

# BEYOND EXPECTATION: CONCENTRATION INEQUALITIES FOR RANDOMIZED ITERATIVE METHODS

TOBY ANDERSON, MAX COLLINS, JAMIE HADDOCK, JACKIE LOK, ELIZAVETA REBROVA

**ABSTRACT.** Stochastic iterative methods are useful in a variety of large-scale numerical linear algebraic, machine learning, and statistical problems, in part due to their low-memory footprint. They are frequently used in a variety of applications, and thus it is imperative to have a thorough theoretical understanding of their behavior. Most theoretical convergence results for stochastic iterative methods provide bounds on the expected error of the iterates, and yield a type of average case analysis. However, understanding the behavior of these methods in the near-worst-case is desirable. For stochastic methods, this motivates providing bounds on the variance and concentration of their error, which can be used to generate confidence intervals around the bounds on their expected error.

Here, we provide upper bounds for the concentration and variance of the error of a general class of linear stochastic iterative methods, including the randomized Kaczmarz method and the randomized Gauss–Seidel method, and a more general class of nonlinear stochastic iterative methods, including the randomized Kaczmarz method for systems of linear inequalities.

## 1. INTRODUCTION

Stochastic or randomized iterative methods have become increasingly popular approaches for a variety of large-scale data problems as these methods typically have low-memory footprint and are accompanied by attractive theoretical guarantees [35]. Indeed, the scale of modern problems often make application of direct or non-iterative methods challenging or infeasible. Examples of randomized iterative methods that have found popularity in recent years include the stochastic gradient descent method [40] and the randomized Kaczmarz method [42].

Theoretical guarantees that accompany randomized iterative methods tend to focus upon bounding the expected error of the sequence of iterates [13, 34, 31]. For example, the seminal work of [42] proved that when applied to a consistent linear system  $\mathbf{Ax} = \mathbf{b}$  with unique solution  $\mathbf{x}^*$ , the randomized Kaczmarz method (with a specific sampling distribution, see Section 2.1.1 below for details) converges at least linearly in expectation with the guarantee

$$(1) \quad \mathbb{E}\|\mathbf{e}_k\|^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}\right)^k \|\mathbf{e}_0\|^2.$$

Here,  $\mathbf{e}_k := \mathbf{x}_k - \mathbf{x}^*$  denotes the error vector between the  $k$ th iterate and the solution, and  $\sigma_{\min}(\mathbf{A})$  is the minimum singular value of the matrix  $\mathbf{A}$ . However, understanding the average case behavior of a randomized method can be an insufficient measure of how well it performs. Indeed, much effort is typically put into understanding the worst-case behavior of even simple algorithms [26]. In the context of randomized methods, one may interpolate between the average and worst cases by proving bounds on the *concentration* of the error. These results provide upper bounds on the probability that the error deviates significantly above its mean (or an upper bound for its mean).

Besides yielding a better understanding of the behavior of randomized methods, concentration inequalities bounding the error of randomized methods can be used as an algorithmic tool to detect data inconsistency, outliers, and poor objective landscape geometry. For instance, if the residual error of a consistent linear system should be less than  $\epsilon$  after  $k_\epsilon$  iterations of a randomized method with high probability, and one encounters a system that has residual error well above  $\epsilon$  after  $k_\epsilon$  iterations, this may suggest that the system is inconsistent. Furthermore, the magnitude of the entries of the residual may yield information about the form of the system inconsistency (e.g., the positions of corruptions) [15].

Our work considers two classes of iterative methods. In Section 2, we examine the concentration and variance of the error of randomized iterative methods whose errors or residual errors in sequential iterations obey a linear relationship,  $\mathbf{e}_j = \mathbf{Y}_j \mathbf{e}_{j-1}$ , where  $\mathbf{Y}_j$  is independently sampled in the  $j$ th iteration from a

family of square matrices. In Subsection 2.1.1, we show that two families of methods for solving systems of linear equations  $\mathbf{Ax} = \mathbf{b}$  satisfy this recursive error relation. The *randomized Kaczmarz (RK)* methods have errors that satisfy

$$\mathbf{e}_j = \mathbf{Y}_j \mathbf{e}_{j-1}, \quad \text{with } \mathbf{e}_j := \mathbf{x}_j - \mathbf{x}^* \text{ and } \mathbf{Y}_j := \mathbf{I} - \frac{\mathbf{a}_{i_j} \mathbf{a}_{i_j}^\top}{\|\mathbf{a}_{i_j}\|^2},$$

where  $\mathbf{a}_{i_j}$  is the  $i_j$ th row of  $\mathbf{A}$ . It is known that the RK errors satisfy (1) [42]. The *randomized Gauss–Seidel (RGS)* methods have errors that satisfy

$$\mathbf{e}_j = \mathbf{Y}_j \mathbf{e}_{j-1}, \quad \text{with } \mathbf{e}_j := \mathbf{Ax}_j - \mathbf{Ax}^* \text{ and } \mathbf{Y}_j := \mathbf{I} - \frac{\mathbf{A}_{i_j} \mathbf{A}_{i_j}^\top}{\|\mathbf{A}_{i_j}\|^2},$$

where  $\mathbf{A}_{i_j}$  is the  $i_j$ th column of  $\mathbf{A}$ . It is known that the RGS errors satisfy (1) [27].

In Section 3, we consider a broader class of potentially nonlinear updating methods in which the iterates are given by  $\mathbf{x}_j = f_{i_j}(\mathbf{x}_{j-1})$ , where  $f_{i_j}$  is an independent and identically distributed (i.i.d.) sample from a fixed set  $F = \{f_1, f_2, \dots, f_m\}$  of updating functions in the  $j$ th iteration. In Subsection 3.1.1, we apply these results to the randomized Kaczmarz method for linear inequalities [27].

**1.1. Simple Markov inequality based bound.** Expectation-based guarantees like (1) are well-established for many randomized iterative methods. However, comparatively little is known about how much randomized iterative methods may deviate from their mean behavior, or how they *concentrate* around their mean. A natural first approach to deriving bounds on the concentration of the error of randomized methods is to apply Markov’s inequality to extend a convergence bound in expectation to one in probability. Doing so yields the following general lemma:

**Lemma 1.1.** *Consider a stochastic process  $\{\mathbf{x}_k : k \in \mathbb{N}\}$  approximating an element of nonempty convex  $S \subset \mathbb{R}^n$ , where  $\mathbf{x}_k = f_{i_k}(\mathbf{x}_{k-1})$  and  $f_{i_k}$  is independently and randomly selected from a set  $F = \{f_1, f_2, \dots, f_m\}$  at each time  $k$  according to a fixed distribution  $\mathcal{D}$ . Suppose that*

$$\mathbb{E}[d(\mathbf{x}_k, S)^2] \leq r^k d(\mathbf{x}_0, S)^2 \quad \text{for some } r \in (0, 1),$$

*and  $d(\mathbf{x}, S) := \inf_{\mathbf{s} \in S} \|\mathbf{x} - \mathbf{s}\|$  is defined with respect to a vector norm  $\|\cdot\|$ . Then, for any  $t > 0$ , it follows that*

$$(2) \quad \mathbb{P}(d(\mathbf{x}_k, S)^2 - \mathbb{E}[d(\mathbf{x}_k, S)^2] \geq t) \leq \frac{r^k d(\mathbf{x}_0, S)^2}{t}.$$

We compare our main results in Subsection 1.2 to this elementary result in our numerical experiments in Subsection 2.1.4. We have found that this bound is surprisingly difficult to outperform for small values of  $t$ . We note that Lemma 1.1, unlike those presented in Subsection 1.2, provides a one-sided bound, and thus is weaker. However, the case where the error deviates above its mean is likely of most interest practically.

**1.2. Contributions.** In this paper, we are interested in understanding the behavior of randomized iterative methods *beyond the average case*. Results for randomized methods often consider bounding the error *in expectation*, but less often provide bounds for how far the error of these methods can deviate above their average case bound. We provide bounds on the variance and concentration of commonly studied methods in the area of randomized numerical linear algebra and optimization, and additionally consider some high-probability bounds for the error.

Our first main result bounds the variance and concentration of the squared norm of the error of randomized methods whose error obeys a linear recurrence relation.

**Theorem 1.2.** *Let  $\mathbf{e}_k = \mathbf{Y}_k \mathbf{e}_{k-1}$ , where  $\mathbf{Y}_k \sim \mathbf{Y}$  is sampled i.i.d., and define  $\mu := \|\mathbb{E}[(\mathbf{Y}^\top \mathbf{Y})^{\otimes 2}]\|$  and  $\eta := \lambda_{\min}(\mathbb{E}[\mathbf{Y}^\top \mathbf{Y}])$ . Then for all  $k$ ,*

$$(3) \quad \text{Var}(\|\mathbf{e}_k\|^2) \leq (\mu^k - \eta^k) \cdot \|\mathbf{e}_0\|^4.$$

We prove this result in Subsection 2.1. It automatically extends to a concentration result of type (2) via Chebyshev’s inequality, and provides confidence intervals for the trajectories of randomized iterative methods: see Remark 1. We note that this confidence interval is two-sided (see Figure 1), but the upper bound on the error is most interesting from an algorithmic perspective. We showcase the improvement of

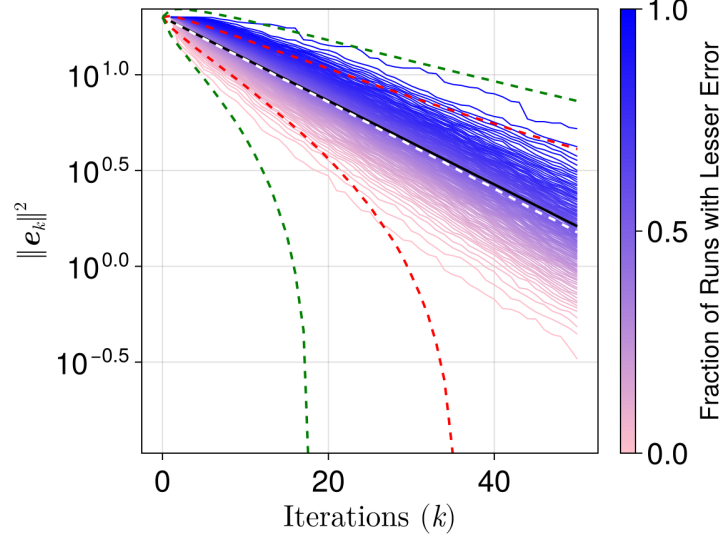


FIGURE 1. Visualization of errors of 500 independent trials of RK. Empirical mean error (white dashed line), bound (1) (black solid line), and the 75% (red dashed lines) and 95% (green dashed lines) confidence intervals for the error derived by combining Chebyshev's inequality with Theorem 1.2 and (4) are plotted.

the refined variance analysis on some standard methods, including RK and RGS, in Subsection 2.1.1. For both of these methods, one can bound

$$(4) \quad \mu^k - \eta^k \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}\right)^k - \left(1 - \frac{\sigma_{\max}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}\right)^k.$$

As a motivating example, we visualize the empirical concentration of the error of the RK method applied to a consistent system of equations defined by a well-conditioned, randomly-generated matrix  $\mathbf{A} \in \mathbb{R}^{1000 \times 20}$  with linearly decreasing singular values over 500 trials in Figure 1, along with the 75% and 95% confidence intervals for the error derived by combining Chebyshev's inequality with Theorem 1.2 and (4). We also plot the empirical mean of the error across the independent trials and the bound on the mean error (1). More numerical experiments are provided in Subsection 2.1.3.

Next, we prove high-probability results for the error of randomized methods that satisfy a linear recurrence relation. Unlike what one can get by simply applying Markov's inequality to the bound in expectation, these results crucially hold for the *whole random trajectory*, rather than for any fixed iteration.

**Theorem 1.3.** *Let  $\mathbf{e}_k = \mathbf{Y}_k \mathbf{e}_{k-1}$ , where each  $\mathbf{Y}_k$  is independently sampled from a family of  $n \times n$  matrices such that  $\sup_k \|\mathbb{E}[\mathbf{Y}_k^\top \mathbf{Y}_k]\| \leq \rho$ . Then, the following hold:*

(a) *For any  $\epsilon \in (0, 1]$ ,*

$$(5) \quad \mathbb{P}\left(\forall k \geq 0 : \|\mathbf{e}_k\|^2 \leq \epsilon^{-1} \rho^k \|\mathbf{e}_0\|^2\right) \geq 1 - \epsilon.$$

(b) *If moreover,  $\sup_k \|\mathbf{Y}_k^\top \mathbf{Y}_k\| \leq \alpha$  almost surely for some  $\alpha \geq 1$  (e.g., if  $\mathbf{Y}_k^\top \mathbf{Y}_k$  are contraction operators, then  $\alpha = 1$ ), then, for any  $k \geq 0$ ,*

$$(6) \quad \mathbb{P}\left[\sup_{0 \leq t \leq k} \|\mathbf{e}_t\|^2 \leq \exp\left(-k \cdot (1 - \rho) + \alpha \cdot \sqrt{2k \log(\epsilon^{-1})}\right) \|\mathbf{e}_0\|^2\right] \geq 1 - \epsilon.$$

This theorem is proved in Subsection 2.3. We note that the bound (5) is generally better for larger  $k$  or when  $\epsilon$  not too small. However, (6) can provide a stronger bound if one is interested in the first few iterations or in tight probabilistic bounds: e.g., with an exponentially small tolerance for the failure probability  $\epsilon = e^{-O(n)}$  with respect to the size of the data, (6) provides a stronger bound for the first  $O(n)$  iterations of the algorithm.

Theorem 1.3 applies to the errors of the standard RK and RGS methods with parameters  $\rho = 1 - \sigma_{\min}^2(\mathbf{A})/\|\mathbf{A}\|_F^2$  and  $\alpha = 1$ .

Finally, we show that the approach is not limited to methods whose error obeys a linear recurrence relation and provide an upper bound on the variance of error of nonlinear randomized methods as follows.

**Theorem 1.4.** *Consider a stochastic process  $\{\mathbf{x}_k : k \in \mathbb{N}\}$  approximating an element of a nonempty convex set  $S \subset \mathbb{R}^n$ , where  $\mathbf{x}_k = f_{i_k}(\mathbf{x}_{k-1})$  and  $f_{i_k}$  is independently selected from a set  $F = \{f_1, f_2, \dots, f_m\}$  at each time  $k$  according to a fixed distribution  $\mathcal{D}$ . Let  $d(\mathbf{x}, S) := \inf_{\mathbf{s} \in S} \|\mathbf{x} - \mathbf{s}\|$  denote the distance to  $S$  with respect to a vector norm  $\|\cdot\|$ . Suppose that*

$$\mathbb{E}[d(\mathbf{x}_k, S)^2] \leq r^k d(\mathbf{x}_0, S)^2 \quad \text{for some } r \in (0, 1),$$

and that there exists some  $D$  such that  $\sup_{k \in \mathbb{N}} d(\mathbf{x}_k, S) \leq D$ . Then, it follows that

$$(7) \quad \text{Var}(d(\mathbf{x}_k, S)^2) \leq D^2 r^k d(\mathbf{x}_0, S)^2.$$

This result provides confidence intervals through an application of Chebyshev's inequality. One can apply this result for the randomized Kaczmarz method for solving a system of linear inequalities [27] by setting  $D = d(\mathbf{x}_0, S)$  and  $r = 1 - 1/(L^2 \|\mathbf{A}\|_F^2)$ , where  $L$  is the Hoffman constant for the system [20]. We show this in Section 3 and apply these results to the randomized Kaczmarz method for linear feasibility in Subsection 3.1.1.

**1.3. Related work.** In the literature on stochastic gradient descent (SGD) and related methods, there has been some work proving bounds on the concentration of the error of SGD or high-probability convergence results for SGD and variants [37, 23, 11]. Most of these results apply to variants of stochastic gradient descent which average iterates to reduce the effect of noise and variance [10, 18, 33, 30]. These results are challenging to apply or generalize to the methods we consider in this paper due to their assumptions on the step size schedule and application to averaging methods. In [12], the authors mention the application of Markov's inequality to prove high-probability convergence. The latter paper also includes a nice survey of high-probability convergence results for a variety of SGD variants with differing assumptions on the problem to be solved, as well as new results for clipped variants of SGD on composite and distributed problems. In [6], the authors use Markov's bound for the complexity analysis of a block Kaczmarz-type algorithm.

In [8], the authors utilize Azuma's inequality, a concentration inequality for sequences of martingale random variables with bounded sequential differences, to bound the concentration and variance of the errors of a variety of consensus protocols. These protocols are well-studied in the discrete dynamical systems community and have applications in a variety of areas including distributed computing, opinion dynamics modelling, gene network models, and control theory. It has been recently noted that these discrete updates can be viewed as iterations of common iterative methods in numerical linear algebra applied to linear systems encoding the consensus problem [29, 14].

Some results for stronger (e.g., almost sure) convergence guarantees for RK in the streaming setting with independent measurement vectors are derived in [4, 28].

The following two approaches are most related to ours and we describe them in more detail:

**1. Matrix concentration bounds.** Many iterative methods of interest in randomized numerical linear algebra (e.g., variants of the randomized Kaczmarz and randomized Gauss–Seidel methods) can be interpreted as a product of projection matrices applied recursively to the iterates. These projection matrices are sampled from a fixed set of update matrices, usually defined by the rows or columns of the matrix defining a linear system or regression problem. For this reason, results providing tight concentration results for products of random matrices could be used to provide concentration bounds for the error of these randomized iterative methods. This line of research is explored in [21, 19, 25].

For example, we can apply [21, Theorem 7.1] to analyze the concentration of the error of the randomized Kaczmarz method [42]. Let  $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}^*$  be the error in the  $k$ th iteration of the RK method applied to a consistent, full-rank system  $\mathbf{Ax} = \mathbf{b}$  with  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Define  $\mathbf{Y}_i = \mathbf{I} - \mathbf{a}_{j_i} \mathbf{a}_{j_i}^\top / \|\mathbf{a}_{j_i}\|^2$  to be the random contraction sampled in the  $i$ th iteration, and define  $\mathbf{Z}_k = \mathbf{Y}_k \mathbf{Y}_{k-1} \cdots \mathbf{Y}_1$  so that  $\mathbf{e}_k = \mathbf{Z}_k \mathbf{e}_0$ . If

$\rho = 1 - \sigma_{\min}^2(\mathbf{A})/\|\mathbf{A}\|_F^2$ , then applying [21, Theorem 7.1] yields

$$\begin{aligned} \mathbb{P}[\|\mathbf{e}_k - \mathbb{E}\mathbf{e}_k\|^2 \geq t^2] &= \mathbb{P}[\|\mathbf{Z}_k \mathbf{e}_0 - \mathbb{E}\mathbf{Z}_k \mathbf{e}_0\|^2 \geq t^2] \\ &\leq \mathbb{P}[\|\mathbf{Z}_k - \mathbb{E}\mathbf{Z}_k\|^2 \|\mathbf{e}_0\|^2 \geq t^2] \\ (8) \quad &\leq n\rho^k \exp\left(\frac{-t^2}{2ek\|\mathbf{e}_0\|^2(\rho + 2 + 1/\rho)}\right) \end{aligned}$$

when  $t^2 \geq 2ek\|\mathbf{e}_0\|^2(\rho + 2 + 1/\rho)$ . Note that this result may only be applied for  $t = \Omega(\sqrt{k})$  where  $k$  is the number of iterations. We further note that this bound is qualitatively different from those that we will primarily consider in this paper; this concentration result bounds  $\|\mathbf{e}_k - \mathbb{E}\mathbf{e}_k\|^2$ , while we will consider bounds on  $\|\mathbf{e}_k\|^2 - \mathbb{E}\|\mathbf{e}_k\|^2$ . Regardless, we compare the upper bounds offered by our main results in Subsection 1.2 to the bound (8) in our numerical experiments in Subsection 2.1.4.

*2. Moment bounds.* Another line of research has sought to provide bounds on the rate of convergence of general moments of the error of variants of the Kaczmarz method. In [39], the authors consider the Generalized Block Randomized Kaczmarz (GBRK) methods, which are a general class of iterative methods which encompass the usual randomized Kaczmarz (RK) methods [42] and block RK methods [36]. They show that after a given stopping time, the  $d$ th moment converges with exponential rate. Our Theorem 2.1 is most closely related to their result [39, Theorem 3], but provides a simpler analysis, more specific bounds on the rate of convergence of the moments, and applies to a broader class of methods.

**1.4. Notation.** We use boldfaced lower-case Latin letters (e.g.,  $\mathbf{x}$ ) to denote vectors, and boldfaced upper-case Latin letters (e.g.,  $\mathbf{A}$ ) to denote matrices. We use unbolded lower-case Latin and Roman letters (e.g.,  $t$  and  $\mu$ ) to denote scalars. We denote by  $\mathbf{A}_j$  the  $j$ th column of matrix  $\mathbf{A}$  and by  $\mathbf{a}_i$  the  $i$ th row vector of matrix  $\mathbf{A}$ . We let  $[m]$  denote the set  $\{1, 2, \dots, m\}$ . The notation  $\|\mathbf{v}\|$  denotes the Euclidean norm of a vector  $\mathbf{v}$ , and  $\|\mathbf{A}\|$  the operator norm and  $\|\mathbf{A}\|_F$  the Frobenius norm of a matrix  $\mathbf{A}$ . We denote by  $\sigma_{\min}(\mathbf{A})$  and  $\sigma_{\max}(\mathbf{A})$  the smallest and largest singular value of the matrix  $\mathbf{A}$  respectively. We use

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix} \in \mathbb{R}^{mp \times nq}$$

to denote the Kronecker product of matrices  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} \in \mathbb{R}^{p \times q}$ , and the notation

$$\mathbf{A}^{\otimes p} = \mathbf{A} \otimes \mathbf{A} \otimes \cdots \otimes \mathbf{A}$$

to denote the Kronecker product of  $\mathbf{A}$  with itself  $p$  times.

## 2. LINEAR METHODS

In this section, we bound the moments, variance, and concentration of the error of randomized iterative methods whose errors obey a sequential linear relationship,  $\mathbf{e}_k = \mathbf{Y}_k \mathbf{e}_{k-1}$ , where  $\mathbf{Y}_k$  is sampled from a family of  $n \times n$  matrices. Specifically, in Section 2.1, we obtain new bounds based on the concentration of the moments techniques via a tensor lifting approach, and in Section 2.3, further high-probability bounds on the whole trajectory are obtained using martingale techniques.

As noted in Section 1, the important randomized Kaczmarz (RK) and randomized Gauss–Seidel (RGS) methods for solving consistent systems of linear equations fall into this category, and we specialize the results to these popular methods in Section 2.1.1 and explore these bounds empirically with numerical experiments in Sections 2.1.2, 2.1.3, and 2.1.4. Finally, simple lower bounds illustrate that our bounds have the right shape in terms of the iteration count (Section 2.2).

**2.1. Bounds on the moments, variance, and concentration of error.** Our analysis of the variance of the error for randomized linear iterative methods uses a new bound on the second-moment of the error. It was shown in [1] (also see [44, 2]) that an exact expression for the mean squared error of the randomized Kaczmarz algorithm can be written using the matrix Kronecker product by using a “tensor lifting” trick. Here, we show that this approach can be generalized and used to generate bounds on the higher-order even moments of the error of linear randomized iterative methods by using the well-known “kernel trick”.

**Theorem 2.1.** Let  $\mathbf{e}_k = \mathbf{Y}_k \mathbf{e}_{k-1}$  where  $\mathbf{Y}_k \sim D_k$ . Denote  $\mathbb{E}_i$  to be the expectation, conditional on the choices of  $\mathbf{Y}_1, \dots, \mathbf{Y}_i$ . Then, for any  $p = 1, 2, \dots$ , we have

$$\prod_{i=1}^k \lambda_{\min}(\mathbb{E}_{i-1}[(\mathbf{Y}_i^\top \mathbf{Y}_i)^{\otimes p}]) \cdot \|\mathbf{e}_0\|^{2p} \leq \mathbb{E}\|\mathbf{e}_k\|^{2p} \leq \prod_{i=1}^k \|\mathbb{E}_{i-1}[(\mathbf{Y}_i^\top \mathbf{Y}_i)^{\otimes p}]\| \cdot \|\mathbf{e}_0\|^{2p}.$$

In particular, if the  $\mathbf{Y}_k \sim \mathbf{Y}$  are identically distributed and we define  $\mu_p := \|\mathbb{E}[(\mathbf{Y}^\top \mathbf{Y})^{\otimes p}]\|$  and  $\eta_p := \lambda_{\min}(\mathbb{E}[(\mathbf{Y}^\top \mathbf{Y})^{\otimes p}])$ , then

$$\eta_p^k \cdot \|\mathbf{e}_0\|^{2p} \leq \mathbb{E}\|\mathbf{e}_k\|^{2p} \leq \mu_p^k \cdot \|\mathbf{e}_0\|^{2p}.$$

*Proof.* First, we note that

$$(9) \quad \|\mathbf{e}_k\|^{2p} = \langle \mathbf{Y}_k^\top \mathbf{Y}_k \mathbf{e}_{k-1}, \mathbf{e}_{k-1} \rangle^p = \langle (\mathbf{Y}_k^\top \mathbf{Y}_k \mathbf{e}_{k-1})^{\otimes p}, \mathbf{e}_{k-1}^{\otimes p} \rangle$$

by the kernel trick (e.g., [43, Exercise 3.7.4]). Recall that  $\mathbb{E}_{k-1}$  denotes the expectation operator conditioned on the choices of  $\mathbf{Y}_1, \dots, \mathbf{Y}_{k-1}$ . By the law of total expectation, we obtain

$$(10) \quad \begin{aligned} \mathbb{E}\|\mathbf{e}_k\|^{2p} &= \mathbb{E}[\mathbb{E}_{k-1} \langle (\mathbf{Y}_k^\top \mathbf{Y}_k \mathbf{e}_{k-1})^{\otimes p}, \mathbf{e}_{k-1}^{\otimes p} \rangle] \\ &= \mathbb{E}[\langle \mathbb{E}_{k-1}[(\mathbf{Y}_k^\top \mathbf{Y}_k)^{\otimes p}] \mathbf{e}_{k-1}^{\otimes p}, \mathbf{e}_{k-1}^{\otimes p} \rangle], \end{aligned}$$

where the second equality follows from the mixed-product property of the Kronecker product. Hence, by applying the min-max variational theorem for the positive semidefinite matrix  $\mathbb{E}_{k-1}[(\mathbf{Y}_k^\top \mathbf{Y}_k)^{\otimes p}]$  in (10), we deduce that the following upper bound holds:

$$\mathbb{E}\|\mathbf{e}_k\|^{2p} \leq \mathbb{E}[\|\mathbb{E}_{k-1}[(\mathbf{Y}_k^\top \mathbf{Y}_k)^{\otimes p}]\| \cdot \|\mathbf{e}_{k-1}\|^{2p}] = \|\mathbb{E}_{k-1}[(\mathbf{Y}_k^\top \mathbf{Y}_k)^{\otimes p}]\| \cdot \mathbb{E}\|\mathbf{e}_{k-1}\|^{2p}.$$

Similarly, the following lower bound holds:

$$\mathbb{E}\|\mathbf{e}_k\|^{2p} \geq \lambda_{\min}(\mathbb{E}_{k-1}[(\mathbf{Y}_k^\top \mathbf{Y}_k)^{\otimes p}]) \cdot \mathbb{E}\|\mathbf{e}_{k-1}\|^{2p}.$$

The result follows by iterating these bounds and applying the law of total expectation.  $\square$

While Theorem 2.1 provides a bound on every even moment of the squared-error, our focus will be on  $\mu_2$  and  $\eta_1$ , so we drop the subscript and simply write  $\mu := \mu_2$  and  $\eta := \eta_1$ . In particular, one may use Theorem 2.1 in this case to prove Theorem 1.2.

*Proof of Theorem 1.2.* Recall that the variance of  $\|\mathbf{e}_k\|^2$  is given by  $\text{Var}(\|\mathbf{e}_k\|^2) = \mathbb{E}\|\mathbf{e}_k\|^4 - (\mathbb{E}\|\mathbf{e}_k\|^2)^2$ . By applying Theorem 2.1 with  $p = 2$ , we obtain the upper bound

$$\mathbb{E}\|\mathbf{e}_k\|^4 \leq \mu^k \cdot \|\mathbf{e}_0\|^4$$

with  $\mu = \|\mathbb{E}[(\mathbf{Y}^\top \mathbf{Y})^{\otimes 2}]\|$ . Similarly, by applying the same result with  $p = 1$ , we obtain the lower bound

$$\mathbb{E}\|\mathbf{e}_k\|^2 \geq \eta^k \cdot \|\mathbf{e}_0\|^2$$

with  $\eta = \lambda_{\min}(\mathbb{E}[\mathbf{Y}^\top \mathbf{Y}])$ . Combining these bounds implies that  $\text{Var}(\|\mathbf{e}_k\|^2) \leq (\mu^k - \eta^k) \cdot \|\mathbf{e}_0\|^4$ .  $\square$

**Remark 1.** The bound on the variance from Theorem 1.2 immediately implies the following concentration result by applying Chebyshev's inequality:

$$(11) \quad \mathbb{P}(|\|\mathbf{e}_k\|^2 - \mathbb{E}\|\mathbf{e}_k\|^2| \geq t) \leq \frac{\mu^k - \eta^k}{t^2}.$$

In particular, this implies that for any  $\epsilon \in (0, 1)$ , we have

$$(12) \quad \mathbb{P}\left(|\|\mathbf{e}_k\|^2 - \mathbb{E}\|\mathbf{e}_k\|^2| \geq \sqrt{\frac{\mu^k - \eta^k}{\epsilon}} \|\mathbf{e}_0\|^2\right) \leq \epsilon.$$

Hence, with probability at least  $1 - \epsilon$ , the squared error norm lies in the interval  $\mathbb{E}\|\mathbf{e}_k\|^2 \pm \sqrt{(\mu^k - \eta^k)\epsilon^{-1}} \|\mathbf{e}_0\|^2$ .

**Remark 2.** Since Theorem 2.1 does not, in general, assume that the  $\mathbf{Y}_k$  are sampled independently or identically, the bounds are applicable if one is able to estimate bounds for the conditional expectations

$$\mu^{(j)} \geq \|\mathbb{E}_{j-1}[(\mathbf{Y}_j^\top \mathbf{Y}_j)^{\otimes 2}]\| \quad \text{and} \quad \eta^{(j)} \leq \lambda_{\min}(\mathbb{E}_{j-1}[\mathbf{Y}_j^\top \mathbf{Y}_j]).$$

This generalization is relevant for the important case of randomized Kaczmarz or randomized Gauss–Seidel with iteration-dependent step-sizes, or for hybrid greedy and random sampling techniques, like in [32, 22, 5].

Next, we specify the results obtained by applying Theorem 1.2 to commonly studied randomized linear iterative methods.

**2.1.1. Randomized Kaczmarz and randomized Gauss–Seidel.** The *randomized Kaczmarz (RK)* methods are members of the family of Kaczmarz methods, classical examples of *row-action* iterative methods. These methods consist of sequential orthogonal projections towards the solution set of a single equation [24]; the  $j$ th iterate is recursively defined as

$$(13) \quad \mathbf{x}_j = \mathbf{x}_{j-1} - \frac{\mathbf{a}_{i_j}^\top \mathbf{x}_{j-1} - b_{i_j}}{\|\mathbf{a}_{i_j}\|^2} \mathbf{a}_{i_j},$$

where  $\mathbf{a}_{i_j}^\top$  is the  $i_j$ th row of the matrix  $\mathbf{A}$  and  $b_{i_j}$  is the  $i_j$ th entry of  $\mathbf{b}$ . As mentioned previously, the RK method for solving a linear system  $\mathbf{Ax} = \mathbf{b}$  has error vectors  $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}^*$  that satisfy the recursive relation

$$\mathbf{e}_j = \mathbf{x}_{j-1} - \mathbf{x}^* - \frac{\mathbf{a}_{i_j}^\top \mathbf{x}_{j-1} - \mathbf{a}_{i_j}^\top \mathbf{x}^*}{\|\mathbf{a}_{i_j}\|^2} \mathbf{a}_{i_j} = \left( \mathbf{I} - \frac{\mathbf{a}_{i_j} \mathbf{a}_{i_j}^\top}{\|\mathbf{a}_{i_j}\|^2} \right) (\mathbf{x}_{j-1} - \mathbf{x}^*) = \mathbf{Y}_j \mathbf{e}_{j-1},$$

where  $\mathbf{Y}_j = \mathbf{I} - \mathbf{a}_{i_j} \mathbf{a}_{i_j}^\top / \|\mathbf{a}_{i_j}\|^2$  is a random orthogonal projection matrix corresponding to a projection of the error onto the subspace orthogonal to row  $\mathbf{a}_{i_j}$ . The RK methods saw a renewed surge of interest after the elegant convergence analysis of the RK method in [42]. The authors showed that for a consistent system with unique solution  $\mathbf{x}^*$ , if the row  $i_j$  is sampled with probability  $\|\mathbf{a}_{i_j}\|^2 / \|\mathbf{A}\|_F^2$  in each iteration, then RK converges at least linearly in expectation with the guarantee (1).

The *randomized Gauss–Seidel (RGS)* methods are a related family of *column-action* iterative methods that focus on updating a single coordinate (or subset of coordinates) in each iteration to minimize the residual error; see e.g., [31]. The  $j$ th iterate is recursively defined as

$$(14) \quad \mathbf{x}_j = \mathbf{x}_{j-1} - \frac{\mathbf{A}_{i_j}^\top (\mathbf{Ax}_{j-1} - \mathbf{b})}{\|\mathbf{A}_{i_j}\|^2} \mathbf{c}_{i_j},$$

where  $\mathbf{A}_{i_j}$  is the  $i_j$ th column of  $\mathbf{A}$  and  $\mathbf{c}_{i_j}$  is the  $i_j$ th standard basis vector. The RGS residual errors  $\mathbf{e}_k = \mathbf{Ax}_k - \mathbf{Ax}^*$  satisfy the recursive relation

$$\mathbf{e}_j = \mathbf{Ax}_{j-1} - \mathbf{Ax}^* - \frac{\mathbf{A}_{i_j}^\top (\mathbf{Ax}_{j-1} - \mathbf{Ax}^*)}{\|\mathbf{A}_{i_j}\|^2} \mathbf{Ac}_{i_j} = \left( \mathbf{I} - \frac{\mathbf{A}_{i_j} \mathbf{A}_{i_j}^\top}{\|\mathbf{A}_{i_j}\|^2} \right) \mathbf{A} (\mathbf{x}_{j-1} - \mathbf{x}^*) = \mathbf{Y}_j \mathbf{e}_{j-1},$$

where  $\mathbf{Y}_j = \mathbf{I} - \mathbf{A}_{i_j} \mathbf{A}_{i_j}^\top / \|\mathbf{A}_{i_j}\|^2$  is a randomly sampled projection matrix corresponding to a projection of the residual error onto the subspace orthogonal to column  $\mathbf{A}_{i_j}$ . It was shown in [27] that for a consistent system, if the column  $i_j$  is sampled with probability  $\|\mathbf{A}_{i_j}\|^2 / \|\mathbf{A}\|_F^2$  in each iteration, then RGS also converges at least linearly in expectation with the same guarantee (1) for the residual error  $\mathbf{e}_k = \mathbf{Ax}_k - \mathbf{b}$ .

Note that in both of these cases, the  $\mathbf{Y}_k$  matrices are quite nice: they are independent and identically distributed copies of a random orthogonal projection matrix  $\mathbf{Y}$  ( $\mathbf{Y}^\top \mathbf{Y} = \mathbf{Y}^2 = \mathbf{Y}$ ) that is positive semidefinite ( $\mathbf{Y} \succeq \mathbf{0}$ ) and a contraction ( $\|\mathbf{Y}\| \leq 1$ ). Moreover, for RK, we have the closed-form expression

$$\mathbb{E}[\mathbf{Y}] = \mathbf{I} - \sum_{i=1}^m \frac{\|\mathbf{a}_i\|^2}{\|\mathbf{A}\|_F^2} \frac{\mathbf{a}_i \mathbf{a}_i^\top}{\|\mathbf{a}_i\|^2} = \mathbf{I} - \frac{\mathbf{A}^\top \mathbf{A}}{\|\mathbf{A}\|_F^2}.$$

Similarly, for RGS, we have

$$\mathbb{E}[\mathbf{Y}] = \mathbf{I} - \sum_{j=1}^n \frac{\|\mathbf{A}_j\|^2}{\|\mathbf{A}\|_F^2} \frac{\mathbf{A}_j \mathbf{A}_j^\top}{\|\mathbf{A}_j\|^2} = \mathbf{I} - \frac{\mathbf{A} \mathbf{A}^\top}{\|\mathbf{A}\|_F^2}.$$

Note that by [2, Theorem 4.2], when  $m \geq n$  and  $\mathbf{A}$  is full rank,

$$(15) \quad \mu = \|\mathbb{E}[(\mathbf{Y}^\top \mathbf{Y})^{\otimes 2}]\| \leq 1 - \frac{\sigma_{\min}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}.$$

This follows from the observation that  $\mathbf{Y} \succeq \mathbf{0}$  and  $\mathbf{I} - \mathbf{Y} \succeq \mathbf{0}$  almost surely (because  $\mathbf{Y}$  are positive semidefinite contraction matrices), which implies that  $\mathbf{I} \otimes \mathbf{Y} - \mathbf{Y} \otimes \mathbf{Y} = (\mathbf{I} - \mathbf{Y}) \otimes \mathbf{Y} \succeq \mathbf{0}$  and hence

$$\mu = \|\mathbb{E}[\mathbf{Y} \otimes \mathbf{Y}]\| \leq \|\mathbb{E}[\mathbf{I} \otimes \mathbf{Y}]\| = \|\mathbb{E}[\mathbf{Y}]\| = 1 - \frac{\sigma_{\min}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}.$$

Moreover,

$$(16) \quad \eta = \lambda_{\min}(\mathbb{E}[\mathbf{Y}^\top \mathbf{Y}]) = \lambda_{\min}(\mathbb{E}[\mathbf{Y}]) = \lambda_{\min}\left(\mathbf{I} - \frac{\mathbf{A}^\top \mathbf{A}}{\|\mathbf{A}\|_F^2}\right) = 1 - \frac{\sigma_{\max}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}.$$

Thus, applying Theorem 1.2, we may bound the variance of the squared norm of the error of both the RK and RGS methods by

$$(17) \quad \text{Var}(\|\mathbf{e}_k\|^2) \leq \left( \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}\right)^k - \left(1 - \frac{\sigma_{\max}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}\right)^k \right) \cdot \|\mathbf{e}_0\|^4.$$

**Remark 3.** Using the same notation as Theorem 1.2, note that (15) can be generalized for the RK and RGS methods to show that for all  $p \in \mathbb{N}$ ,

$$\mu_p \leq \mu_{p-1} \leq \dots \leq \mu_2 \leq \|\mathbb{E}\mathbf{Y}_k\| = 1 - \frac{\sigma_{\min}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2},$$

since

$$\mu_p = \|\mathbb{E}[\underbrace{\mathbf{Y} \otimes \dots \otimes \mathbf{Y}}_{p \text{ times}}]\| \leq \|\mathbb{E}[\underbrace{\mathbf{I} \otimes \mathbf{Y} \otimes \dots \otimes \mathbf{Y}}_{p-1 \text{ times}}]\| = \|\mathbb{E}[\underbrace{\mathbf{Y} \otimes \dots \otimes \mathbf{Y}}_{p-1 \text{ times}}]\| = \mu_{p-1}.$$

2.1.2. *The  $\mu$  parameter.* For the RK and RGS methods, we are able to show that

$$\mu = \|\mathbb{E}[(\mathbf{Y}^\top \mathbf{Y})^{\otimes 2}]\| \leq 1 - \frac{\sigma_{\min}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2} =: r,$$

and this bound yields (17). Thus, we know that  $\log_r(\mu) \geq 1$ , and the bound on the variance improves as  $\log_r(\mu)$  increases. We explore the value of  $\log_r(\mu)$  in Figure 2. We compute the parameter

$$(18) \quad \mu = \left\| \sum_{i=1}^m \frac{\|\mathbf{a}_i\|^2}{\|\mathbf{A}\|_F^2} (\mathbf{Y}_i)^{\otimes 2} \right\| \quad \text{with } \mathbf{Y}_i = \mathbf{I} - \frac{\mathbf{a}_i \mathbf{a}_i^\top}{\|\mathbf{a}_i\|^2}$$

for row-normalized Gaussian matrices of various sizes with entries generated independently from  $\mathcal{N}(0, 10)$ , and calculate  $\log_r(\mu)$ . We know by (15) that  $\mu \leq r$  when  $m \geq n$ . We observe that  $\mu$  tends towards  $r$  as the system becomes more overdetermined towards to the lower left corner of the heatmap, but that we have  $\mu \approx r^2$  for systems that are approximately square (i.e., near the diagonal). When  $m < n$ , we have  $r < \mu = 1$ , and thus  $\log_r(\mu) = 0$  in the upper right half of the heatmap.

2.1.3. *Empirical performance of concentration bound.* Given the concentration bounds derived by combining Theorem 1.2 and Chebyshev's inequality, we expect the errors of RK to be more concentrated around their mean when the matrix is well-conditioned. We explore this empirically in the left plots of Figure 3. We plot the empirical squared error of RK over 100 iterations across 500 runs. We color these error curves according to a gradient which indicates what fraction of trials had error below that of the given trial (more blue means more trials had smaller error). We illustrate the bound on the mean error given by (1) (black solid lines) and plot the empirical mean error (white dashed lines). Finally, we plot the 75% (red dashed lines) and 95% (green dashed lines) confidence intervals derived from combining Theorem 1.2 with (4) and Chebyshev's inequality. We note that the confidence interval is centered at the empirical mean of the errors across the trials.

In Figure 3, we generate a matrix  $\mathbf{A} \in \mathbb{R}^{1,000 \times 20}$  with the singular values plotted in the corresponding right plots by generating a matrix  $\tilde{\mathbf{A}} \in \mathbb{R}^{1,000 \times 20}$  with entries sampled i.i.d. from  $\mathcal{N}(0, 1)$ , computing the singular value decomposition  $\tilde{\mathbf{A}} = \mathbf{U}\tilde{\Sigma}\mathbf{V}^\top$ , and defining  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$  where the entries of  $\Sigma$  on the diagonal are as specified in Subfigures 3a and 3c and plotted in the corresponding right plots. In the top row of plots in Figure 3, the experiments are run with  $\mathbf{A}_1$  where  $\sigma_i(\mathbf{A}_1) = 1 - (i - 1)/m$ ; in the middle row of plots, we have  $\mathbf{A}$  with entries sampled i.i.d. from  $\mathcal{N}(0, 1)$ ; and in the bottom row of plots, we have  $\sigma_i(\mathbf{A}_2) = 1/i$ .

We note that, like the bound (1), the concentration bound derived from combining Theorem 1.2 with (4) and Chebyshev's inequality is quite sensitive to the conditioning of the problem-defining matrix.

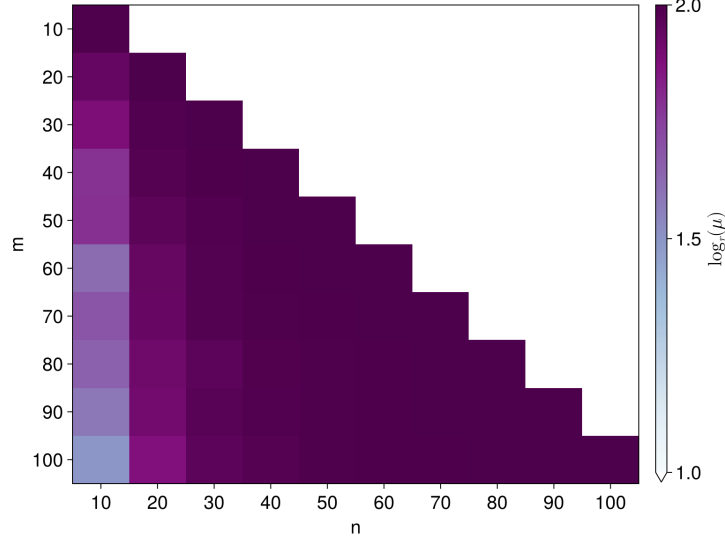


FIGURE 2. The relationship between  $\mu$  and the RK convergence rate,  $r = 1 - \frac{\sigma_{\min}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}$ . For each cell, we initialized five  $m \times n$  row-normalized Gaussian matrices. We compute the  $\mu$  parameter as in (18) and plot the average value of  $\log_r(\mu)$  across the five trials in each cell. When  $n > m$ , we note that  $\log_r(\mu) = 0$ .

**2.1.4. Comparison of concentration bounds.** We now compare the values of the concentration bound offered by Theorem 1.2 for the RK and RGS methods (top), Lemma 1.1 (2) (middle), and the concentration bound (8) (bottom) which is a consequence of [21, Theorem 7.1], for a variety of choices of constant  $t$  and iteration number  $k$ . In these plots, we generate matrices as described in Subsection 2.1.3. We define a well-conditioned matrix  $\mathbf{A}_1 \in \mathbb{R}^{1,000 \times 20}$  (condition number  $\kappa(\mathbf{A}_1) = 1.02$ ), a Gaussian matrix  $\mathbf{A} \in \mathbb{R}^{1,000 \times 20}$  (condition number  $\kappa(\mathbf{A}) = 1.29$ ), and an ill-conditioned matrix  $\mathbf{A}_2 \in \mathbb{R}^{1,000 \times 20}$  (condition number  $\kappa(\mathbf{A}_2) = 20.00$ ). We also take the matrix  $\mathbf{A}_3 \in \mathbb{R}^{1,200 \times 400}$  from a 2D tomography test problem (condition number  $\kappa(\mathbf{A}_3) = 21.53$ ), generated using the Matlab Regularization Toolbox by P.C. Hansen (<http://www.imm.dtu.dk/~pcha/Regutools/>) [17] with  $m = fN^2$  and  $n = N^2$ , where  $N = 20$  and the oversampling factor  $f = 3$ . In Figure 4, we present these heatmaps for the well-conditioned matrix on the top left plot, for the Gaussian matrix on the top right plot, for the ill-conditioned matrix on the bottom left plot, and for the matrix arising from a computed tomography test problem on the bottom right plot.

We note that the bound (8), derived from [21], is worse relative to our concentration bound, derived from Theorem 1.2, as well as the bound from Lemma 1.1 (2). We also observe that in most regimes of  $t$  and  $k$ , our concentration bound derived from Theorem 1.2 is an improvement over the bound derived from Lemma 1.1 (2) using Markov's inequality.

**2.2. Lower bound on concentration.** Note that for the RK method, the random variable  $\|\mathbf{e}_k\|^2$  only takes on a finite set of values. Since  $\mathbf{e}_k = \mathbf{Y}_k \mathbf{e}_{k-1}$  and  $\mathbf{Y}_k$  is sampled i.i.d. from the set  $\{\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_m\}$  at iteration  $k$ , there are at most  $m^k$  possible values for  $\|\mathbf{e}_k\|^2$ , corresponding to each possible sequence of sampled indices. Moreover, observe that we have

$$(19) \quad \mathbf{e}_k = \mathbf{Y}_k \mathbf{e}_{k-1} = \mathbf{Y}_k \mathbf{Y}_{k-1} \mathbf{e}_{k-2} = \mathbf{Y}_{k-1} \mathbf{e}_{k-2} = \mathbf{e}_{k-1} \quad \text{if } i_k = i_{k-1}.$$

To simplify calculations, we assume in this section that  $\mathbf{A}$  is row-normalized, that is  $\|\mathbf{a}_i\| = 1$  for all  $i \in [m]$ . Thus, since the norm of the error  $\|\mathbf{e}_k\|^2$  is non-increasing, by considering the event that  $i_1 = i_2 = \dots = i_k$ , we deduce that

$$\mathbb{P}[\|\mathbf{e}_k\|^2 \geq \|\mathbf{e}_1\|^2] \geq \frac{1}{m} \sum_{j=1}^m \mathbb{P}[i_k = i_{k-1} = \dots = i_2 = i_1 \mid i_1 = j] = \frac{1}{m^{k-1}}.$$

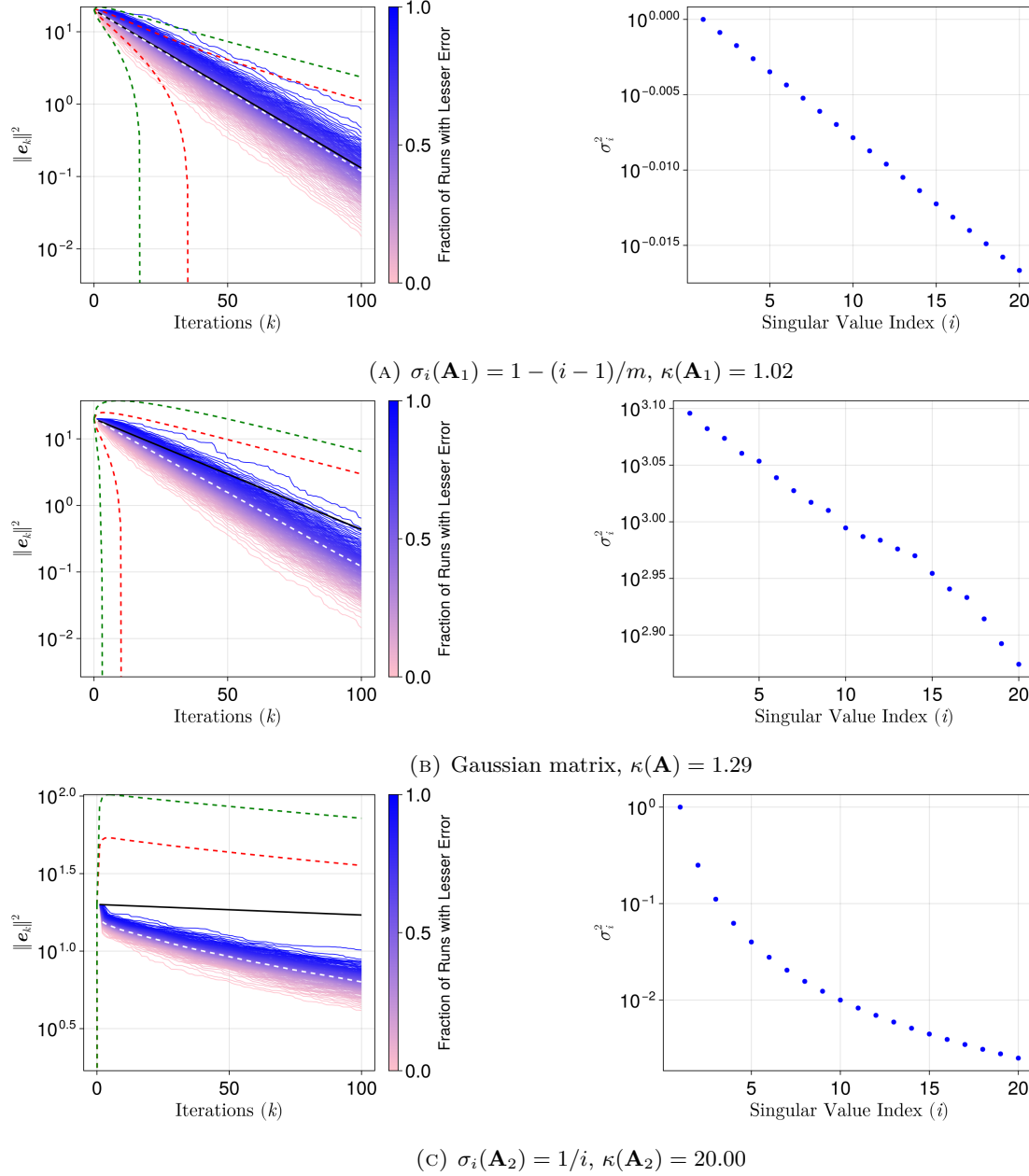


FIGURE 3. (Left) Visualization of errors of 500 independent trials of RK applied to  $\mathbf{A}\mathbf{x} = \mathbf{b}$  where  $\mathbf{A} \in \mathbb{R}^{1000 \times 20}$  has singular values given in subfigure captions (pink-blue gradient indicates quantiles of errors). Empirical mean error (white dashed line), bound (1) (black solid line), and the 75% (red dashed lines) and 95% (green dashed lines) confidence intervals for the error derived by combining Chebyshev's inequality with Theorem 1.2 and (4) are plotted. (Right) Spectral profile for  $\mathbf{A}$ .

Now, for sufficiently small values of  $t$  where  $t \leq \min_{i_1} \|\mathbf{e}_1\|^2 - \mathbb{E}\|\mathbf{e}_k\|^2$ , we have that

$$(20) \quad \mathbb{P}[\|\mathbf{e}_k\|^2 - \mathbb{E}\|\mathbf{e}_k\|^2 \geq t] \geq \mathbb{P}[\|\mathbf{e}_k\|^2 \geq \|\mathbf{e}_1\|^2] \geq \frac{1}{m^{k-1}}.$$

Thus, we note that bounds for the concentration of  $\|\mathbf{e}_k\|^2$  must respect this lower bound.

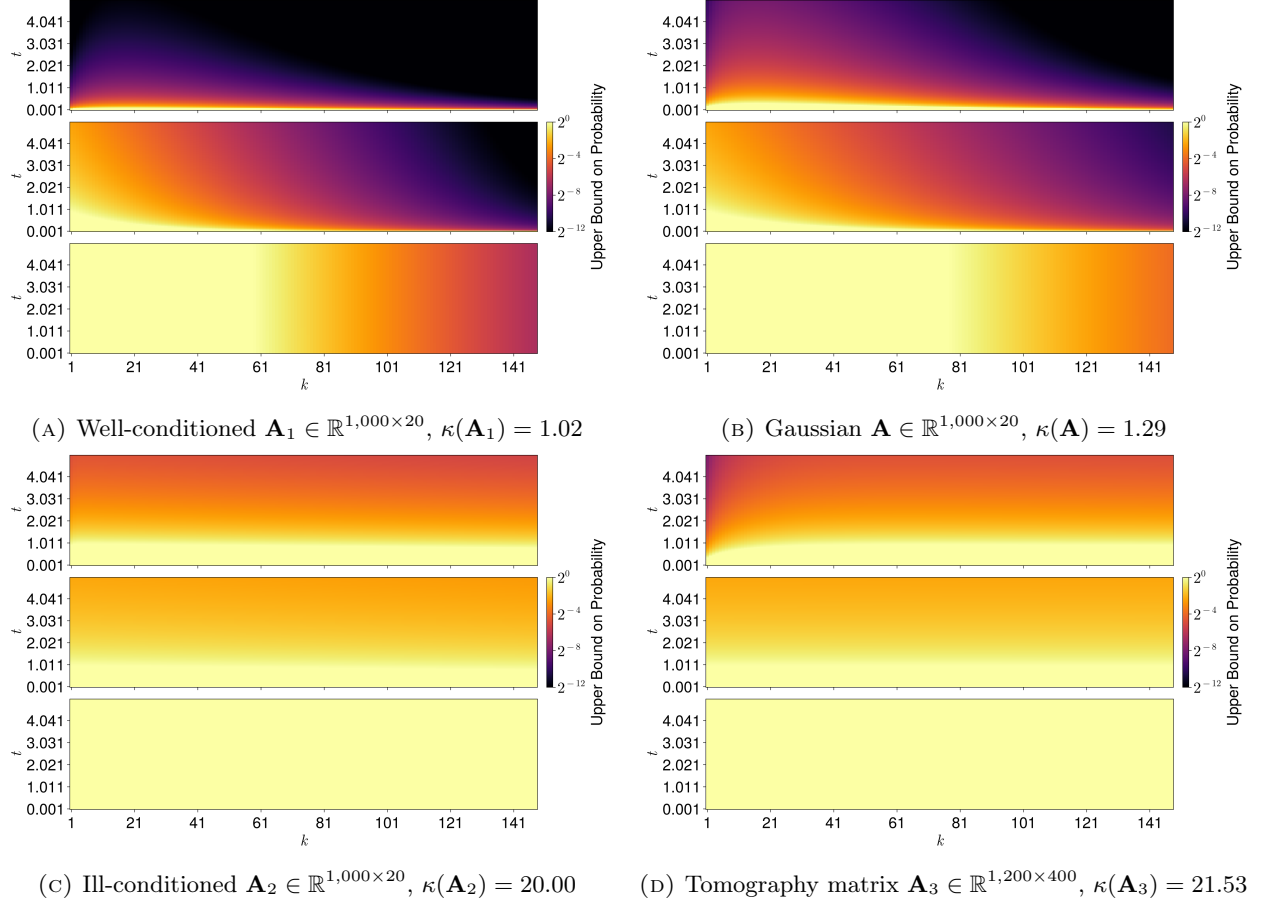


FIGURE 4. Comparison of the upper bounds on the probability  $\mathbb{P}(\|\mathbf{e}_k\|^2 - \mathbb{E}\|\mathbf{e}_k\|^2 \geq t)$ , for various  $k$  and  $t$ , resulting from Theorem 1.2 (top, our contribution), from Lemma 1.1 (2) (middle, simple Markov inequality bound), and from (8) which is a consequence of [21, Theorem 7.1] (bottom, matrix concentration) for RK applied to a various matrices. Each cell corresponds to the upper bound on the probability of the squared-error exceeding its mean by  $t$  after  $k$  iterations. Darker cells correspond to smaller values, i.e., better concentration bounds.

We note that the same logic holds for the RGS method. Again, to simplify calculations, we assume that  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is column-normalized. In this case, for  $t \leq \min_{i_1} \|\mathbf{e}_1\|^2 - \mathbb{E}\|\mathbf{e}_k\|^2$ , we have

$$(21) \quad \mathbb{P}[\|\mathbf{e}_k\|^2 - \mathbb{E}\|\mathbf{e}_k\|^2 \geq t] \geq \mathbb{P}[\|\mathbf{e}_k\|^2 \geq \|\mathbf{e}_1\|^2] \geq \frac{1}{n^{k-1}}.$$

**Remark 4.** We note that the lower bound (20) implies that any upper bound for the concentration of the error of the RK and RGS methods cannot decrease with the number of iterations  $k$  faster than  $e^{-O(k)}$ , where  $O(k) = Ck$  for some constant  $C$ . This implies that the concentration bound (11) cannot be improved with respect to  $k$  beyond constants.

**2.3. High-probability bounds.** In this subsection, we prove Theorem 1.3 that provides high-probability bound for randomized iterative methods whose errors in sequential iterations obey a linear relationship,  $\mathbf{e}_k = \mathbf{Y}_k \mathbf{e}_{k-1}$ , and  $\mathbf{Y}_k$  is independently sampled in the  $k$ th iteration from a family of  $n \times n$  matrices. Part (a) shows that we can upgrade the upper bound implied by Markov's inequality for any fixed iteration number  $k$  to the entire trajectory.

*Proof of Theorem 1.3, Part (a).* Consider the discrete stochastic process

$$Z_k := \frac{\|\mathbf{e}_k\|^2}{\rho^k \|\mathbf{e}_0\|^2}.$$

Note that if  $\mathbb{E}_{k-1}$  denotes the expectation conditional on  $\mathbf{Y}_1, \dots, \mathbf{Y}_{k-1}$ , then

$$(22) \quad \mathbb{E}_{k-1} \|\mathbf{e}_k\|^2 = \mathbb{E}_{k-1} [\mathbf{e}_{k-1}^\top \mathbf{Y}_k^\top \mathbf{Y}_k \mathbf{e}_{k-1}] = \mathbf{e}_{k-1}^\top \mathbb{E}[\mathbf{Y}_k^\top \mathbf{Y}_k] \mathbf{e}_{k-1} \leq \rho \|\mathbf{e}_{k-1}\|^2.$$

Thus,  $Z_k$  is a non-negative supermartingale:  $\mathbb{E}_{k-1}[Z_k] \leq Z_{k-1}$ . Applying Doob's supermartingale inequality (e.g., [7, Exercise 4.8.2]) implies that for all  $\lambda > 0$ ,

$$\mathbb{P} \left[ \sup_{k \geq 0} Z_k \geq \lambda \right] \leq \frac{\mathbb{E}[Z_0]}{\lambda} = \frac{1}{\lambda}.$$

Choosing  $\lambda = \epsilon^{-1}$  and rearranging leads to the claimed result.  $\square$

*Proof of Theorem 1.3, Part (b).* Write

$$\frac{\|\mathbf{e}_k\|^2}{\|\mathbf{e}_0\|^2} = \prod_{j=1}^k \frac{\|\mathbf{e}_j\|^2}{\|\mathbf{e}_{j-1}\|^2}.$$

By taking logarithms and using the inequality  $\log(1+x) \leq x$ , which holds for all  $x \geq -1$ , we deduce that

$$\log \left( \frac{\|\mathbf{e}_k\|^2}{\|\mathbf{e}_0\|^2} \right) = \sum_{j=1}^k \log \left( 1 + \frac{\|\mathbf{e}_j\|^2}{\|\mathbf{e}_{j-1}\|^2} - 1 \right) \leq \sum_{j=1}^k \xi_j, \quad \text{where } \xi_j = \frac{\|\mathbf{e}_j\|^2}{\|\mathbf{e}_{j-1}\|^2} - 1.$$

Now, consider the process

$$\tilde{Z}_k := \sum_{j=1}^k (\xi_j + (1 - \rho)).$$

If  $\mathbb{E}_{k-1}$  denotes the expectation conditional on  $\mathbf{Y}_1, \dots, \mathbf{Y}_{k-1}$ , then

$$\mathbb{E}_{k-1} \|\mathbf{e}_k\|^2 = \mathbf{e}_{k-1}^\top \mathbb{E}[\mathbf{Y}_k^\top \mathbf{Y}_k] \mathbf{e}_{k-1} \leq \|\mathbb{E}[\mathbf{Y}_k^\top \mathbf{Y}_k]\| \cdot \|\mathbf{e}_{k-1}\|^2 \leq \rho \|\mathbf{e}_{k-1}\|^2.$$

This implies that  $\mathbb{E}_{k-1}[\xi_k] \leq \rho - 1$ , and hence

$$\mathbb{E}_{k-1} \tilde{Z}_k = \tilde{Z}_{k-1} + \mathbb{E}_{k-1} [\xi_k + (1 - \rho)] \leq \tilde{Z}_{k-1}.$$

That is,  $\tilde{Z}_k$  is a supermartingale, null at zero, with respect to the natural filtration. Moreover, from the almost sure boundedness assumption,

$$\|\mathbf{e}_k\|^2 = \mathbf{e}_{k-1}^\top \mathbf{Y}_k^\top \mathbf{Y}_k \mathbf{e}_{k-1} \leq \alpha \|\mathbf{e}_{k-1}\|^2.$$

This shows that the process  $\tilde{Z}_k$  has bounded increments, since

$$\tilde{Z}_k - \tilde{Z}_{k-1} = \frac{\|\mathbf{e}_k\|^2}{\|\mathbf{e}_{k-1}\|^2} - 1 + (1 - \rho) \in [-\rho, \alpha - \rho].$$

By the Azuma-Hoeffding inequality for supermartingales with bounded differences (e.g., [9, Corollary 2.1]), this implies that for all  $\lambda \geq 0$ ,

$$\begin{aligned} \mathbb{P} \left[ \sup_{0 \leq t \leq k} \log \left( \frac{\|\mathbf{e}_t\|^2}{\|\mathbf{e}_0\|^2} \right) + t \cdot (1 - \rho) \geq \lambda \right] &\leq \mathbb{P} \left[ \sup_{0 \leq t \leq k} \sum_{j=1}^t \xi_j + t \cdot (1 - \rho) \geq \lambda \right] \\ &= \mathbb{P} \left[ \sup_{0 \leq t \leq k} \tilde{Z}_t \geq \lambda \right] \leq \exp \left( \frac{-\lambda^2}{2k\alpha^2} \right). \end{aligned}$$

Choosing  $\lambda = \alpha \sqrt{2k \log(\epsilon^{-1})}$  in particular, we deduce that with probability at least  $1 - \epsilon$ , we have for all  $0 \leq t \leq k$  simultaneously,

$$\log \left( \frac{\|\mathbf{e}_t\|^2}{\|\mathbf{e}_0\|^2} \right) \leq -t \cdot (1 - \rho) + \alpha \cdot \sqrt{2k \log(\epsilon^{-1})}.$$

Rearranging leads to the claimed result.  $\square$

## 3. NONLINEAR METHODS

As we have seen above, the available convergence guarantees for many iterative methods are typically of the form  $\mathbb{E}[d(\mathbf{x}_k, S)^2] \leq r^k d(\mathbf{x}_0, S)^2$ , where  $d$  measures some distance to the solution set. However, some lack the linear structure that enable the moment and variance bounds provided above. In this section, we bound the concentration and variance of the error of randomized iterative methods whose errors in sequential iterations do not necessarily obey a linear relationship. As described in Section 1, the variant of randomized Kaczmarz (RK) for solving consistent systems of linear *inequalities* falls into this category.

**3.1. Bounds on the variance and concentration of error.** We will prove Theorem 1.4, which provides a bound on the variance of the squared error, and can be combined with Chebyshev's inequality to yield a bound on the concentration.

*Proof of Theorem 1.4.* To begin, we observe that:

$$\begin{aligned} \text{Var}(d(\mathbf{x}_k, S)^2) &= \mathbb{E}[d(\mathbf{x}_k, S)^4] - (\mathbb{E}[d(\mathbf{x}_k, S)^2])^2 \\ &\leq \mathbb{E}[d(\mathbf{x}_k, S)^4] = \mathbb{E}[d(\mathbf{x}_k, S)^2 \cdot d(\mathbf{x}_k, S)^2]. \end{aligned}$$

Then, we invoke the bound  $d(\mathbf{x}_k, S) \leq D$  to bound

$$\text{Var}(d(\mathbf{x}_k, S)^2) \leq D^2 \mathbb{E}[d(\mathbf{x}_k, S)^2] \leq D^2 r^k d(\mathbf{x}_0, S)^2.$$

This completes the proof.  $\square$

**Remark 5.** *The bound on the variance immediately implies the following concentration result by applying Chebyshev's inequality:*

$$\mathbb{P}(|\|\mathbf{e}_k\|^2 - \mathbb{E}\|\mathbf{e}_k\|^2| \geq t) \leq \frac{D^2 r^k d(\mathbf{x}_0, S)^2}{t^2}.$$

*In particular, this implies that for any  $\epsilon \in (0, 1)$ , we have*

$$(23) \quad \mathbb{P}\left(|d(\mathbf{x}_k, S)^2 - \mathbb{E}[d(\mathbf{x}_k, S)^2]| \geq 2D\sqrt{\frac{r^k}{\epsilon}} d(\mathbf{x}_0, S)\right) \leq \epsilon.$$

*Hence, with probability at least  $1 - \epsilon$ , the squared distance to  $S$  lies in the interval  $\mathbb{E}[d(\mathbf{x}_k, S)^2] \pm 2D\sqrt{r^k \epsilon^{-1}} d(\mathbf{x}_0, S)$ .*

Next, we specify the results obtained by applying Theorem 1.4 to a commonly studied randomized linear iterative method for linear feasibility problems.

**3.1.1. Randomized Kaczmarz method for linear feasibility.** As a generalization of the RK methods for linear equations, Leventhal and Lewis [27] proposed a randomized algorithm for solving linear feasibility problems of the form

$$(24) \quad \begin{cases} \mathbf{a}_i^\top \mathbf{x} \leq b_i & (i \in I_{\leq}) \\ \mathbf{a}_i^\top \mathbf{x} = b_i & (i \in I_{=}), \end{cases}$$

where  $\mathbf{a}_i \in \mathbb{R}^n$  for each  $i$ , and the disjoint index sets  $I_{\leq}$  and  $I_{=}$  partition the set  $\{1, 2, \dots, m\}$ . At each iteration  $j$ , the algorithm randomly samples an index  $i_j \in I_{\leq} \cup I_{=}$ , and if  $i_j \in I_{=}$  or if  $i_j \in I_{\leq}$  and  $\mathbf{a}_{i_j}^\top \mathbf{x}_{j-1} > b_{i_j}$ , projects the current iterate  $\mathbf{x}_{j-1}$  onto the hyperplane  $\{\mathbf{x} : \mathbf{a}_{i_j}^\top \mathbf{x} = b_{i_j}\}$ . This update may be written as

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \frac{(\mathbf{a}_{i_k}^\top \mathbf{x}_{k-1} - b_{i_k})^+}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k},$$

where  $z^+ = \max\{z, 0\}$ . Note that since the update defined by a sampled index  $i_j \in I_{\leq}$  depends upon the position of the current iterate  $\mathbf{x}_{j-1}$ , this method does not satisfy the linear relationship,  $\mathbf{e}_j = \mathbf{Y}_j \mathbf{e}_{j-1}$ . In particular, the update matrix  $\mathbf{Y}_j$  cannot be sampled from a set of fixed matrices, but must depend upon  $\mathbf{e}_{j-1}$ . Thus, the results from Section 2 do not immediately apply to this benignly nonlinear method.

It was shown in [27] that this algorithm converges at least linearly in expectation, with the guarantee

$$\mathbb{E}[d(\mathbf{x}_k, S)^2] \leq \left(1 - \frac{1}{L^2 \|\mathbf{A}\|_F^2}\right)^k d(\mathbf{x}_0, S)^2,$$

if the rows are sampled with probability  $\|\mathbf{a}_{i_k}\|^2/\|\mathbf{A}\|_F^2$ , where  $\mathbf{A}$  is the  $m \times n$  matrix whose  $i$ th row is  $\mathbf{a}_i^\top$ ,  $S$  is the feasible region defined by (24),  $d(\mathbf{x}, S) := \min_{\mathbf{y} \in S} \|\mathbf{x} - \mathbf{y}\|$  denotes the Euclidean distance of a point  $\mathbf{x}$  to set  $S$ , and  $L$  is the *Hoffman constant* for the system (24) (see [20]).

Now, noting that this method satisfies  $d(\mathbf{x}_k, S) \leq D := d(\mathbf{x}_0, S)$  for all  $k \geq 0$ , we may use Theorem 1.4 (7) to bound the variance of the squared distance to the feasible set by

$$(25) \quad \text{Var}(d(\mathbf{x}_k, S)^2) \leq \left(1 - \frac{1}{L^2 \|\mathbf{A}\|_F^2}\right)^k d(\mathbf{x}_0, S)^4.$$

We note that the Hoffman constants are difficult to calculate or bound in general [38], so we do not include any experiments evaluating this bound.

#### 4. CONCLUSIONS

In this work, we establish upper bounds on the variance and concentration of the error of general classes of randomized iterative methods. While most previous analysis primarily focused on convergence in expectation, our results illustrate how the error can deviate above (and around) the expected error. For linear iterative methods like the randomized Kaczmarz and randomized Gauss–Seidel methods, we derived higher-order moment bounds using tensor-based analysis, extended these to bounds on the variance and concentration via Chebyshev’s inequality, and provided some additional martingale-based high-probability results that can simultaneously bound entire random trajectory of an algorithm. We also extended our analysis to nonlinear iterative methods, such as the randomized Kaczmarz method for solving systems of linear inequalities, demonstrating similar bounds under mild assumptions.

These theoretical contributions are supported by comprehensive numerical experiments, which illustrate the validity and usefulness of our bounds across a range of problem types, including synthetically generated matrices with varying spectral gaps and structures. Our results offer not only improved understanding of the near-worst-case behavior of stochastic iterative methods, but also practical tools for designing and evaluating algorithms with quantifiable probabilistic guarantees.

Future work includes a closer analysis of the tensor-based parameter  $\mu$  (and more generally the  $\mu_p$  parameters appearing in the higher-order moment bounds), and better understanding its relationship to the matrix  $\mathbf{A}$ . Also, in some adaptive versions of RK, such as the corruption-robust QuantileRK algorithm [16, 41], the error does not satisfy the linear relationship  $\mathbf{e}_k = \mathbf{Y}_k \mathbf{e}_{k-1}$  in every iteration, and yet expectation bounds of the type (1) are available. Extending the concentration analysis to such methods is especially interesting due to the inherent non-monotonicity of the error. Further, there has been significant interest in partially greedy row selection strategies [5, 3]; extending our results to these methods where the  $\mathbf{Y}_k$  samples are not independent would require careful bounds on the conditional expectation of these random variables as described in Remark 2.

#### ACKNOWLEDGEMENTS

TA, MC, and JH were partially supported by NSF DMS #2211318 and JH was partially supported by NSF CAREER #2440040. ER and JL were partially supported by NSF DMS #2309685.

#### REFERENCES

- [1] A. Agaskar, C. Wang, and Y. M. Lu. Randomized Kaczmarz algorithms: Exact MSE analysis and optimal sampling probabilities. In *2014 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 389–393, 2014.
- [2] Z.-Z. Bai and W.-T. Wu. On convergence rate of the randomized Kaczmarz method. *Linear Algebra Appl.*, 553:252–269, 2018.
- [3] Z.-Z. Bai and W.-T. Wu. On greedy randomized Kaczmarz method for solving large sparse linear systems. *SIAM J. Sci. Comput.*, 40(1):A592–A606, 2018.
- [4] X. Chen and A. Powell. Almost sure convergence of the Kaczmarz algorithm with random measurements. *J. Fourier Anal. Appl.*, pages 1–20, 2012. 10.1007/s00041-012-9237-2.
- [5] J. A. De Loera, J. Haddock, and D. Needell. A sampling Kaczmarz-Motzkin algorithm for linear feasibility. *SIAM J. Sci. Comput.*, 39(5):S66–S87, 2017.
- [6] Michal Dereziński, Daniel LeJeune, Deanna Needell, and Elizaveta Rebrova. Fine-grained analysis and faster algorithms for iteratively solving linear systems. *Journal of Machine Learning Research*, 26(144):1–49, 2025.
- [7] R. Durrett. *Probability: Theory and Examples*. Cambridge University Press, 5 edition, 2019.
- [8] F. Fagnani and S. Zampieri. Randomized consensus algorithms over large scale networks. *IEEE J. Sel. Area Comm.*, 26(4):634–649, 2008.

- [9] Xiequan Fan, Ion Grama, and Quansheng Liu. Hoeffding’s inequality for supermartingales. *Stochastic Processes and their Applications*, 122(10):3545–3559, 2012.
- [10] V. Feldman and J. Vondrak. High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. In *Conference on Learning Theory*, pages 1270–1279. PMLR, 2019.
- [11] S. Ghadimi and G. Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization i: A generic algorithmic framework. *SIAM J. Optimiz.*, 22(4):1469–1492, 2012.
- [12] E. Gorbunov, A. Sadiev, M. Danilova, S. Horváth, G. Gidel, P. Dvurechensky, A. Gasnikov, and P. Richtárik. High-probability convergence for composite and distributed stochastic minimization and variational inequalities with heavy-tailed noise. *arXiv preprint arXiv:2310.01860*, 2023.
- [13] R. M. Gower and P. Richtárik. Randomized iterative methods for linear systems. *SIAM J. Matrix Anal. A.*, 36(4):1660–1690, 2015.
- [14] J. Haddock, B. Jarman, and C. Yap. Paving the way for consensus: Convergence of block gossip algorithms. *IEEE T. Inform. Theory*, 68(11):7515–7527, 2022.
- [15] J. Haddock and D. Needell. Randomized projection methods for linear systems with arbitrarily large sparse corruptions. *SIAM J. Sci. Comput.*, 41(5):S19–S36, 2019.
- [16] J. Haddock, D. Needell, E. Rebrova, and W. Swartworth. Quantile-based iterative methods for corrupted systems of linear equations. *SIAM J. Matrix Anal. A.*, 43(2):605–637, 2022.
- [17] P. C. Hansen. Regularization tools version 4.0 for Matlab 7.3. *Numer. Algorithms*, 46:189–194, 2007.
- [18] N. J. A. Harvey, C. Liaw, Y. Plan, and S. Randhawa. Tight analyses for non-smooth stochastic gradient descent. In *Conference on Learning Theory*, pages 1579–1613. PMLR, 2019.
- [19] A. Henriksen and R. Ward. Concentration inequalities for random matrix products. *Linear Algebra Appl.*, 594:81–94, 2020.
- [20] A. J. Hoffman. On approximate solutions of systems of linear inequalities. *J. Res. Nat. Bur. Stand.*, 49:263–265, 1952.
- [21] D. Huang, J. Niles-Weed, J. A. Tropp, and R. Ward. Matrix concentration for products. *Found. Comput. Math.*, 22(6):1767–1799, 2022.
- [22] Halyun Jeong, Deanna Needell, and Elizaveta Rebrova. Stochastic gradient descent for streaming linear and rectified linear systems with adversarial corruptions. *SIAM Journal on Mathematics of Data Science*, 7(2):516–541, 2025.
- [23] A. Juditsky, A. Nemirovski, and C. Tauvel. Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems*, 1(1):17–58, 2011.
- [24] S. Kaczmarz. Angenäherte auflösung von systemen linearer gleichungen. *Bull. Int. Acad. Polon. Sci. Lett. Ser. A*, pages 335–357, 1937.
- [25] T. Kathuria, S. Mukherjee, and N. Srivastava. On concentration inequalities for random matrix products. *arXiv preprint arXiv:2003.06319*, 2020.
- [26] J. Kleinberg and E. Tardos. *Algorithm design*. Pearson Education India, 2006.
- [27] D. Leventhal and A. S. Lewis. Randomized methods for linear constraints: convergence rates and conditioning. *Math. Oper. Res.*, 35(3):641–654, 2010.
- [28] Junhong Lin and Ding-Xuan Zhou. Learning theory of randomized kaczmarz algorithm. *The Journal of Machine Learning Research*, 16(1):3341–3365, 2015.
- [29] N. Loizou and P. Richtárik. Revisiting randomized gossip algorithms: General framework, convergence rates and novel block and accelerated protocols. *IEEE T. Inform. Theory*, 67(12):8300–8324, 2021.
- [30] Z. Lou, W. Zhu, and W. B. Wu. Beyond sub-gaussian noises: Sharp concentration analysis for stochastic gradient descent. *J. Mach. Learn. Res.*, 23(1):2227–2248, 2022.
- [31] A. Ma, D. Needell, and A. Ramdas. Convergence properties of the randomized extended Gauss–Seidel and Kaczmarz methods. *SIAM J. Matrix Anal. A.*, 36(4):1590–1604, 2015.
- [32] Nicholas F Marshall and Oscar Mickelin. An optimal scheduled learning rate for a randomized kaczmarz algorithm. *SIAM Journal on Matrix Analysis and Applications*, 44(1):312–330, 2023.
- [33] W. Mou, C. J. Li, M. J. Wainwright, P. L. Bartlett, and M. I. Jordan. On linear stochastic approximation: Fine-grained Polyak-Ruppert and non-asymptotic concentration. In *Conference on Learning Theory*, pages 2947–2997. PMLR, 2020.
- [34] E. Moulines and F. Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. *Adv. Neur. In.*, 24, 2011.
- [35] R. Murray, J. Demmel, M. W. Mahoney, N. B. Erichson, M. Melnichenko, O. A. Malik, L. Grigori, P. Luszczek, M. Derez-ński, M. E. Lopes, et al. Randomized numerical linear algebra: A perspective on the field with an eye to software. *arXiv preprint arXiv:2302.11474*, 2023.
- [36] D. Needell and J. A. Tropp. Paved with good intentions: Analysis of a randomized block Kaczmarz method. *Linear Algebra Appl.*, 2013.
- [37] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM J. Optimiz.*, 19(4):1574–1609, 2009.
- [38] Javier F Peña. An easily computable upper bound on the hoffman constant for homogeneous inequality systems. *Computational Optimization and Applications*, 87(1):323–335, 2024.
- [39] N. Pritchard and V. Patel. Solving, tracking and stopping streaming linear inverse problems. *Inverse Probl.*, 2024.
- [40] H. Robbins and S. Monro. A stochastic approximation method. *Ann. Math. Statist.*, 22:400–407, 1951.
- [41] S. Steinerberger. Quantile-based random Kaczmarz for corrupted linear systems of equations. *Inform. Inf.*, 12(1):448–465, 2023.
- [42] T. Strohmer and R. Vershynin. A randomized Kaczmarz algorithm with exponential convergence. *J. Fourier Anal. Appl.*, 15(2):262, 2009.

- [43] R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.
- [44] C. Wang, A. Agaskar, and Y. M. Lu. Randomized Kaczmarz algorithm for inconsistent linear systems: An exact MSE analysis. In *2015 International Conference on Sampling Theory and Applications (SampTA)*, pages 498–502, 2015.