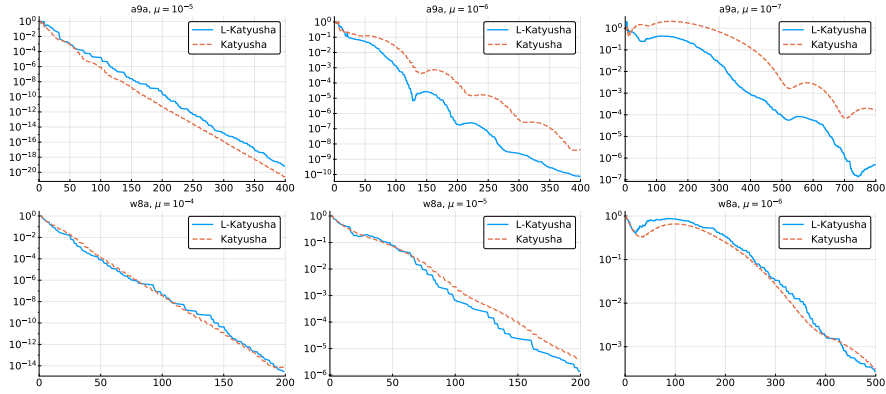


Figure 2: Comparison of Katyusha & L-Katyusha for different datasets and regularizer weights μ .

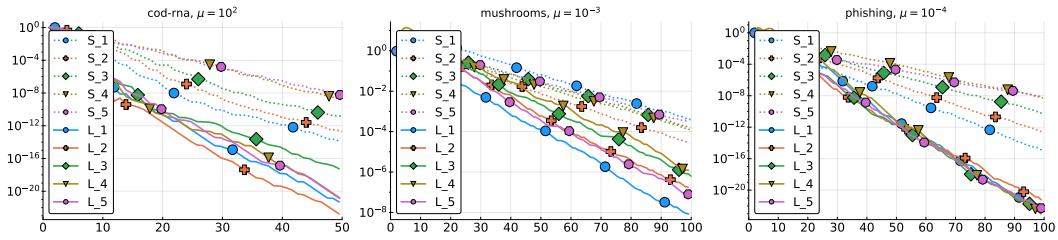


Figure 3: Comparison of SVRG (S) and L-SVRG (L) for several choices of expected outer loop length (L-SVRG) or deterministic outer loop length (SVRG). Numbers 1–5 correspondent to loop-lengths $n, \sqrt[4]{\kappa n^3}, \sqrt{\kappa n}, \sqrt[4]{\kappa^3 n}, \kappa$, respectively, where $\kappa = L/\mu$.

number of epochs. The x -axis displays the number of epochs (full gradient evaluations). That is, n computations of $\nabla f_i(x)$ equals one epoch.

Superior practical behaviour of the loopless approach. Here we show that L-SVRG and L-Katyusha perform better in experiments than their loopy variants. In terms of theoretical iteration complexity, both the loopy and the loopless methods are the same. However, as we can see from Figure 1, the improvement of the loopless approach can be significant. One can see that for these datasets, L-SVRG is always better than SVRG, and can be faster by several orders of magnitude! Looking at Figure 2, we see that the performance of L-Katyusha is at least as good as that of Katyusha, and can be significantly faster in some cases. All parameters of the methods were chosen as suggested by theory. For L-SVRG and L-Katyusha they are chosen based on Theorems 5 and 11, respectively. For SVRG and Katyusha we also choose the parameters based on the theory, as described in the original papers. The initial point x^0 is chosen to be the origin.

Different choices of probability/ outer loop size. We now compare several choices of the probability p of updating the full gradient for SVRG and several outer loop sizes m for SVRG. Since our analysis guarantees the optimal rate for any choice of p between $1/n$ and μ/L for well condition problems, we decided to perform experiments for p within this range. More precisely, we choose 5 values of p , uniformly distributed in logarithmic scale across this interval, and thus our choices are $n, \sqrt[4]{\kappa n^3}, \sqrt{\kappa n}, \sqrt[4]{\kappa^3 n},$ and κ , where $\kappa = L/\mu$, denoted in the figures by 1, 2, 3, 4, 5, respectively. Since the expected “outer loop” length (length for which reference point stays the same) is $1/p$, for

SVRG AND KATYUSHA ARE BETTER WITHOUT THE OUTER LOOP

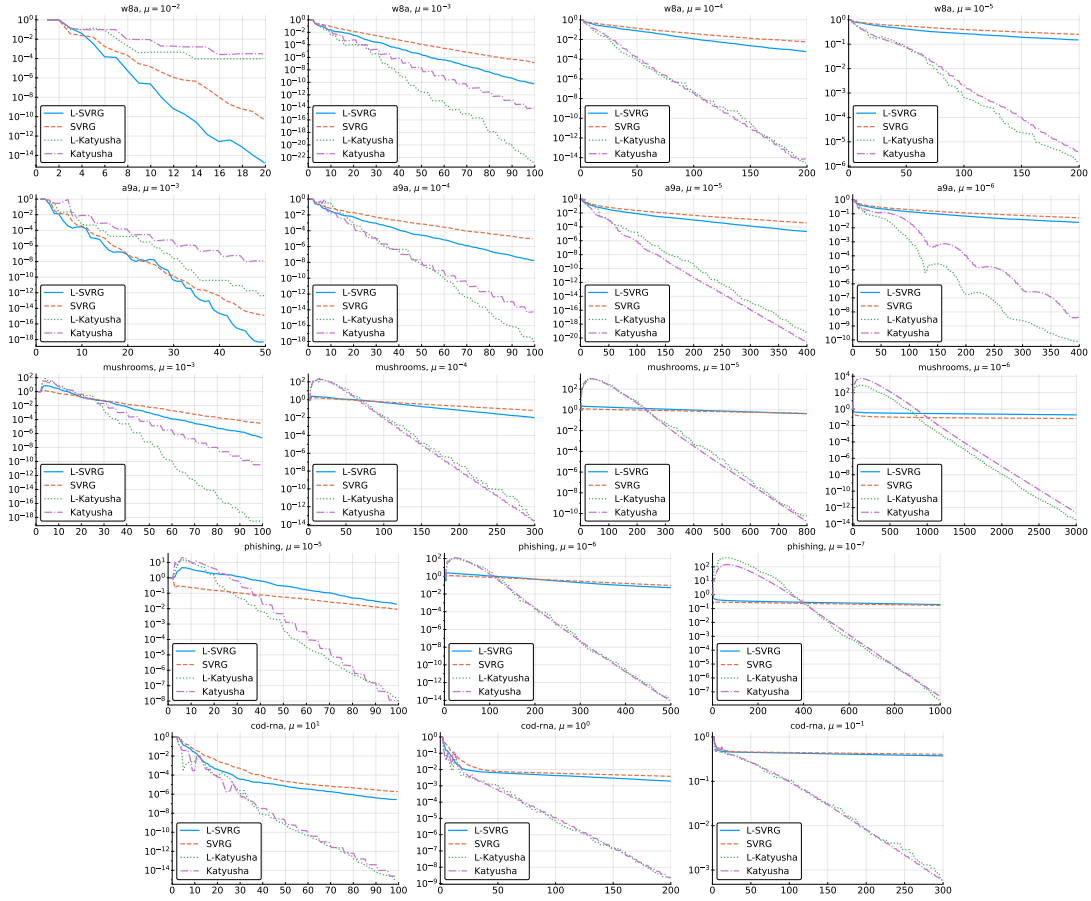


Figure 4: All methods together for different datasets and different regularizer weights.

SVRG we choose $m = 1/p$. Looking at Figure 3, one can see that L-SVRG is very *robust* to the choice of p from the “optimal interval” predicted by our theory. Moreover, *even the worst case for L-SVRG outperforms the best case for SVRG*.

All methods together. Finally, we provide all algorithms together in one plot for different datasets with different regularizer weight, thus with different condition numbers, displayed in Figure 4. As for the previous experiments, loopless methods are not worse and sometimes significantly better.

References

Zeyuan Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1200–1205. ACM, 2017.

Adel Bibi, Alibek Sailanbayev, Bernard Ghanem, Robert Mansel Gower, and Peter Richtárik. Improving SAGA via a probabilistic interpolation with gradient descent. *arXiv: 1806.05633*, 2018.

- Dominik Csiba and Peter Richtárik. Primal method for ERM with flexible mini-batching schemes and non-convex losses. *arXiv:1506.02227*, 2015.
- Dominik Csiba, Zheng Qu, and Peter Richtárik. Stochastic dual coordinate ascent with adaptive probabilities. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 674–683, 2015.
- Aaron Defazio. A simple practical accelerated method for finite sums. In *Advances in Neural Information Processing Systems*, pages 676–684, 2016.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2014a.
- Aaron Defazio, Tiberio Caetano, and Justin Domke. Finito: A faster, permutable incremental gradient method for Big Data problems. *The 31st International Conference on Machine Learning*, 2014b.
- Nikita Doikov and Peter Richtárik. Randomized block cubic Newton method. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- Robert Mansel Gower, Donald Goldfarb, and Peter Richtárik. Stochastic block BFGS: squeezing more curvature out of data. In *33rd International Conference on Machine Learning*, pages 1869–1878, 2016.
- Robert Mansel Gower, Peter Richtárik, and Francis Bach. Stochastic quasi-gradient methods: variance reduction via Jacobian sketching. *arXiv:1805.02632*, 2018.
- Thomas Hofmann, Aurelien Lucchi, Simon Lacoste-Julien, and Brian McWilliams. Variance reduced stochastic gradient descent with neighbors. In *Advances in Neural Information Processing Systems*, pages 2305–2313, 2015.
- Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
- Jakub Konečný and Peter Richtárik. S2GD: Semi-stochastic gradient descent methods. *Frontiers in Applied Mathematics and Statistics*, pages 1–14, 2017.
- Jakub Konečný, Jie Lu, Peter Richtárik, and Martin Takáč. Mini-batch semi-stochastic gradient descent in the proximal setting. *IEEE Journal of Selected Topics in Signal Processing*, 10(2): 242–255, 2016.
- Lihua Lei, Cheng Ju, Jianbo Chen, and Michael I Jordan. Non-convex finite-sum optimization via CSSG methods. In *Advances in Neural Information Processing Systems*, pages 2348–2358, 2017.
- Julien Mairal. Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal on Optimization*, 25(2):829–855, 2015.
- A Nemirovski, A Juditsky, G Lan, and A Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.

- Arkadi Nemirovsky and David B. Yudin. *Problem complexity and method efficiency in optimization*. Wiley, New York, 1983.
- Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *International Conference on Machine Learning*, pages 2613–2621, 2017.
- Zheng Qu, Peter Richtárik, and Tong Zhang. Quartz: Randomized dual coordinate ascent with arbitrary sampling. In *Advances in Neural Information Processing Systems 28*, pages 865–873, 2015.
- Zheng Qu, Peter Richtárik, Martin Takáč, and Olivier Fercoq. SDNA: Stochastic dual Newton ascent for empirical risk minimization. In *The 33rd International Conference on Machine Learning*, pages 1823–1832, 2016.
- Nicolas Le Roux, Mark Schmidt, and Francis Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In *Advances in Neural Information Processing Systems*, pages 2663–2671, 2012.
- Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1-2):83–112, 2017.
- Shai Shalev-Shwartz. SDCA without duality, regularization, and individual convexity. In *International Conference on Machine Learning*, pages 747–754, 2016.
- Shai Shalev-Shwartz and Tong Zhang. Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization. In *International Conference on Machine Learning*, pages 64–72, 2014.
- Martin Takáč, Avleen Bijral, Peter Richtárik, and Nathan Srebro. Mini-batch primal and dual methods for SVMs. In *30th International Conference on Machine Learning*, pages 537–552, 2013.
- Kaiwen Zhou. Direct acceleration of SAGA using sampled negative momentum. *arXiv preprint arXiv:1806.11048*, 2018.
- Kaiwen Zhou, Fanhua Shang, and James Cheng. A simple stochastic variance reduced algorithm with fast convergence rates. *arXiv preprint arXiv:1806.11027*, 2018.

Appendix

Appendix A. Auxiliary Lemmas

Lemma 12 For random vector $x \in \mathbb{R}^d$ and any $y \in \mathbb{R}^d$, the variance of y can be decomposed as

$$\mathbb{E} \left[\|x - \mathbb{E}[x]\|^2 \right] = \mathbb{E} \left[\|x - y\|^2 \right] - \mathbb{E} \left[\|\mathbb{E}[x] - y\|^2 \right]. \quad (19)$$

The next lemma is a consequence of Jensen's inequality applied to $x \mapsto \|x\|^2$.

Lemma 13 For any vectors $a_1, a_2, \dots, a_k \in \mathbb{R}^d$, the following inequality holds:

$$\left\| \sum_{i=1}^k a_i \right\|^2 \leq k \sum_{i=1}^k \|a_i\|^2. \quad (20)$$

Appendix B. Proofs for Algorithm 1 (L-SVRG)

In all proofs below, we will for simplicity write $f^* \stackrel{\text{def}}{=} f(x^*)$.

B.1. Proof of Lemma 1

Definition of x^{k+1} and unbiasedness of g^k guarantee that

$$\begin{aligned} \mathbb{E} \left[\left\| x^{k+1} - x^* \right\|^2 \right] &= \mathbb{E} \left[\left\| x^k - x^* - \eta g^k \right\|^2 \right] \\ &\stackrel{\text{Alg. 1}}{=} \left\| x^k - x^* \right\|^2 + \mathbb{E} \left[2\eta \langle g^k, x^* - x^k \rangle \right] + \eta^2 \mathbb{E} \left[\left\| g^k \right\|^2 \right] \\ &\stackrel{(2)}{=} \left\| x^k - x^* \right\|^2 + 2\eta \langle \nabla f(x^k), x^* - x^k \rangle + \eta^2 \mathbb{E} \left[\left\| g^k \right\|^2 \right] \\ &\stackrel{(4)}{\leq} \left\| x^k - x^* \right\|^2 + 2\eta \left(f^* - f(x^k) - \frac{\mu}{2} \left\| x^k - x^* \right\| \right) + \eta^2 \mathbb{E} \left[\left\| g^k \right\|^2 \right] \\ &= \left\| x^k - x^* \right\|^2 (1 - \eta\mu) + 2\eta \left(f^* - f(x^k) \right) + \eta^2 \mathbb{E} \left[\left\| g^k \right\|^2 \right]. \end{aligned}$$

B.2. Proof of Lemma 2

Using definition of g^k

$$\begin{aligned} \mathbb{E} \left[\left\| g^k \right\|^2 \right] &\stackrel{\text{Alg. 1}}{=} \mathbb{E} \left[\left\| \nabla f_i(x^k) - \nabla f_i(x^*) + \nabla f_i(x^*) - \nabla f_i(w^k) + \nabla f(w^k) \right\|^2 \right] \\ &\stackrel{(20)}{\leq} 2\mathbb{E} \left[\left\| \nabla f_i(x^k) - \nabla f_i(x^*) \right\|^2 \right] + \\ &\quad 2\mathbb{E} \left[\left\| \nabla f_i(x^*) - \nabla f_i(w^k) - \mathbb{E} \left[\nabla f_i(x^*) - \nabla f_i(w^k) \right] \right\|^2 \right] \\ &\stackrel{(3)+(19)}{\leq} 4L(f(x^k) - f^*) + 2\mathbb{E} \left[\left\| \nabla f_i(w^k) - \nabla f_i(x^*) \right\|^2 \right] \\ &\stackrel{(8)}{=} 4L(f(x^k) - f^*) + \frac{p}{2\eta^2} \mathcal{D}^k. \end{aligned}$$

B.3. Proof of Lemma 3

$$\begin{aligned}
 \mathbb{E} \left[\mathcal{D}^{k+1} \right] &\stackrel{\text{Alg. 1}}{=} (1-p)\mathcal{D}^k + p \frac{4\eta^2}{pn} \sum_{i=1}^n \left\| \nabla f(x^k) - \nabla f(x^*) \right\|^2 \\
 &\stackrel{(3)}{\leq} (1-p)\mathcal{D}^k + 8L\eta^2(f(x^k) - f^*).
 \end{aligned}$$

B.4. Proof of Lemma 4

Combining Lemmas 1 and 3 we obtain

$$\begin{aligned}
 \mathbb{E} \left[\left\| x^{k+1} - x^* \right\|^2 + \mathcal{D}^{k+1} \right] &\stackrel{(9)+(11)}{\leq} (1-\mu\eta) \left\| x^k - x^* \right\|^2 + 2\eta(f^* - f(x^k)) + \eta^2 \mathbb{E} \left[\left\| g^k \right\|^2 \right] \\
 &\quad + (1-p)\mathcal{D}^k + 8L\eta^2(f(x^k) - f^*) \\
 &\stackrel{(10)}{\leq} (1-\mu\eta) \left\| x^k - x^* \right\|^2 + (1-p)\mathcal{D}^k + (2\eta - 8L\eta^2)(f^* - f(x^k)) \\
 &\quad + \eta^2 \left(4L(f(x^k) - f^*) + \frac{p}{2\eta^2} \mathcal{D}^k \right) \\
 &= (1-\mu\eta) \left\| x^k - x^* \right\|^2 + \left(1 - \frac{p}{2} \right) \mathcal{D}^k + (2\eta - 12L\eta^2)(f^* - f(x^k)).
 \end{aligned}$$

Now we use the fact that $\eta \leq \frac{1}{6L}$ and obtain the desired inequality:

$$\mathbb{E} \left[\left\| x^{k+1} - x^* \right\|^2 + \mathcal{D}^{k+1} \right] \leq (1-\mu\eta) \left\| x^k - x^* \right\|^2 + \left(1 - \frac{p}{2} \right) \mathcal{D}^k.$$

Appendix C. Proofs for Algorithm 2 (L-Katyusha)

C.1. Proof of Lemma 6

To upper bound the variance of g^k we first uses its definition

$$\begin{aligned}
 \mathbb{E} \left[\left\| g^k - \nabla f(x^k) \right\|^2 \right] &\stackrel{\text{Alg. 2}}{=} \mathbb{E} \left[\left\| \nabla f_i(x^k) - \nabla f_i(w^k) - \mathbb{E} \left[\nabla f_i(x^k) - \nabla f_i(w^k) \right] \right\|^2 \right] \\
 &\stackrel{(19)}{\leq} \mathbb{E} \left[\left\| \nabla f_i(x^k) - \nabla f_i(w^k) \right\|^2 \right] \\
 &\stackrel{(3)}{\leq} 2L \left(f(w^k) - f(x^k) - \left\langle \nabla f(x^k), w^k - x^k \right\rangle \right).
 \end{aligned}$$

C.2. Proof of Lemma 7

We start with the definition of z^{k+1}

$$z^{k+1} \stackrel{\text{Alg. 2}}{=} \frac{1}{1+\eta\sigma} \left(\eta\sigma x^k + z^k - \frac{\eta}{L} g^k \right),$$

which implies $\frac{\eta}{L}g^k = \eta\sigma(x^k - z^{k+1}) + (z^k - z^{k+1})$, which further implies that

$$\begin{aligned}
 \langle g^k, z^{k+1} - x^* \rangle &= \mu \langle x^k - z^{k+1}, z^{k+1} - x^* \rangle + \frac{L}{\eta} \langle z^k - z^{k+1}, z^{k+1} - x^* \rangle \\
 &= \frac{\mu}{2} \left(\|x^k - x^*\|^2 - \|x^k - z^{k+1}\|^2 - \|z^{k+1} - x^*\|^2 \right) \\
 &\quad + \frac{L}{2\eta} \left(\|z^k - x^*\|^2 - \|z^k - z^{k+1}\|^2 - \|z^{k+1} - x^*\|^2 \right) \\
 &\leq \frac{\mu}{2} \|x^k - x^*\|^2 + \frac{L}{2\eta} \left(\|z^k - x^*\|^2 - (1 + \eta\sigma) \|z^{k+1} - x^*\|^2 \right) \\
 &\quad - \frac{L}{2\eta} \|z^k - z^{k+1}\|^2.
 \end{aligned}$$

C.3. Proof of Lemma 8

$$\begin{aligned}
 &\frac{L}{2\eta} \|z^{k+1} - z^k\|^2 + \langle g^k, z^{k+1} - z^k \rangle \\
 &= \frac{1}{\theta_1} \left(\frac{L}{2\eta\theta_1} \|\theta_1(z^{k+1} - z^k)\|^2 + \langle g^k, \theta_1(z^{k+1} - z^k) \rangle \right) \\
 \stackrel{\text{Alg. 2}}{=} &\frac{1}{\theta_1} \left(\frac{L}{2\eta\theta_1} \|y^{k+1} - x^k\|^2 + \langle g^k, y^{k+1} - x^k \rangle \right) \\
 &= \frac{1}{\theta_1} \left(\frac{L}{2\eta\theta_1} \|y^{k+1} - x^k\|^2 + \langle \nabla f(x^k), y^{k+1} - x^k \rangle + \langle g^k - \nabla f(x^k), y^{k+1} - x^k \rangle \right) \\
 &= \frac{1}{\theta_1} \left(\frac{L}{2} \|y^{k+1} - x^k\|^2 + \langle \nabla f(x^k), y^{k+1} - x^k \rangle + \frac{L}{2} \left(\frac{1}{\eta\theta_1} - 1 \right) \|y^{k+1} - x^k\|^2 \right) \\
 &\quad + \frac{1}{\theta_1} \left(\langle g^k - \nabla f(x^k), y^{k+1} - x^k \rangle \right) \\
 \stackrel{(3)}{\geq} &\frac{1}{\theta_1} \left(f(y^{k+1}) - f(x^k) + \frac{L}{2} \left(\frac{1}{\eta\theta_1} - 1 \right) \|y^{k+1} - x^k\|^2 + \langle g^k - \nabla f(x^k), y^{k+1} - x^k \rangle \right) \\
 \geq &\frac{1}{\theta_1} \left(f(y^{k+1}) - f(x^k) - \frac{\eta\theta_1}{2L(1-\eta\theta_1)} \|g^k - \nabla f(x^k)\|^2 \right) \\
 = &\frac{1}{\theta_1} \left(f(y^{k+1}) - f(x^k) - \frac{\theta_2}{2L} \|g^k - \nabla f(x^k)\|^2 \right),
 \end{aligned}$$

where the last inequality uses the Young's inequality in the form of $\langle a, b \rangle \geq -\frac{\|a\|^2}{2\beta} - \frac{\beta\|b\|^2}{2}$ for $\beta = \frac{\eta\theta_1}{L(1-\eta\theta_1)}$, which concludes the proof.

C.4. Proof of Lemma 9

From the definition of w^{k+1} in Algorithm 2 we have

$$\mathbb{E} \left[f(w^{k+1}) \right] \stackrel{\text{Alg. 2}}{=} (1-p)f(w^k) + pf(y^k). \quad (21)$$

The rest of proof follows from the definition of \mathcal{W}^k (17).

C.5. Proof of Lemma 10

Combining all the previous lemmas together, we obtain

$$\begin{aligned}
f^* &\stackrel{(4)}{\geq} f(x^k) + \langle \nabla f(x^k), x^* - x^k \rangle + \frac{\mu}{2} \|x^k - x^*\|^2 \\
&= f(x^k) + \frac{\mu}{2} \|x^k - x^*\|^2 + \langle \nabla f(x^k), x^* - z^k + z^k - x^k \rangle \\
&\stackrel{\text{Alg. 2}}{=} f(x^k) + \frac{\mu}{2} \|x^k - x^*\|^2 + \langle \nabla f(x^k), x^* - z^k \rangle + \frac{\theta_2}{\theta_1} \langle \nabla f(x^k), x^k - w^k \rangle \\
&\quad + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} \langle \nabla f(x^k), x^k - y^k \rangle \\
&\stackrel{(2)}{\geq} f(x^k) + \frac{\theta_2}{\theta_1} \langle \nabla f(x^k), x^k - w^k \rangle + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(x^k) - f(y^k)) \\
&\quad + \mathbb{E} \left[\frac{\mu}{2} \|x^k - x^*\|^2 + \langle g^k, x^* - z^{k+1} \rangle + \langle g^k, z^{k+1} - z^k \rangle \right] \\
&\stackrel{(15)}{\geq} f(x^k) + \frac{\theta_2}{\theta_1} \langle \nabla f(x^k), x^k - w^k \rangle + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(x^k) - f(y^k)) \\
&\quad + \mathbb{E} \left[\mathcal{Z}^{k+1} - \frac{1}{1 + \eta\sigma} \mathcal{Z}^k \right] + \mathbb{E} \left[\langle g^k, z^{k+1} - z^k \rangle + \frac{L}{2\eta} \|z^k - z^{k+1}\|^2 \right] \\
&\stackrel{(16)}{\geq} f(x^k) + \frac{\theta_2}{\theta_1} \langle \nabla f(x^k), x^k - w^k \rangle + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(x^k) - f(y^k)) \\
&\quad + \mathbb{E} \left[\mathcal{Z}^{k+1} - \frac{1}{1 + \eta\sigma} \mathcal{Z}^k \right] + \mathbb{E} \left[\frac{1}{\theta_1} (f(y^{k+1}) - f(x^k)) - \frac{\theta_2}{2L\theta_1} \|g^k - \nabla f(x^k)\|^2 \right] \\
&\stackrel{(14)}{\geq} f(x^k) + \frac{\theta_2}{\theta_1} \langle \nabla f(x^k), x^k - w^k \rangle + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(x^k) - f(y^k)) \\
&\quad + \mathbb{E} \left[\mathcal{Z}^{k+1} - \frac{1}{1 + \eta\sigma} \mathcal{Z}^k + \frac{1}{\theta_1} (f(y^{k+1}) - f(x^k)) \right] \\
&\quad - \frac{\theta_2}{\theta_1} \mathbb{E} \left[f(w^k) - f(x^k) - \langle \nabla f(x^k), w^k - x^k \rangle \right] \\
&= f(x^k) + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(x^k) - f(y^k)) - \frac{1}{1 + \eta\sigma} \mathcal{Z}^k - \frac{\theta_2}{\theta_1} (f(w^k) - f(x^k)) \\
&\quad + \mathbb{E} \left[\mathcal{Z}^{k+1} + \frac{1}{\theta_1} (f(y^{k+1}) - f(x^k)) \right],
\end{aligned}$$

where in the second inequality we use also convexity of $f(x)$.

$$\begin{aligned}
x^k &\stackrel{\text{Alg. 2}}{=} \theta_1 z^k + \theta_2 w^k + (1 - \theta_1 - \theta_2) y^k \\
z^k - x^k &= \frac{\theta_2}{\theta_1} (x^k - w^k) + \frac{1 - \theta_1 - \theta_2}{\theta_1} (x^k - y^k).
\end{aligned}$$

After rearranging we get

$$\frac{1}{1 + \eta\sigma} \mathcal{Z}^k + (1 - \theta_1 - \theta_2) \mathcal{Y}^k + \frac{\theta_2}{\theta_1} (f(w^k) - f^*) \geq \mathbb{E} \left[\mathcal{Z}^{k+1} + \mathcal{Y}^{k+1} \right].$$

Using definition of \mathcal{W}^k we get

$$\mathbb{E} \left[\mathcal{Z}^{k+1} + \mathcal{Y}^{k+1} \right] \leq \frac{1}{1 + \eta\sigma} \mathcal{Z}^k + (1 - \theta_1 - \theta_2) \mathcal{Y}^k + \frac{p}{(1 + \theta_1)} \mathcal{W}^k. \quad (22)$$

Finally, using Lemma 9 we get

$$\begin{aligned} \mathbb{E} \left[\mathcal{Z}^{k+1} + \mathcal{Y}^{k+1} + \mathcal{W}^{k+1} \right] &\leq \frac{1}{1 + \eta\sigma} \mathcal{Z}^k + (1 - \theta_1 - \theta_2) \mathcal{Y}^k + \frac{p}{(1 + \theta_1)} \mathcal{W}^k \\ &\quad + (1 - p) \mathcal{W}^k + \theta_2 (1 + \theta_1) \mathcal{Y}^k \\ &= \frac{1}{1 + \eta\sigma} \mathcal{Z}^k + (1 - \theta_1 (1 - \theta_2)) \mathcal{Y}^k + \left(1 - \frac{p\theta_1}{1 + \theta_1} \right) \mathcal{W}^k, \end{aligned}$$

which concludes the proof.