

ITERATIVE METHODS FOR SOLVING FACTORIZED LINEAR SYSTEMS

A. MA, D. NEEDELL, A. RAMDAS

ABSTRACT. Stochastic iterative algorithms such as the Kaczmarz and Gauss-Seidel methods have gained recent attention because of their speed, simplicity, and the ability to approximately solve large-scale linear systems of equations without needing to access the entire matrix. In this work, we consider the setting where we wish to solve a linear system in a large matrix \mathbf{X} that is stored in a factorized form, $\mathbf{X} = \mathbf{UV}$; this setting either arises naturally in many applications or may be imposed when working with large low-rank datasets for reasons of space required for storage. We propose a variant of the randomized Kaczmarz method for such systems that takes advantage of the factored form, and avoids computing \mathbf{X} . We prove an exponential convergence rate and supplement our theoretical guarantees with experimental evidence demonstrating that the factored variant yields significant acceleration in convergence.

1. INTRODUCTION

Recently, revived interest in stochastic iterative methods like the Kaczmarz [11, 25, 22, 23] and Gauss-Seidel [6, 19] methods has grown due to the need for large-scale approaches for solving linear systems of equations. Such methods utilize simple projections and require access to only a single row in a given iteration, hence having a low memory footprint. For this reason, they are very efficient and practical for solving extremely large, usually highly overdetermined, linear systems. In this work, we consider algorithms for solving linear systems when the matrix is available in a factorized form. As we discuss below, such a factorization may arise naturally in the application, or may be constructed explicitly for efficient storage and computation. We seek a solution to the original system directly from its factorized form, without the need to perform matrix multiplication.

To that end, suppose we want to solve the linear system $\mathbf{X}\boldsymbol{\beta} = \mathbf{y}$ with $\mathbf{X} \in \mathbb{C}^{m \times n}$. However, instead of the full system \mathbf{X} , we only have access to \mathbf{U}, \mathbf{V} such that $\mathbf{X} = \mathbf{UV}$. In this case, we want to solve the linear system:

$$\mathbf{UV}\boldsymbol{\beta} = \mathbf{y}, \tag{1}$$

where $\mathbf{U} \in \mathbb{C}^{m \times k}$ and $\mathbf{V} \in \mathbb{C}^{k \times n}$. Instead of taking the product of \mathbf{U} and \mathbf{V} , to form \mathbf{X} , which may not be desirable, we approach this problem using stochastic iterative methods to solve the individual subsystems

$$\mathbf{U}\mathbf{x} = \mathbf{y} \tag{2}$$

$$\mathbf{V}\boldsymbol{\beta} = \mathbf{x}, \tag{3}$$

in an alternating fashion. Notice that if we substitute (3) into (2), we acquire the full linear system (1). We will often refer to (1) as the “full system” and (2) and (3) as “subsystems”, and say that a system is consistent if it has at least one solution (and inconsistent otherwise).

There are some situations when approximately knowing \mathbf{x} would suffice. We assume that (for reasons of interpretability, or for downstream usage) the scientist is genuinely interested in solving the full system, i.e. she is interested in the vector $\boldsymbol{\beta}$, not in \mathbf{x} .

It is arguably of practical interest to give special importance to the case of $k < \min(m, n)$, which arises in modern data science as motivated by the following examples, but we discuss other settings later.

1.1. **Motivation.** If \mathbf{X} is large and low-rank, one may have many reasons to work with a factorization of \mathbf{X} . We shall discuss three reasons below — algorithmic, infrastructural, and statistical.

Consider data matrices encountered in “recommender systems” in machine learning [2, 14, 20, 26, 27]. For concreteness, consider the Netflix (or Amazon, or Yelp) problem, where one has a users-by-movies matrix whose entries correspond to ratings given by users to movies. \mathbf{X} is usually quite well approximated by low-rank matrices — intuitively, many rows and columns are redundant because every row is usually similar to many other rows (corresponding to users with similar tastes), and every column is usually similar to many other columns (corresponding to similar quality movies in the same genre). Usually we have observed only a few entries of \mathbf{X} , and wish to infer the unseen ratings in order to provide recommendations to different users based on their tastes. Algorithms for “low-rank matrix completion” have proved to be quite successful in the applied and theoretical machine learning community [5, 13, 28, 12]. One popular algorithm, alternating-minimization [10], chooses a (small) target rank k , and tries to find \mathbf{U}, \mathbf{V} such that $\mathbf{X}_{ij} \approx (\mathbf{UV})_{ij}$ for all the observed entries (i, j) of \mathbf{X} . As its name suggests, the algorithm alternates between solving for \mathbf{U} keeping \mathbf{V} fixed and then solving for \mathbf{V} keeping \mathbf{U} fixed. In this case, at no point does the algorithm even form the entire completed (inferred) matrix \mathbf{X} , and the algorithm only has access to factors \mathbf{U}, \mathbf{V} simply due to *algorithmic* choices.

There may be other instances where a data scientist may have access to the full matrix \mathbf{X} , but in order to reduce the memory storage footprint, or to communicate the data, may explicitly choose to decompose $\mathbf{X} \approx \mathbf{UV}$ and discard \mathbf{X} to work with the smaller matrices instead.

Consider an example motivated by “topic modeling” of text data. Suppose Google has scraped the internet for English documents (or maybe a subset of documents like news articles), to form a document-by-word matrix \mathbf{X} , where each entry of the matrix indicates the number of times that word occurred in that document. Since many documents could be quite similar in their content (like articles about the same incident covered by different newspapers), this matrix is easily seen to be low-rank. This is a classic setting for applying a machine learning technique called “non-negative matrix factorization” [15, 29, 18], where one decomposes \mathbf{X} as the product of two low-rank non-negative matrices \mathbf{U}, \mathbf{V} ; the non-negativity is imposed for human interpretability, so that \mathbf{U} can be interpreted as a documents-by-topics matrix, and \mathbf{V} as a topics-by-words. In this case, we do not have access to \mathbf{X} as a result of *systems infrastructure* constraints (memory/storage/communication).

Often, even for modestly sized data matrices, the relevant “signal” is contained in the leading singular vectors corresponding to large singular values, and the tail of small eigenvalues is often deemed to be “noise”. This is precisely the idea behind the classical topic of principal component analysis (PCA), and the modern machine learning literature has proposed and analyzed a variety of algorithms to approximate the top k left and right singular vectors in a streaming/stochastic/online fashion [7]. Hence, the factorization may arise from a purely *statistical* motivation.

Given a vector \mathbf{y} (representing age, or document popularity, for example), suppose the data scientist is interested in regressing \mathbf{X} onto \mathbf{y} , for the purpose of scientific understanding or to take future actions. Can we utilize the available factorization efficiently, designing methods that work directly on the lower dimensional factors \mathbf{U} and \mathbf{V} rather than computing the full system \mathbf{X} ?

Our goal will be to propose iterative methods that work directly on the factored system, eliminating the need for a full matrix product and potentially saving computations on the much larger full system.

1.2. **Main contribution.** We propose two stochastic iterative methods for solving system (1) without computing the product of \mathbf{U} and \mathbf{V} . Both methods utilize iterates of well studied algorithms for solving linear systems. When the full system is consistent, the first method, called RK-RK, interlaces iterates of the Randomized Kaczmarz (RK) algorithm to solve each subsystem and finds the optimal solution. When the full system is inconsistent, we introduce the REK-RK

method, an interlacing of Randomized Extended Kaczmarz (REK) iterates to solve (2) and RK iterates to solve (3), that converges to the so-called ordinary least squares solution. We prove linear (“exponential”) convergence to the solution in both cases.

1.3. Outline. In the next section, we provide background and discuss existing work on stochastic methods that solve linear systems. In particular, we describe the RK and REK algorithms as well as the Randomized Gauss-Seidel (RGS) and Randomized Extended Gauss-Seidel (REGS) algorithms. In Section 3 we investigate variations of settings for subsystems (2) and (3) that arise depending on the consistency and size of \mathbf{X} . Section 4 introduces our proposed methods, RK-RK and REK-RK. We provide theory that shows linear convergence in expectation to the optimal solution for both methods. Finally, we present experiments in Section 5 and conclude with final remarks and future work in Section 6.

1.4. Notation. Here and throughout the paper, matrices and vectors are denoted with boldface letters (uppercase for matrices and lowercase for vectors). We call \mathbf{X}^i the i^{th} row of the matrix \mathbf{X} and $\mathbf{X}_{(j)}$ the j^{th} column of \mathbf{X} . The euclidean norm is denoted by $\|\cdot\|_2$ and the Frobenius norm by $\|\cdot\|_F$. Lastly, \mathbf{X}^* denotes the adjoint (conjugate transpose) of the matrix \mathbf{X} . Motivated by applications, we allow \mathbf{X} to be rank deficient and assume that \mathbf{U} and \mathbf{V} are full rank.

2. BACKGROUND AND EXISTING WORK

In this section we summarize existing work on stochastic iterative methods and different variations of linear systems.

2.1. Linear Systems. Linear systems take on one of three settings determined by the size of the system, rank of the matrix \mathbf{X} , and the existence of a solution. First we discuss solutions to systems with full rank matrices \mathbf{X} then remark on how rank deficiency affects the desired solution.

In the full rank underdetermined case, $m < n$ and the system has infinitely many solutions; here, we often want to find the least Euclidean norm solution to (1):

$$\beta_{LN} := \mathbf{X}^*(\mathbf{X}\mathbf{X}^*)^{-1}\mathbf{y}. \quad (4)$$

Clearly, $\mathbf{X}\beta_{LN} = \mathbf{y}$, and all other solutions to an underdetermined system can be written as $\mathbf{b} = \beta_{LN} + \mathbf{z}$ where $\mathbf{X}\mathbf{z} = 0$.

A system with a unique solution is called an overdetermined *consistent* system while a system with no exact solution is called an overdetermined *inconsistent* system. In the overdetermined consistent setting, when \mathbf{X} is full rank, the optimal solution is the unique β_{uniq} such that $\mathbf{X}\beta_{\text{uniq}} = \mathbf{y}$:

$$\beta_{\text{uniq}} := (\mathbf{X}^*\mathbf{X})^{-1}\mathbf{X}^*\mathbf{y}. \quad (5)$$

When a system is inconsistent, we often seek to minimize the sum of squared residuals, i.e. to find the ordinary least squares solution

$$\beta_{LS} := (\mathbf{X}^*\mathbf{X})^{-1}\mathbf{X}^*\mathbf{y}. \quad (6)$$

The residual can be written as $\mathbf{r} = \mathbf{X}\beta_{LS} - \mathbf{y}$. Note that $\mathbf{X}^*\mathbf{r} = 0$, which can be easily seen by substituting $\mathbf{y} = \mathbf{X}\beta_{LS} + \mathbf{r}$ into (6). For simplicity, we will refer to the matrix \mathbf{X} of a linear system as consistent or inconsistent when the system itself is consistent or inconsistent.

If the matrix \mathbf{X} in the linear system $\mathbf{X}\beta = \mathbf{y}$ is rank deficient, then there are infinitely many solutions to the system regardless the size of m and n . In this case, we again want the least norm solution in the underdetermined case and the “least-norm least-squares” solution in the overdetermined case,

$$\beta_{LN} := \begin{cases} \mathbf{X}^*(\mathbf{X}\mathbf{X}^*)^\dagger\mathbf{y}, & \text{if } m < n \\ (\mathbf{X}^*\mathbf{X})^\dagger\mathbf{X}^*\mathbf{y}, & \text{if } m \geq n, \end{cases} \quad (7)$$

where $(\cdot)^\dagger$ is the pseudo-inverse. General solutions to the linear system can be written as $\mathbf{b} = \boldsymbol{\beta}_{LN} + \mathbf{z}$ where $\mathbf{X}\mathbf{z} = 0$ —note that $\mathbf{y} = \mathbf{X}\mathbf{b} = \mathbf{X}\boldsymbol{\beta}_{LN} + \mathbf{X}\mathbf{z} = \mathbf{X}\boldsymbol{\beta}_{LN}$. Similar to the full rank case, when the low-rank system is inconsistent, we can write $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_{LN} + \mathbf{r}$, again where $\mathbf{X}^*\mathbf{r} = 0$.

2.2. Randomized Kaczmarz and its Extension. The Kaczmarz Algorithm [11] solves a linear system $\mathbf{X}\boldsymbol{\beta} = \mathbf{y}$ by cycling through rows of \mathbf{X} and projects the estimate onto the solution space given by the chosen row. It was initially proposed by Kaczmarz [11] and has recently regained interest in the setting of computer tomography where it is known as the Algebraic Reconstruction Technique [8, 21, 4, 9]. The randomized variant of the Kaczmarz method introduced by Strohmer and Vershynin [25] was proven to converge linearly in expectation for consistent systems. Formally, given \mathbf{X} and \mathbf{y} of (1), RK chooses row $i \in \{1, 2, \dots, m\}$ of \mathbf{X} with probability $\frac{\|\mathbf{X}^i\|_2^2}{\|\mathbf{X}\|_F^2}$, and projects the previous estimate onto that row with the update

$$\boldsymbol{\beta}_t := \boldsymbol{\beta}_{t-1} + \frac{(\mathbf{y}_i - \mathbf{X}^i\boldsymbol{\beta}_{t-1})}{\|\mathbf{X}^i\|_2^2}(\mathbf{X}^i)^*.$$

Needell [22] later studied the inconsistent case and showed that RK does not converge to the least squares solution for inconsistent systems, but rather converges linearly to some convergence radius of the solution. To remedy this, Zouzias and Freris [30] proposed the Randomized Extended Kaczmarz (REK) algorithm to solve linear systems in all settings. For REK, row $i \in \{1, 2, \dots, m\}$ and column $j \in \{1, \dots, n\}$ of \mathbf{X} are chosen at random with probability

$$P(\text{row} = i) = \frac{\|\mathbf{X}^i\|_2^2}{\|\mathbf{X}\|_F^2}, \quad P(\text{column} = j) = \frac{\|\mathbf{X}_{(j)}\|_2^2}{\|\mathbf{X}\|_F^2},$$

and starting from $\boldsymbol{\beta}_0 = 0$ and $\mathbf{z}_0 = \mathbf{y}$, every iteration computes

$$\boldsymbol{\beta}_t := \boldsymbol{\beta}_{t-1} + \frac{(\mathbf{y}^i - \mathbf{z}_t^i - \mathbf{X}^i\boldsymbol{\beta}_{t-1})}{\|\mathbf{X}^i\|_2^2}(\mathbf{X}^i)^*, \quad \mathbf{z}_t := \mathbf{z}_{t-1} - \frac{\langle \mathbf{X}_{(j)}, \mathbf{z}_{t-1} \rangle}{\|\mathbf{X}_{(j)}\|_2^2} \mathbf{X}_{(j)}.$$

REK finds the optimal solution in all linear system settings. In the consistent setting, it behaves as RK. In the *overdetermined* inconsistent setting, \mathbf{z} estimates the residual vector \mathbf{r} and allows $\boldsymbol{\beta}_t$ to converge to the true least squares solution of the system. REK was shown to converge linearly in expectation to the least-squares solution by Zouzias and Freris [30].

2.3. Randomized Gauss-Seidel and its Extension. The Gauss-Seidel method was originally published by Seidel but it was later discovered that Gauss had studied this method in a letter to his student [3]. Instead of relying on rows of a matrix, the Gauss-Seidel method relies on columns of \mathbf{X} . The randomized variant was studied by Leventhal and Lewis [16] shortly after RK was published.

The randomized variant (RGS) requires a column j to be chosen randomly with probability $\frac{\|\mathbf{X}_{(j)}\|_2^2}{\|\mathbf{X}\|_F^2}$, and updates at every iteration

$$\boldsymbol{\beta}_t := \boldsymbol{\beta}_{t-1} + \frac{\mathbf{X}_{(j)}^*(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_{t-1})}{\|\mathbf{X}_{(j)}\|_2^2} \mathbf{e}_{(j)}, \quad (8)$$

where $\mathbf{e}_{(j)}$ is the j^{th} basis vector (a vector with 1 in the j^{th} position and 0 elsewhere). Leventhal and Lewis [16] showed that RGS converges linearly in expectation when \mathbf{X} is overdetermined. However, it fails to find the least norm solution for an *underdetermined* linear system [19]. The Randomized Extended Gauss-Seidel (REGS) resolves this problem, much like REK did for RK in the case of

overdetermined systems. The method chooses a random row and column of \mathbf{X} exactly as in REK, and then updates at every iteration

$$\begin{aligned}\beta_t &:= \beta_{t-1} + \frac{\mathbf{X}_{(j)}^*(\mathbf{y} - \mathbf{X}\beta_{t-1})}{\|\mathbf{X}_{(j)}\|_2^2} \mathbf{e}_{(j)}, \\ \mathbf{P}_i &:= \mathbf{Id}_n - \frac{(\mathbf{X}^i)^* \mathbf{X}^i}{\|\mathbf{X}^i\|_2^2}, \\ \mathbf{z}_t &:= \mathbf{P}_i(\mathbf{z}_{t-1} + \beta_t - \beta_{t-1}),\end{aligned}$$

and at any fixed time t , outputs $\beta_t - \mathbf{z}_t$ as the estimated solution to $\mathbf{X}\beta = \mathbf{y}$. Here, \mathbf{Id}_n denotes the $n \times n$ identity matrix. This extension works for all variations of linear systems and was proven to converge linearly in expectation by Ma et al. [19].

The RK and RGS methods along with their extensions are extensively studied and compared in [19]. Table 1 summarizes the convergence properties of each of the randomized methods and their extensions.

Method	Overdetermined, consistent : convergence to β_{uniq} ?	Overdetermined, inconsistent : convergence to β_{LS} ?	Underdetermined : convergence to β_{LN} ?
RK	Yes [25]	No [22]	Yes [19]
REK	Yes [30]	Yes [30]	Yes [19]
RGS	Yes [16]	Yes [16]	No [19]
REGS	Yes [19]	Yes [19]	Yes [19]

TABLE 1. Summary of convergence properties of randomized methods under all settings.

In this paper, we focus on using combinations of RK and REK but also discuss RGS and REGS for comparison. We choose to focus on RK and REK because their updates consist only of scalar operations and inner products as opposed to REGS which requires an outer product. The methods proposed are easily extendable to RGS and REGS.

3. VARIATIONS OF FACTORED LINEAR SYSTEMS

Our proposed methods rely on interleaving solution estimates to subsystem (2) and subsystem (3). Because the convergence of RK, RGS, REK, and REGS are heavily dependent on the number of rows and columns in the linear system, it is important to discuss how the settings of (2) and (3) are determined by \mathbf{X} . In this section, we will discuss when we can expect our methods to solve the full system.

Linear System	Optimal Solution
$\mathbf{X}\beta = \mathbf{y}$ (1)	β_\star
$\mathbf{U}\mathbf{x} = \mathbf{y}$ (2)	\mathbf{x}_\star
$\mathbf{V}\mathbf{b} = \mathbf{x}$ (3)	\mathbf{b}_\star

TABLE 2. Summary of notation for linear systems discussed and their solutions

For simplicity in notation, we will denote β_\star , \mathbf{x}_\star , and \mathbf{b}_\star as the “optimal” solution of (1), (2), and (3) respectively, as summarized in Table 2. By “optimal” solution for (2) and (3), we mean the unique, least norm, or the least squares solution, depending on the type of system (overdetermined consistent, underdetermined, overdetermined inconsistent). Since we assume that \mathbf{UV} may be

low-rank, β_\star is going to be the least norm solution as described in (7). Table 3 presents such a summary depending on the size of k with respect to m and n .

\mathbf{X}	$k < m, n$	$m < k < n$ $n < k < m$	$k > m, n$
Underdetermined	$\mathbf{U} = \text{Over, Consis.}$ $\mathbf{V} = \text{Under}$ (S1)	$\mathbf{U} = \text{Under}$ $\mathbf{V} = \text{Under}$ (S2)	$\mathbf{U} = \text{Under}$ $\mathbf{V} = \text{Over, Consis.}$ (S2)
Overdetermined Consis.	$\mathbf{U} = \text{Over, Consis.}$ $\mathbf{V} = \text{Under}$ (S1)	$\mathbf{U} = \text{Over, Consis.}$ $\mathbf{V} = \text{Over, Consis.}$ (S1)	$\mathbf{U} = \text{Under}$ $\mathbf{V} = \text{Over, Consis.}$ (S2)
Overdetermined Inonsis.	$\mathbf{U} = \text{Over, Incon.}$ $\mathbf{V} = \text{Under}$ (S3b)	$\mathbf{U} = \text{Over, Incon.}$ $\mathbf{V} = \text{Over, Consis.}$ (S3a)	$\mathbf{U} = \text{Under}$ $\mathbf{V} = \text{Over, Consis.}$ (S2)

TABLE 3. Summary of types of matrices \mathbf{U} and \mathbf{V} for given m, n , and k relations. The cells in white indicate when our proposed methods will converge to the solution of the full system, gray cells indicate when our methods will not. Recall that $\mathbf{X} \in \mathbb{C}^{m \times n}$, $\mathbf{U} \in \mathbb{C}^{m \times k}$ and $\mathbf{V} \in \mathbb{C}^{k \times n}$. Arguably, the $k < m, n$ setting is most practically relevant, and in this case our methods do recover the optimal solution.

We spend the rest of this section justifying the observations in Table 3. For this, we split Table 3 into three scenarios : (S1) \mathbf{U} is overdetermined and consistent, (S2) \mathbf{U} is underdetermined, (S3a and S3b) \mathbf{X} is overdetermined and inconsistent. This section provides the intuition on when we should expect our methods (or similar ones based on interleaving solutions to the subsystems) to work. However, one may also skip ahead to the next section where formally we present our algorithm and main results.

- **Scenario S1: \mathbf{U} overdetermined, consistent.** When \mathbf{U} is overdetermined and consistent, we find that solving (2) and (3) gives us the optimal solution of (1).

Indeed, in the case where \mathbf{V} is overdetermined and consistent, we have:

$$\begin{aligned}
 \mathbf{b}_\star &= (\mathbf{V}^* \mathbf{V})^{-1} \mathbf{V}^* (\mathbf{U}^* \mathbf{U})^{-1} \mathbf{U}^* \mathbf{y} \\
 &= (\mathbf{V}^* \mathbf{V})^{-1} \mathbf{V}^* (\mathbf{U}^* \mathbf{U})^{-1} \mathbf{U}^* \mathbf{U} \mathbf{V} \beta_\star \\
 &= (\mathbf{V}^* \mathbf{V})^{-1} \mathbf{V}^* \mathbf{V} \beta_\star \\
 &= \beta_\star
 \end{aligned}$$

In the case where \mathbf{V} is underdetermined, we have:

$$\begin{aligned}
 \mathbf{V} \mathbf{b}_\star &= \mathbf{V} \mathbf{V}^* (\mathbf{V} \mathbf{V}^*)^{-1} (\mathbf{U}^* \mathbf{U})^{-1} \mathbf{U}^* \mathbf{y} \\
 &= \mathbf{V} \mathbf{V}^* (\mathbf{V} \mathbf{V}^*)^{-1} (\mathbf{U}^* \mathbf{U})^{-1} \mathbf{U}^* \mathbf{U} \mathbf{V} \beta_\star \\
 &= \mathbf{V} \beta_\star
 \end{aligned}$$

Since \mathbf{X} is possibly low-rank, we still need to argue that this implies that $\mathbf{b}_\star = \beta_\star$, i.e. \mathbf{b}_\star is indeed the least norm solution to the full system. Suppose towards a contradiction that β_\star is the least norm solution of the full system but not subsystem (3); in other words assume that $\beta_\star = \mathbf{b}_\star + \mathbf{b}$ where $\mathbf{V} \mathbf{b} = 0$ and \mathbf{b} is nontrivial. Multiplying both sides by \mathbf{X} , we see that since β_\star has a nontrivial component \mathbf{b} such that $\mathbf{X} \mathbf{b} = 0$, it cannot be the least norm solution to the full system as assumed, reaching a contradiction. *Therefore, in the consistent case when \mathbf{U} is overdetermined, we have that $\mathbf{b}_\star = \beta_\star$, and may hope that our proposed methods will be able to solve the full system (1) utilizing the subsystems (2) and (3).*

- **Scenario S2: \mathbf{U} underdetermined.** When \mathbf{U} is underdetermined, solving (2) and (3) for their optimal solutions does not guarantee the optimal solution of the full system.

Intuitively, (2) has infinitely many solutions and $\mathbf{x}_\star = \mathbf{x}_{LN} \neq \mathbf{V}\beta_\star$. Mathematically, investigating \mathbf{b}_\star , we find that

$$\begin{aligned}
\mathbf{b}_\star &= (\mathbf{V}^*\mathbf{V})^{-1}\mathbf{V}^*\mathbf{U}^*(\mathbf{U}\mathbf{U}^*)^{-1}\mathbf{y} \\
&= (\mathbf{V}^*\mathbf{V})^{-1}\mathbf{V}^*\mathbf{U}^*(\mathbf{U}\mathbf{U}^*)^{-1}\mathbf{U}\mathbf{V}\beta_\star \\
&\neq \beta_\star \text{ if } \mathbf{V} \text{ overdetermined, consistent} \\
\mathbf{V}\mathbf{b}_\star &= \mathbf{V}\mathbf{V}^*(\mathbf{V}\mathbf{V}^*)^{-1}\mathbf{U}^*(\mathbf{U}\mathbf{U}^*)^{-1}\mathbf{y} \neq \beta_\star \\
&= \mathbf{V}\mathbf{V}^*(\mathbf{V}\mathbf{V}^*)^{-1}\mathbf{U}^*(\mathbf{U}\mathbf{U}^*)^{-1}\mathbf{U}\mathbf{V}\beta_\star \\
&= \mathbf{U}^*(\mathbf{U}\mathbf{U}^*)^{-1}\mathbf{U}\mathbf{V}\beta_\star \\
&\neq \mathbf{V}\beta_\star \text{ if } \mathbf{V} \text{ underdetermined.}
\end{aligned}$$

Therefore, we do not expect our proposed methods to succeed when \mathbf{U} is underdetermined. Fortunately, this case seems to be of little practical interest, since factoring an underdetermined system does not typically save any computation.

- **Scenario S3: \mathbf{X} inconsistent.** Before we discuss whether it's possible to recover the optimal solution to the full system, we must first discuss what \mathbf{X} being inconsistent implies about the subsystems (2) and (3). In particular, one needs to determine whether inconsistency in the full system creates inconsistencies in the individual subsystems. If \mathbf{X} is inconsistent then we have $\mathbf{X}\beta_\star + \mathbf{r} = \mathbf{y}$ where β_\star is the optimal solution of (1) and $\mathbf{X}^*\mathbf{r} = 0$. Now, consider decomposing $\mathbf{r} = \mathbf{r}_1 + \mathbf{r}_2$ where $\mathbf{U}^*\mathbf{r}_1 = 0$, $\mathbf{U}^*\mathbf{r}_2 \neq 0$, and $\mathbf{V}^*\mathbf{U}^*\mathbf{r}_2 = 0$. Notice that $\mathbf{X}^*\mathbf{r} = \mathbf{V}^*\mathbf{U}^*(\mathbf{r}_1 + \mathbf{r}_2) = \mathbf{V}^*\mathbf{U}^*\mathbf{r}_1 + \mathbf{V}^*\mathbf{U}^*\mathbf{r}_2 = 0$, as desired. We want to decompose the full system $\mathbf{X}\beta_\star + \mathbf{r} = \mathbf{y}$ into two subsystems. Following a similar thought process as before, we choose to decompose our full system into the following:

$$\mathbf{U}\mathbf{x} + \mathbf{r}_1 + \mathbf{r}_2 = \mathbf{y} \quad (9)$$

$$\mathbf{V}\mathbf{b} = \mathbf{x}. \quad (10)$$

Clearly, (9) is inconsistent since $\mathbf{U}^*\mathbf{r}_1 = 0$ and $\mathbf{U}^*\mathbf{r}_2 \neq 0$. To determine if (10) is inconsistent, we must see whether \mathbf{x}_\star has a nontrivial component which lies in the null space of \mathbf{V}^* . Note that $\mathbf{x}_\star = (\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*(\mathbf{y} - \mathbf{r}_1 - \mathbf{r}_2) = (\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*(\mathbf{y} - \mathbf{r}_2)$ is the least squares solution since \mathbf{U} must be overdetermined for \mathbf{X} to be inconsistent.

- **Case S3a: \mathbf{V} overdetermined.** Note that the second subsystem (10) is consistent (since there is no component in the null space of \mathbf{V}^*). In this case, we have

$$\begin{aligned}
\mathbf{b}_\star &= (\mathbf{V}^*\mathbf{V})^{-1}\mathbf{V}^*(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*(\mathbf{y} - \mathbf{r}_1 - \mathbf{r}_2) \\
&= (\mathbf{V}^*\mathbf{V})^{-1}\mathbf{V}^*(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*(\mathbf{y} - \mathbf{r}_2) \\
&= \beta_\star - (\mathbf{V}^*\mathbf{V})^{-1}\mathbf{V}^*(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*\mathbf{r}_2 \\
&\neq \beta_\star.
\end{aligned}$$

Therefore, in this case, we do not expect to find the optimal solution to (1).

- **Case S3b: \mathbf{V} underdetermined.** In this case, $\mathbf{r}_2 = 0$ and solving (9) and (10) obtains the optimal solution to the full system since

$$\begin{aligned}
\mathbf{V}\mathbf{b}_\star &= \mathbf{V}\mathbf{V}^*(\mathbf{V}\mathbf{V}^*)^{-1}(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*(\mathbf{y} - \mathbf{r}_1) \\
&= \mathbf{V}\mathbf{V}^*(\mathbf{V}\mathbf{V}^*)^{-1}(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*\mathbf{y} \\
&= \mathbf{V}\mathbf{V}^*(\mathbf{V}\mathbf{V}^*)^{-1}(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*\mathbf{U}\mathbf{V}\beta_\star \\
&= \mathbf{V}\mathbf{V}^*(\mathbf{V}\mathbf{V}^*)^{-1}\mathbf{V}\beta_\star \\
&= \mathbf{V}\beta_\star
\end{aligned}$$

Following the same argument as in Scenario S1 when \mathbf{V} is underdetermined, we reach the conclusion that $\mathbf{b}_\star = \beta_\star$. Thus, in this case our methods have the potential to solve the full system.

These three scenarios fully explain the observations made in Table 3. The focus of the remainder of this paper will be the case in which $k < m, n$ (i.e. left column of Table 3) since, as mentioned, it is more practically useful.

4. METHODS AND MAIN RESULTS

Our approach intertwines two iterative methods to solve subsystem (2) followed by subsystem (3). For the consistent setting, we propose Algorithm 1 which uses an iterate of RK on (2) intertwined with an iterate of RK to solve (3). For the inconsistent setting, we propose using REK to solve subsystem (9) followed by RK to solve subsystem (10) as shown in Algorithm 2. Recall that \mathbf{x}_t^p is the p^{th} element in the vector \mathbf{x}_t . We propose an approach that interlaces solving subsystems (2) and (3); this has a couple of advantages over solving each subsystem separately in a sequential manner. First, if we are given some tolerance ϵ that we allow on the full system, it is unclear when we should stop the iterates of the first subsystem to obtain such an error — if solving the first subsystem is terminated prematurely, the error may propagate through iterates when solving the second subsystem. Second, the interlacing allows for opportunities to implement these algorithms in parallel. We leave the specifics of such an implementation as future work as it is outside the scope of this paper.

Algorithm 1: RK-RK

Input: $\mathbf{U}, \mathbf{V}, \mathbf{y}$
 Choose row \mathbf{U}^i with probability $\frac{\|\mathbf{U}^i\|_2^2}{\|\mathbf{U}\|_F^2}$
 Update $\mathbf{x}_t := \mathbf{x}_{t-1} + \frac{(\mathbf{y}^i - \mathbf{U}^i \mathbf{x}_{t-1})}{\|\mathbf{U}^i\|_2^2} (\mathbf{U}^i)^*$
 Choose row \mathbf{V}^p with probability $\frac{\|\mathbf{V}^p\|_2^2}{\|\mathbf{V}\|_F^2}$
 Update $\mathbf{b}_t := \mathbf{b}_{t-1} + \frac{(\mathbf{x}_t^p - \mathbf{V}^p \mathbf{b}_{t-1})}{\|\mathbf{V}^p\|_2^2} (\mathbf{V}^p)^*$

Algorithm 2: REK-RK

Input: $\mathbf{U}, \mathbf{V}, \mathbf{y}$
 Choose row \mathbf{U}^i with probability $\frac{\|\mathbf{U}^i\|_2^2}{\|\mathbf{U}\|_F^2}$
 Choose column $\mathbf{U}_{(j)}$ with probability $\frac{\|\mathbf{U}_{(j)}\|_2^2}{\|\mathbf{U}\|_F^2}$
 Update $\mathbf{z}_t := \mathbf{z}_{t-1} - \frac{\mathbf{U}_{(j)}^* \mathbf{z}_{t-1}}{\|\mathbf{U}_{(j)}\|_2^2} \mathbf{U}_{(j)}$
 Update $\mathbf{x}_t := \mathbf{x}_{t-1} + \frac{(\mathbf{y}^i - \mathbf{z}_t^i + \mathbf{U}^i \mathbf{x}_{t-1})}{\|\mathbf{U}^i\|_2^2} (\mathbf{U}^i)^*$
 Choose row \mathbf{V}^p with probability $\frac{\|\mathbf{V}^p\|_2^2}{\|\mathbf{V}\|_F^2}$
 Update $\mathbf{b}_t := \mathbf{b}_{t-1} + \frac{(\mathbf{x}_t^p - \mathbf{V}^p \mathbf{b}_{t-1})}{\|\mathbf{V}^p\|_2^2} (\mathbf{V}^p)^*$

4.1. Main result. Our main result shows that Algorithm 1 and Algorithm 2 converge linearly to the desired solution. The convergence rate, as expected, is a function of the conditioning of the subsystems, and hence we introduce the following notation. Here and throughout, for any matrix \mathbf{A} we write

$$\alpha_A := 1 - \frac{\sigma_{\min}^2(\mathbf{A})}{\|\mathbf{A}\|_F^2}, \quad (11)$$

$$\kappa_A^2 := \frac{\sigma_{\min}^2(\mathbf{A})}{\sigma_{\max}^2(\mathbf{A})}, \quad (12)$$

$$\theta_A := \frac{1}{\sigma_{\min}^2(\mathbf{A})}, \quad (13)$$

where $\sigma_{\min}^2(\mathbf{A})$ is the smallest non-zero singular value of \mathbf{A} , and κ_A^2 is the condition number of \mathbf{A} . Unless otherwise mentioned, in this section we assume $k < m, n$ since, as mentioned earlier, this is arguably the most practical setting. Recall that the *optimal solution* to a system is either the least-norm, unique, or least-squares solution depending on whether the system is underdetermined, overdetermined consistent, or overdetermined inconsistent, respectively.

Theorem 1. *If systems (1), (2), and (3) have optimal solutions $\beta_\star, \mathbf{x}_\star$ and \mathbf{b}_\star respectively, and (a) if system (1) is consistent, then $\mathbf{b}_\star = \beta_\star$ and Algorithm 1 converges with expected error*

$$\mathbb{E}\|\mathbf{b}_t - \beta_\star\|^2 \leq \alpha_V^t \|\mathbf{b}_\star\|^2 + \theta_V \alpha_U^t \|\mathbf{x}_\star\|^2,$$

(b) if system (1) is inconsistent, then $\mathbf{b}_\star = \beta_\star$ and Algorithm 2 converges with expected error

$$\mathbb{E}\|\mathbf{b}_t - \beta_\star\|^2 \leq \alpha_V^t \|\mathbf{b}_\star\|^2 + \theta_V \alpha_U^{\lfloor t/2 \rfloor} (1 + 2\kappa_U^2) \|\mathbf{x}_\star\|^2.$$

where α_U, α_V are as defined in (11), θ_V in (13), and κ_U^2 in (12).

Remarks.

1. Theorem 1(a) also applies to the setting in which \mathbf{X} is overdetermined, consistent and $n < k < m$.

2. The individual subsystems (2) and (3) are better conditioned than the full subsystem (1). It can easily be verified that $\alpha_V, \alpha_U \leq \alpha_X$. Empirically, our experiments in the next section suggest that \mathbf{U} and \mathbf{V} can be substantially better conditioned than \mathbf{X} .

4.2. Supporting results. To prepare for the proof of the above theorem (the central theoretical result of the paper), we state a few supporting results which will help simplify the presentation of the proof. We begin by stating known results on the convergence of RK and REK on linear systems. Let \mathbb{E}_b denote the expected value taken over the choice of rows in \mathbf{V} and \mathbb{E}_x the expected value taken over the choice of rows in \mathbf{U} and when necessary the choice of columns in \mathbf{U} . Also, let \mathbb{E} denote the full expected value (over all random variables and iterations) and \mathbb{E}^{t-1} be the expectation conditional on the first $t - 1$ iterations.

Proposition 1. ([24, Theorem 2]) *Given a consistent linear system $\mathbf{X}\beta = \mathbf{y}$, the Randomized Kaczmarz algorithm as described in Section 2.2 converges to the optimal solution β_\star with expected error*

$$\mathbb{E}\|\beta_t - \beta_\star\|^2 \leq \alpha_X^t \|\beta_\star\|^2,$$

where α_X is as defined in (11).

Proposition 2. ([30, Theorem 8]) *Given a linear system $\mathbf{X}\beta = \mathbf{y}$, the Randomized Extended Kaczmarz algorithm as described in Section 2.2 converges to the optimal solution β_\star with expected error*

$$\mathbb{E}\|\beta_t - \beta_\star\|^2 \leq \alpha_X^{\lfloor t/2 \rfloor} (1 + 2\kappa_X^2) \|\beta_\star\|^2,$$

where α_X is as defined in (11) and κ_X^2 is as defined in (12).

The proof of Theorem 1 builds directly on two useful lemmas. Lemma 1 addresses the impact of intertwining the algorithms. In particular, it shows useful relationships involving $\tilde{\mathbf{b}}_t$, the RK update solving the linear system $\mathbf{V}\mathbf{b} = \mathbf{x}_*$ at the t^{th} iteration (with \mathbf{b}_{t-1} as the previous estimate), and our update \mathbf{b}_t . Lemma 2 states that conditional on the first $t-1$ iterations, we can split the norm squared error $\|\mathbf{b}_t - \mathbf{b}_*\|^2$ into two terms relating to the error from solving subsystem (2) and the error from solving subsystem (3). To complete the proof of Theorem 1, we bound the error from solving (2) depending on whether we use RK (as in Algorithm 1) or REK (as in Algorithm 2) then apply the law of iterated expectations to bound the error from solving (3). We now state the aforementioned lemmas, and then formally prove the theorem.

Lemma 1. *Let $\tilde{\mathbf{b}}_t = \mathbf{b}_{t-1} + \frac{\mathbf{x}_*^p - \mathbf{V}^p \mathbf{b}_{t-1}}{\|\mathbf{V}^p\|^2} (\mathbf{V}^p)^*$. In Algorithm 1 and Algorithm 2 we have that:*

- (a) $\mathbb{E}_b^{t-1} \langle \mathbf{b}_t - \tilde{\mathbf{b}}_t, \tilde{\mathbf{b}}_t - \mathbf{b}_* \rangle = 0$,
- (b) $\|\tilde{\mathbf{b}}_t - \mathbf{b}_*\|^2 = \|\mathbf{b}_{t-1} - \mathbf{b}_*\|^2 - \|\tilde{\mathbf{b}}_t - \mathbf{b}_{t-1}\|^2$.

In words, part (a) states that the difference between an RK iterate solving the exact linear system $\mathbf{V}\mathbf{b} = \mathbf{x}_*$ and our RK iterate (which solves the linear system resulting from intertwining $\mathbf{V}\mathbf{b} = \mathbf{x}_t$), is orthogonal to $\tilde{\mathbf{b}}_t - \mathbf{b}_*$. This will come in handy in Lemma 2. Part (b) is a Pythagoras-style statement, which follows from well-known orthogonality properties of RK updates, included here for simplicity and completeness.

Proof. We prove statement (a) by direct substitution and expansion, as follows:

$$\begin{aligned}
\mathbb{E}_b^{t-1} \langle \mathbf{b}_t - \tilde{\mathbf{b}}_t, \tilde{\mathbf{b}}_t - \mathbf{b}_* \rangle &= \mathbb{E}_b^{t-1} \left\langle \frac{\mathbf{x}_t^p - \mathbf{x}_*^p}{\|\mathbf{V}^p\|^2} (\mathbf{V}^p)^*, \mathbf{b}_{t-1} - \mathbf{b}_* + \frac{\mathbf{x}_*^p - \mathbf{V}^p \mathbf{b}_{t-1}}{\|\mathbf{V}^p\|^2} (\mathbf{V}^p)^* \right\rangle \\
&\stackrel{(i)}{=} \mathbb{E}_b^{t-1} \left\langle \frac{\mathbf{x}_t^p - \mathbf{x}_*^p}{\|\mathbf{V}^p\|^2} (\mathbf{V}^p)^*, \frac{\mathbf{x}_*^p - \mathbf{V}^p \mathbf{b}_{t-1}}{\|\mathbf{V}^p\|^2} (\mathbf{V}^p)^* \right\rangle + \mathbb{E}_b^{t-1} \left\langle \frac{\mathbf{x}_t^p - \mathbf{x}_*^p}{\|\mathbf{V}^p\|^2} (\mathbf{V}^p)^*, \mathbf{b}_{t-1} - \mathbf{b}_* \right\rangle \\
&\stackrel{(ii)}{=} \mathbb{E}_b^{t-1} \frac{(\mathbf{x}_t^p - \mathbf{x}_*^p)(\mathbf{x}_*^p - \mathbf{V}^p \mathbf{b}_{t-1})}{\|\mathbf{V}^p\|^2} + \left\langle \mathbb{E}_b^{t-1} \frac{\mathbf{x}_t^p - \mathbf{x}_*^p}{\|\mathbf{V}^p\|^2} (\mathbf{V}^p)^*, \mathbf{b}_{t-1} - \mathbf{b}_* \right\rangle \\
&\stackrel{(iii)}{=} \sum_p \frac{(\mathbf{x}_t^p - \mathbf{x}_*^p)(\mathbf{x}_*^p - \mathbf{V}^p \mathbf{b}_{t-1})}{\|\mathbf{V}^p\|^2} \frac{\|\mathbf{V}^p\|^2}{\|\mathbf{V}\|_F^2} + \left\langle \sum_p \frac{(\mathbf{x}_t^p - \mathbf{x}_*^p)(\mathbf{V}^p)^*}{\|\mathbf{V}^p\|^2} \frac{\|\mathbf{V}^p\|^2}{\|\mathbf{V}\|_F^2}, \mathbf{b}_{t-1} - \mathbf{b}_* \right\rangle \\
&= \frac{(\mathbf{x}_t - \mathbf{x}_*)^* (\mathbf{x}_* - \mathbf{V} \mathbf{b}_{t-1})}{\|\mathbf{V}\|_F^2} + \left\langle \frac{\mathbf{V}^* (\mathbf{x}_t - \mathbf{x}_*)}{\|\mathbf{V}\|_F^2}, \mathbf{b}_{t-1} - \mathbf{b}_* \right\rangle \\
&\stackrel{(iv)}{=} \left\langle \frac{\mathbf{x}_t - \mathbf{x}_*}{\|\mathbf{V}\|_F^2}, \mathbf{V} (\mathbf{b}_* - \mathbf{b}_{t-1}) \right\rangle + \left\langle \frac{\mathbf{x}_t - \mathbf{x}_*}{\|\mathbf{V}\|_F^2}, \mathbf{V} (\mathbf{b}_{t-1} - \mathbf{b}_*) \right\rangle \\
&= 0.
\end{aligned}$$

Step (i) follows from linearity of inner products, step (ii) simplifies the inner product of two parallel vectors, and step (iii) computes the expectation over all possible choices of rows of \mathbf{V} . In step (iv), we use the fact that for $k < m, n$, subsystem (3) is always consistent (since \mathbf{V} is underdetermined) to make the substitution $\mathbf{x}_* = \mathbf{V}\mathbf{b}_*$. To prove statement (b), we note that $(\tilde{\mathbf{b}}_t - \mathbf{b}_{t-1})$ is parallel to \mathbf{V}^p and $(\tilde{\mathbf{b}}_t - \mathbf{b}_*)$ is perpendicular to \mathbf{V}^p since $\mathbf{V}^p(\tilde{\mathbf{b}}_t - \mathbf{b}_*) = \mathbf{V}^p(\mathbf{b}_{t-1} + \frac{\mathbf{x}_*^p - \mathbf{V}^p \mathbf{b}_{t-1}}{\|\mathbf{V}^p\|^2} (\mathbf{V}^p)^* - \mathbf{b}_*) = \mathbf{V}^p \mathbf{b}_{t-1} + \mathbf{x}_*^p - \mathbf{V}^p \mathbf{b}_{t-1} - \mathbf{x}_*^p = 0$. We apply the Pythagorean Theorem to obtain the desired result. \square

Lemma 2. In Algorithm 1 and Algorithm 2, we can bound the expected norm squared error of $\mathbf{b}_t - \mathbf{b}_\star$ as

$$\mathbb{E}^{t-1} \|\mathbf{b}_t - \mathbf{b}_\star\|^2 \leq \alpha_V \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 + \mathbb{E}_x^{t-1} \frac{\|\mathbf{x}_t - \mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2}.$$

We investigate the expectation of the norm squared error of $\mathbf{b}_t - \mathbf{b}_\star$ conditional on the first $t-1$ iterations and over the choice of rows of \mathbf{V} . We keep \mathbb{E}_x^{t-1} in our bound as this expectation will depend on whether Algorithm 1 or Algorithm 2 is being used.

Proof.

$$\begin{aligned} \mathbb{E}^{t-1} \|\mathbf{b}_t - \mathbf{b}_\star\|^2 &= \mathbb{E}^{t-1} \|\mathbf{b}_t - \mathbf{b}_\star + \tilde{\mathbf{b}}_t - \tilde{\mathbf{b}}_t\|^2 \\ &= \mathbb{E}^{t-1} \|\tilde{\mathbf{b}}_t - \mathbf{b}_\star\|^2 + \mathbb{E}^{t-1} \|\mathbf{b}_t - \tilde{\mathbf{b}}_t\|^2 + 2\mathbb{E}^{t-1} \left\langle \tilde{\mathbf{b}}_t - \mathbf{b}_\star, \mathbf{b}_t - \tilde{\mathbf{b}}_t \right\rangle \\ &\stackrel{(iii)}{=} \mathbb{E}^{t-1} \|\tilde{\mathbf{b}}_t - \mathbf{b}_\star\|^2 + \mathbb{E}^{t-1} \|\mathbf{b}_t - \tilde{\mathbf{b}}_t\|^2 \\ &\stackrel{(iv)}{=} \mathbb{E}^{t-1} \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 - \mathbb{E}^{t-1} \|\tilde{\mathbf{b}}_t - \mathbf{b}_{t-1}\|^2 + \mathbb{E}^{t-1} \|\mathbf{b}_t - \tilde{\mathbf{b}}_t\|^2 \\ &\stackrel{(v)}{=} \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 - \mathbb{E}^{t-1} \left\| \frac{(\mathbf{x}_\star^p - \mathbf{V}^p \mathbf{b}_{t-1})}{\|\mathbf{V}^p\|^2} (\mathbf{V}^p)^* \right\|^2 + \mathbb{E}^{t-1} \left\| \frac{(\mathbf{x}_t^p - \mathbf{x}_\star^p)}{\|\mathbf{V}^p\|^2} \right\|^2 (\mathbf{V}^p)^* \\ &= \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 - \mathbb{E}^{t-1} \left[\frac{|\mathbf{V}^p \mathbf{b}_\star - \mathbf{V}^p \mathbf{b}_{t-1}|^2}{\|\mathbf{V}^p\|^2} \right] + \mathbb{E}^{t-1} \left[\frac{|\mathbf{x}_t^p - \mathbf{x}_\star^p|^2}{\|\mathbf{V}^p\|^2} \right]. \end{aligned}$$

Steps (iii) and (iv) are applications of Lemma 1(a) and Lemma 1(b) respectively, and step (v) follows from the definition of each term and simplification using the fact that $\mathbf{V} \mathbf{b}_\star = \mathbf{x}_\star$.

Now, we evaluate the conditional expectation on the choices of rows of \mathbf{V} to complete the proof:

$$\begin{aligned} \mathbb{E}^{t-1} \|\mathbf{b}_t - \mathbf{b}_\star\|^2 &\stackrel{(vi)}{=} \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 - \mathbb{E}_b^{t-1} \left[\frac{|\mathbf{V}^p \mathbf{b}_\star - \mathbf{V}^p \mathbf{b}_{t-1}|^2}{\|\mathbf{V}^p\|^2} \right] + \mathbb{E}_x^{t-1} \mathbb{E}_b^{t-1} \left[\frac{|\mathbf{x}_t^p - \mathbf{x}_\star^p|^2}{\|\mathbf{V}^p\|^2} \right] \\ &= \|\mathbf{b}_{t-1} - \mathbf{x}_\star\|^2 - \sum_{p=1}^k \frac{|\mathbf{V}^p \mathbf{b}_\star - \mathbf{V}^p \mathbf{b}_{t-1}|^2 \|\mathbf{V}^p\|^2}{\|\mathbf{V}^p\|^2 \|\mathbf{V}\|_F^2} + \mathbb{E}_x^{t-1} \sum_{p=1}^k \frac{|\mathbf{x}_t^p - \mathbf{x}_\star^p|^2 \|\mathbf{V}^p\|^2}{\|\mathbf{V}^p\|^2 \|\mathbf{V}\|_F^2} \\ &= \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 - \frac{\|\mathbf{V} \mathbf{b}_\star - \mathbf{V} \mathbf{b}_{t-1}\|^2}{\|\mathbf{V}\|_F^2} + \mathbb{E}_x^{t-1} \left[\frac{\|\mathbf{x}_t - \mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2} \right] \\ &\stackrel{(vii)}{\leq} \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 - \frac{\sigma_{\min}^2(\mathbf{V}) \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2}{\|\mathbf{V}\|_F^2} + \mathbb{E}_x^{t-1} \left[\frac{\|\mathbf{x}_t - \mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2} \right] \\ &= \alpha_V \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 + \mathbb{E}_x^{t-1} \left[\frac{\|\mathbf{x}_t - \mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2} \right]. \end{aligned}$$

In step (vi), we use iterated expectations to split the expected value $\mathbb{E}^{t-1} = \mathbb{E}_x^{t-1} \mathbb{E}_b^{t-1}$. Step (vii) uses the fact that $\|\mathbf{V}(\mathbf{b}_{t-1} - \mathbf{b}_\star)\|^2 \geq \sigma_{\min}^2(\mathbf{V}) \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2$ since $\mathbf{b}_{t-1} - \mathbf{b}_\star$ are in the row span of \mathbf{V} for all t . We simplify and obtain the desired bound. \square

4.3. Proof of main result. We now have all the ingredients we need to prove Theorem 1, which we now proceed to below.

Proof of Theorem 1. The fact that $\mathbf{b}_\star = \beta_\star$ was already argued in scenarios S1 and S3(b) in the previous section, so we do not reproduce its argument here. Given this fact, to prove Theorem 1, we only need to invoke the statement of Lemma 2 and bound the term $\mathbb{E}_x^{t-1} \|\mathbf{x}_t - \mathbf{x}_\star\|^2$ using Proposition 1 or Proposition 2 depending on whether we are using Algorithm 1 or Algorithm 2, respectively.

(a) For Algorithm 1, plugging Proposition 1 into the statement of Lemma 2 yields

$$\mathbb{E}^{t-1} \|\mathbf{b}_t - \mathbf{b}_\star\|^2 \leq \alpha_V \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 + \alpha_U^t \frac{\|\mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2}$$

Taking expectations over the remaining randomness, we have

$$\begin{aligned} \mathbb{E} \|\mathbf{b}_t - \mathbf{b}_\star\|^2 &\leq \alpha_V^t \|\mathbf{b}_\star\|^2 + \alpha_U^t \frac{\|\mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2} \sum_{h=0}^{t-1} \alpha_V^h \\ &\leq \alpha_V^t \|\mathbf{b}_\star\|^2 + \alpha_U^t \frac{\|\mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2} \frac{1}{1 - \alpha_V} \\ &= \alpha_V^t \|\mathbf{b}_\star\|^2 + \theta_V \alpha_U^t \|\mathbf{x}_\star\|^2 \end{aligned}$$

(b) For Algorithm 2, plugging Proposition 2 into the statement of Lemma 2 yields

$$\mathbb{E}^{t-1} \|\mathbf{b}_t - \mathbf{b}_\star\|^2 \leq \alpha_V \|\mathbf{b}_{t-1} - \mathbf{b}_\star\|^2 + (1 + 2\kappa_U^2) \alpha_U^{\lfloor t/2 \rfloor} \frac{\|\mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2}.$$

Taking expectations over the remaining randomness, we have

$$\begin{aligned} \mathbb{E} \|\mathbf{b}_t - \mathbf{b}_\star\|^2 &\leq \alpha_V^t \|\mathbf{b}_\star\|^2 + \alpha_U^{\lfloor t/2 \rfloor} (1 + 2\kappa_U^2) \frac{\|\mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2} \sum_{h=0}^{t-1} \alpha_V^h \\ &\leq \alpha_V^t \|\mathbf{b}_\star\|^2 + \alpha_U^{\lfloor t/2 \rfloor} (1 + 2\kappa_U^2) \frac{\|\mathbf{x}_\star\|^2}{\|\mathbf{V}\|_F^2} \frac{1}{1 - \alpha_V} \\ &= \alpha_V^t \|\mathbf{b}_\star\|^2 + \theta_V \alpha_U^{\lfloor t/2 \rfloor} (1 + 2\kappa_U^2) \|\mathbf{x}_\star\|^2 \end{aligned}$$

This concludes the proof of the theorem. \square

5. EXPERIMENTS

In this section we discuss experiments done on both simulated and real data using different algorithms in different settings. The naming convention for the remainder of the paper will be to refer to ALG1-ALG2 as an interlaced algorithm where ALG1 is the algorithm iterate used to solve subsystem (2) and ALG2 is the algorithm used to solve subsystem (3). When an algorithm's name is used alone, we imply applying the algorithm on the full system (1).

In Figure 1 we show our first set of experiments. Entries of \mathbf{U} , \mathbf{V} , and $\boldsymbol{\beta}$ are drawn from a standard Gaussian distribution. We set $\mathbf{X} = \mathbf{U}\mathbf{V}$ and $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}$ if \mathbf{X} is consistent and $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{r}$ where $\mathbf{r} \in \text{null}(\mathbf{X}^*)$ (computed in Matlab using `null()` function). In this first set of experiments, $m, n, k = \{100, 150, 200\}$ depending on the desired size of k with respect to the over or underdetermined-ness of \mathbf{X} . The plots show iteration vs ℓ_2 -error, $\|\mathbf{b}_t - \boldsymbol{\beta}_\star\|^2$, of each method averaged over 60 runs and allowing each algorithm to run 5×10^4 iterations. The layout of Figure 1 is exactly as in Table 3. For each row, we have a different setting for \mathbf{X} and for each column, we vary the size of k depending on the size of \mathbf{X} . Looking at the overall trends, we see that when $k < m, n$ in addition to when \mathbf{X} is overdetermined, consistent and $n < k < m$, there is a method that obtains the optimal solution for the system. These results align with the expectations set in Table 3. Looking at each individual subplot, we also find what one would expect according to Table 1. In other words, if \mathbf{U} or \mathbf{V} is in one of the settings where RK or RGS are expected to fail then RK-RK or RGS-RGS fail as well.

When \mathbf{X} is overdetermined, inconsistent and $k < m, n$ we have that \mathbf{V} is underdetermined. In this case, we don't need to interlace iterates of REK and REK together. To work on an underdetermined system, using RK is enough to find the optimal solution of that subsystem. This motivated interlacing iterates of RK with REK. Figure 2 has the same set up as discussed in the previous

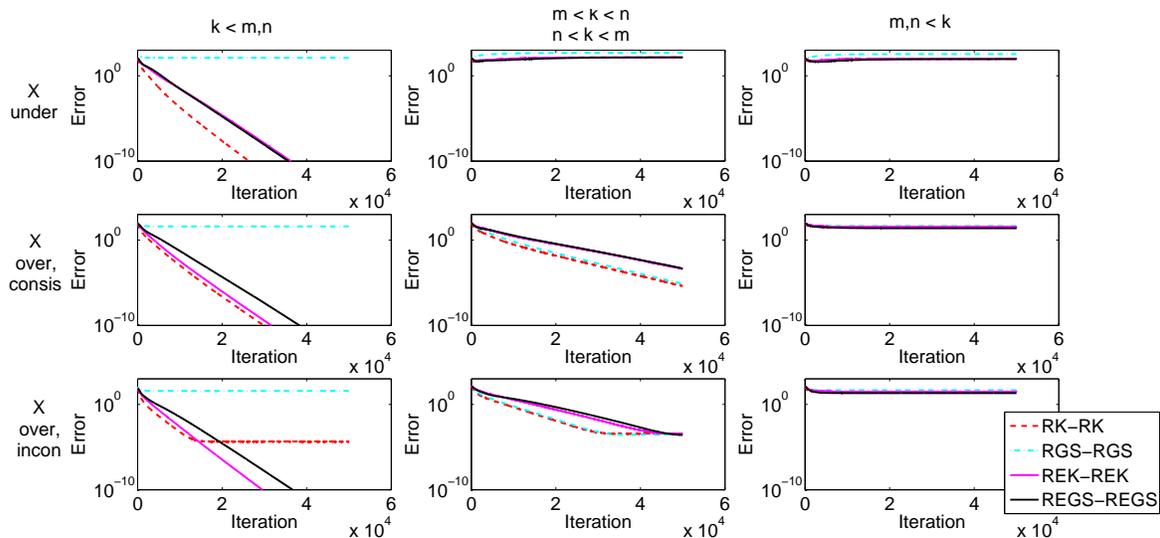


FIGURE 1. This figure shows a summary of convergence for all methods under the variation of possible settings. In the right column, we have that none of the methods convergence. In the middle column, when \mathbf{X} is underdetermined or inconsistent none of the methods converge either. In all other variants, the convergence depends on the general behavior of the standard (RK, RGS) algorithms.

experiment with the exception of using larger random matrices with $\mathbf{X} : 1200 \times 750$ and $k = 500$. In Figure 2a we plot iteration vs ℓ_2 -error and in Figure 2b we plot FLOPS vs ℓ_2 -error. In this experiment we see that REK-RK and REK-REK perform comparably in error and that REK-RK is more efficient in FLOPS.

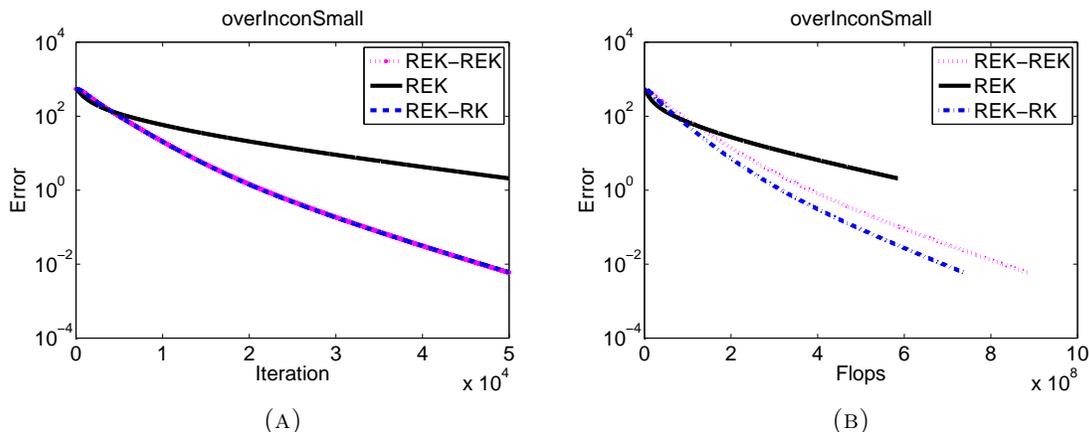


FIGURE 2. When \mathbf{X} is overdetermined, inconsistent and $k < m, n$ we propose interleaving iterates of REK to solve subsystem (2) and RK to solve subsystem (3). This figure demonstrates the advantage of using REK-RK on an inconsistent system as opposed to REK-REK. REK-RK requires less FLOPS to achieve the same accuracy.

In addition to simulated experiments, we also show the usefulness of these algorithms on real world data sets on wine quality, bike rental data, and Yelp reviews. In all following experiments, we plot the average ℓ_2 -error at the t^{th} iteration over 20 runs. The data sets on wine quality and

bike rental data are obtained from the UCI Machine Learning Repository [17]. The wine data set is a sample of $m = 1599$ red wines with $n = 11$ physio-chemical properties of each wine. We choose $k = 5$ and compute \mathbf{U} and \mathbf{V} using Matlab’s `nnmf()` function for nonnegative matrix factorization (recall the motivations from the first section). Figure 3a shows the results from this experiment. The conditioning of \mathbf{X} , \mathbf{U} , and \mathbf{V} are $\kappa_X^2 = 2.46 \times 10^3$, $\kappa_U^2 = 25.96$, and $\kappa_V^2 = 4.20$ respectively. We plot the ℓ_2 -error averaged over 40 runs. Since \mathbf{X} has such a large condition number, this impacts the convergence of REK on \mathbf{X} negatively as shown by the seemingly horizontal line (the error is actually decreasing, but incredibly slowly). We also see that REK-RK and REK-REK are working comparably and significantly faster than REK alone. This can be explained by the better conditioning on \mathbf{U} and \mathbf{V} .

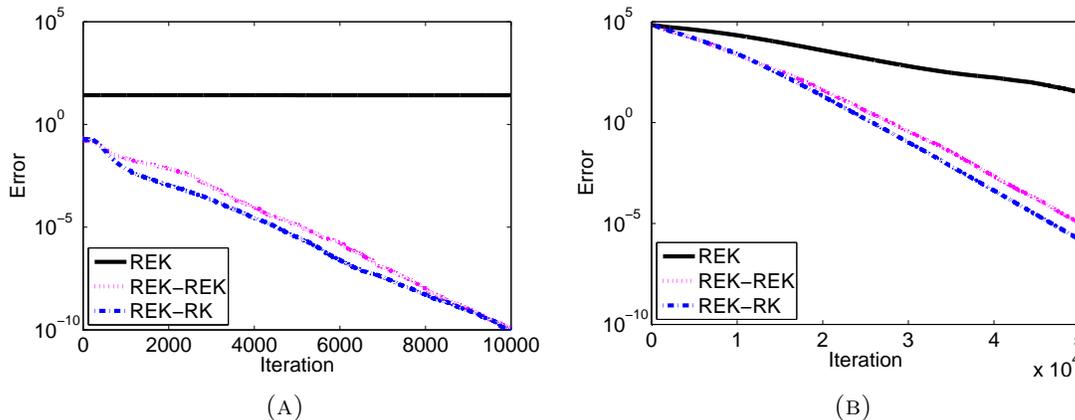


FIGURE 3. In this experiment, we compare the performance of REK, REK-REK, and REK-RK on real world data. Figure 3a shows the performance of these methods on the wine data set and Figure 3b shows performance on the bike data set.

The bike data set contains hourly counts of rental bikes in a bike share system. The data sets contains date as well as weather and seasonal data. There are $m = 17379$ samples and $n = 9$ attributes per sample. We choose $k = 8$ and compute \mathbf{U} and \mathbf{V} in the same way as with the wine data set. Figure 3b shows the results from this experiment. The conditioning of \mathbf{X} , \mathbf{U} , and \mathbf{V} are $\kappa_X^2 = 94.27$, $\kappa_U^2 = 54.91$, and $\kappa_V^2 = 2.99$ respectively. Similar to Figure 3a, we see that the convergence of REK suffers from the poorly conditioned matrix \mathbf{X} . We also see again that REK-REK and REK-RK behave similarly and outperform REK.

To show the advantage of our algorithms on large systems, we create extremely large standard Gaussian matrices $\mathbf{U} : 10^6 \times 10^3$ and $\mathbf{V} : 10^3 \times 10^4$. These matrices are so large that the matrix product \mathbf{UV} cannot be computed in Matlab due to memory constraints. These results are shown in Figure 4. We see that without needing to do the matrix computation, we are still able to find the solution to the linear system $\mathbf{UV}\beta = \mathbf{y}$.

Lastly, we present the performance of our methods on a large real world data set. We use the Yelp challenge data set [1]. In our setting, we let $\mathbf{X} : 10^5 \times 10^4$ be a document term frequency matrix where each row represents a Yelp review and each column represents a word feature. The elements of \mathbf{X} contain the frequency at which the word is used in the review. We only use a subset of the amount of data available due Matlab memory constraints. In this setting, \mathbf{y} is a vector that represents the number of stars a review received. We choose $k = 5000$. Figure 4b shows the results from this experiment using REK, REK-REK, and REK-RK. The conditioning of \mathbf{X} , \mathbf{U} , and \mathbf{V} are $\kappa_X^2 = 127.3592$, $\kappa_U^2 = 24.274$, and $\kappa_V^2 = 19.096$ respectively. In this large real world data set, we can again see the usefulness of our proposed methods when we are given $\mathbf{X} = \mathbf{UV}$.

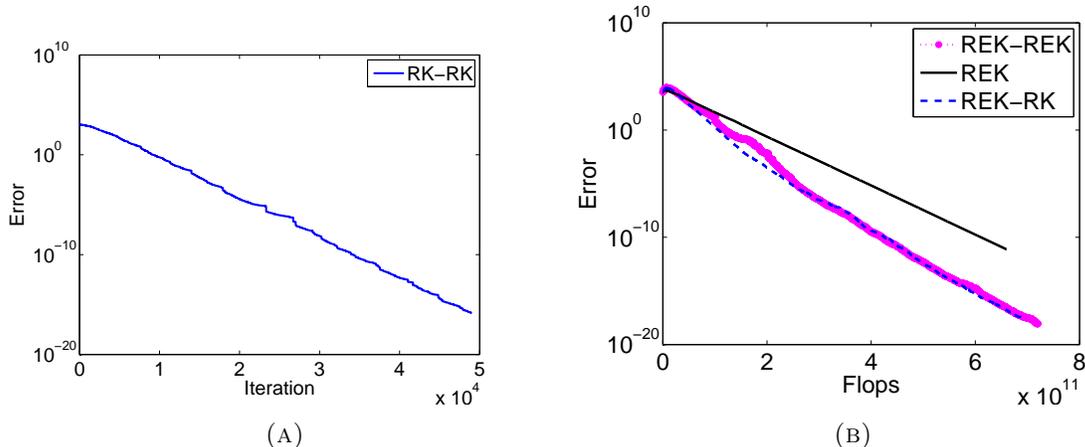


FIGURE 4. We compare the performance of REK, REK-REK, and REK-RK on extremely large datasets. Figure 4a shows results from an experiment that pushes the limits of memory in Matlab. Note that in this experiment, we cannot perform RK on the full system as the matrix product requires too much memory and cannot be formed in Matlab. Figure 4b shows the performance of our method on a large real world dataset.

These experiments complement and verify our theoretical findings. In settings which we expect to fail to obtain the least squares or least norm solutions, our experiments show that they do indeed fail. Additionally, where we expect that the optimal solution is obtainable, the experiments show the proposed methods can obtain such solutions and in many instances outperform the original algorithm on the full system. We see that empirically, subsystems are better conditioned than full systems, thus explaining their better performance.

6. CONCLUSION

We have proposed two methods interlacing Kaczmarz updates to solve factored systems. For large-scale applications in which the system is stored in factored form for efficiency or the factorization arises naturally, our methods allow one to solve the system without the need to perform the large-scale matrix product first. Our main result proves that our methods provide linear convergence in expectation to the (least-squares or least-norm) solution of (overdetermined or underdetermined) linear systems. Our experiments support these results, and show that our methods provide significant computational advantages for factored systems. The interlaced structure of our methods suggests they can be implemented in parallel which would lead to even further computational gains. We leave such details for future work. Additional future work includes the design and analysis of methods that converge to the solution in the settings not covered in this paper, i.e. the gray cells of Table 3. Although its practical implications are not immediately clear to us, these may still be of theoretical interest.

ACKNOWLEDGMENTS

Needell was partially supported by NSF CAREER grant #1348721, and the Alfred P. Sloan Fellowship. Ma was supported in part by NSF CAREER grant #1348721, the CSRC Intellis Fellowship, and the Edison International Scholarship. The authors also thank the Institute of Pure and Applied Mathematics (IPAM) where this collaboration started.

REFERENCES

- [1] *Yelp dataset challenge*, https://www.yelp.com/dataset_challenge.
- [2] G. ADOMAVICIUS AND A. TUZHILIN, *Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions*, IEEE transactions on knowledge and data engineering, 17 (2005), pp. 734–749.
- [3] M. BENZI, *Key moments in the history of numerical analysis*. Presented at 2009 SIAM Applied Linear Algebra Conference, Oct. 2009.
- [4] C. L. BYRNE, *Applied iterative methods*, A K Peters Ltd., Wellesley, MA, 2008.
- [5] E. J. CANDÈS AND B. RECHT, *Exact matrix completion via convex optimization*, Foundations of Computational mathematics, 9 (2009), pp. 717–772.
- [6] B. DUMITRESCU, *On the relation between the randomized extended kaczmarz algorithm and coordinate descent*, BIT Numerical Mathematics, (2014), pp. 1–11.
- [7] J. GOES, T. ZHANG, R. ARORA, AND G. LERMAN, *Robust stochastic principal component analysis.*, in AISTATS, 2014, pp. 266–274.
- [8] R. GORDON, R. BENDER, AND G. T. HERMAN, *Algebraic reconstruction techniques (ART) for three-dimensional electron microscopy and X-ray photography*, J. Theoret. Biol., 29 (1970), pp. 471–481.
- [9] G. T. HERMAN, *Fundamentals of computerized tomography: image reconstruction from projections*, Springer, 2009.
- [10] P. JAIN, P. NETRAPALLI, AND S. SANGHAVI, *Low-rank matrix completion using alternating minimization*, in Proceedings of the forty-fifth annual ACM symposium on Theory of computing, ACM, 2013, pp. 665–674.
- [11] S. KACZMARZ, *Angenäherte auflösung von systemen linearer gleichungen*, Bull. Int. Acad. Polon. Sci. Lett. Ser. A, (1937), pp. 335–357.
- [12] R. H. KESHAVAN, A. MONTANARI, AND S. OH, *Matrix completion from noisy entries*, Journal of Machine Learning Research, 11 (2010), pp. 2057–2078.
- [13] V. KOLTCHINSKII, K. LOUNICI, AND A. B. TSYBAKOV, *Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion*, The Annals of Statistics, (2011), pp. 2302–2329.
- [14] Y. KOREN, R. BELL, C. VOLINSKY, ET AL., *Matrix factorization techniques for recommender systems*, Computer, 42 (2009), pp. 30–37.
- [15] D. D. LEE AND H. S. SEUNG, *Algorithms for non-negative matrix factorization*, in Advances in neural information processing systems, 2001, pp. 556–562.
- [16] D. LEVENTHAL AND A. S. LEWIS, *Randomized methods for linear constraints: convergence rates and conditioning*, Math. Oper. Res., 35 (2010), pp. 641–654, <https://doi.org/10.1287/moor.1100.0456>, <http://dx.doi.org/10.1287/moor.1100.0456>.
- [17] M. LICHMAN, *UCI machine learning repository*, 2013, <http://archive.ics.uci.edu/ml>.
- [18] A. MA, A. FLENNER, D. NEEDELL, AND A. G. PERCUS, *Improving image clustering using sparse text and the wisdom of the crowds*, in 2014 48th Asilomar Conference on Signals, Systems and Computers, IEEE, 2014, pp. 1555–1557.
- [19] A. MA, D. NEEDELL, AND A. RAMDAS, *Convergence properties of the randomized extended gauss–seidel and kaczmarz methods*, SIAM Journal on Matrix Analysis and Applications, 36 (2015), pp. 1590–1604.
- [20] H. MA, D. ZHOU, C. LIU, M. R. LYU, AND I. KING, *Recommender systems with social regularization*, in Proceedings of the fourth ACM international conference on Web search and data mining, ACM, 2011, pp. 287–296.

- [21] F. NATTERER, *The mathematics of computerized tomography*, vol. 32 of Classics in Applied Mathematics, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2001, <https://doi.org/10.1137/1.9780898719284>, <http://dx.doi.org/10.1137/1.9780898719284>. Reprint of the 1986 original.
- [22] D. NEEDELL, *Randomized Kaczmarz solver for noisy linear systems*, BIT, 50 (2010), pp. 395–403, <https://doi.org/10.1007/s10543-010-0265-5>, <http://dx.doi.org/10.1007/s10543-010-0265-5>.
- [23] D. NEEDELL, N. SBRERO, AND R. WARD, *Stochastic gradient descent and the randomized kaczmarz algorithm*, Math. Program. Series A, (2014). to appear.
- [24] T. STROHMER AND R. VERSHYNIN, *Comments on the randomized Kaczmarz method*, J. Fourier Anal. Appl., 15 (2009), pp. 437–440, <https://doi.org/10.1007/s00041-009-9082-0>, <http://dx.doi.org/10.1007/s00041-009-9082-0>.
- [25] T. STROHMER AND R. VERSHYNIN, *A randomized Kaczmarz algorithm with exponential convergence*, J. Fourier Anal. Appl., 15 (2009), pp. 262–278, <https://doi.org/10.1007/s00041-008-9030-4>, <http://dx.doi.org/10.1007/s00041-008-9030-4>.
- [26] G. TAKÁCS, I. PILÁSZY, B. NÉMETH, AND D. TIKK, *Investigation of various matrix factorization methods for large recommender systems*, in 2008 IEEE International Conference on Data Mining Workshops, IEEE, 2008, pp. 553–562.
- [27] K. VERBERT, N. MANOUSELIS, X. OCHOA, M. WOLPERS, H. DRACHSLER, I. BOSNIC, AND E. DUVAL, *Context-aware recommender systems for learning: a survey and future challenges*, IEEE Transactions on Learning Technologies, 5 (2012), pp. 318–335.
- [28] J. WRIGHT, A. GANESH, S. RAO, Y. PENG, AND Y. MA, *Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization*, in Advances in neural information processing systems, 2009, pp. 2080–2088.
- [29] W. XU, X. LIU, AND Y. GONG, *Document clustering based on non-negative matrix factorization*, in Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval, ACM, 2003, pp. 267–273.
- [30] A. ZOUZIAS AND N. M. FRERIS, *Randomized extended Kaczmarz for solving least squares*, SIAM Journal on Matrix Analysis and Applications, 34 (2013), pp. 773–793.