

Greed Works: An Improved Analysis of Sampling Kaczmarz-Motzkin*

Jamie Haddock[†] and Anna Ma[‡]

Abstract. Stochastic iterative algorithms have gained recent interest in machine learning and signal processing for solving large-scale systems of equations, $A\mathbf{x} = \mathbf{b}$. One such example is the Randomized Kaczmarz (RK) algorithm, which acts only on single rows of the matrix A at a time. While RK randomly selects a row of A to work with, Motzkin’s Method (MM) employs a greedy row selection. Connections between the two algorithms resulted in the Sampling Kaczmarz-Motzkin (SKM) algorithm which samples a random subset of β rows of A and then greedily selects the best row of the subset. Despite their variable computational costs, all three algorithms have been proven to have the same theoretical upper bound on the convergence rate. In this work, an improved analysis of the range of random (RK) to greedy (MM) methods is presented. This analysis improves upon previous known convergence bounds for SKM, capturing the benefit of partially greedy selection schemes. This work also further generalizes previous known results, removing the theoretical assumptions that β must be fixed at every iteration and that A must have normalized rows.

Key words. Kaczmarz method, iterative methods, greedy methods, randomization

AMS subject classifications. 65F10, 65F20

1. Introduction. Large-scale systems of equations arise in many areas of data science, including in machine learning and as subroutines of several optimization methods [12]. We consider solving these large systems of linear equations, $A\mathbf{x} = \mathbf{b}$, where $A \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$, and $m \gg n$. Iterative methods which use a small portion of the data in each iteration are typically employed in this domain. These methods offer a small memory footprint and good convergence guarantees. The Kaczmarz method [36] is such an iterative method that consists of sequential orthogonal projections towards the solution set of a single equation (or subsystem). Given the system $A\mathbf{x} = \mathbf{b}$, the method computes iterates by projecting onto the hyperplane defined by the equation $\mathbf{a}_i^T \mathbf{x} = b_i$ where \mathbf{a}_i^T is a selected row of the matrix A and b_i is the corresponding entry of \mathbf{b} . The iterates are recursively defined as

$$\mathbf{x}_{j+1} = \mathbf{x}_j + \frac{b_i - \mathbf{a}_i^T \mathbf{x}_j}{\|\mathbf{a}_i\|^2} \mathbf{a}_i.$$

We assume that $A\mathbf{x} = \mathbf{b}$ is consistent and $m > n$, but make no assumption on $\text{rank}(A)$. We will use $\mathbf{r}_j := A\mathbf{x}_j - \mathbf{b}$ to represent the j th residual and $\mathbf{e}_j := \mathbf{x}_j - \mathbf{x}^*$ to represent the j th error term. We let A^\dagger denote the Moore-Penrose pseudoinverse of the matrix A . Additionally,

*Submitted to the editors X.

Funding: This material was supported by the NSF DMS-1440140 while the authors were in residence at the Mathematical Science Research Institute in Berkeley, California, during the Fall 2017 semester. Both authors were partially supported by NSF CAREER DMS-1348721 and NSF BIGDATA 1740325. JH was partially funded by NSF DMS-1522158 and NSF DMS-1818969. AM was partially supported the U.S. Air Force Award FA9550-18-1-0031 led by Roman Vershynin.

[†]Department of Mathematics, University of California, Los Angeles, Los Angeles, CA (jhaddock@math.ucla.edu).

[‡]Department of Mathematics, University of California, Irvine, Irvine, CA (anna.ma@uci.edu).

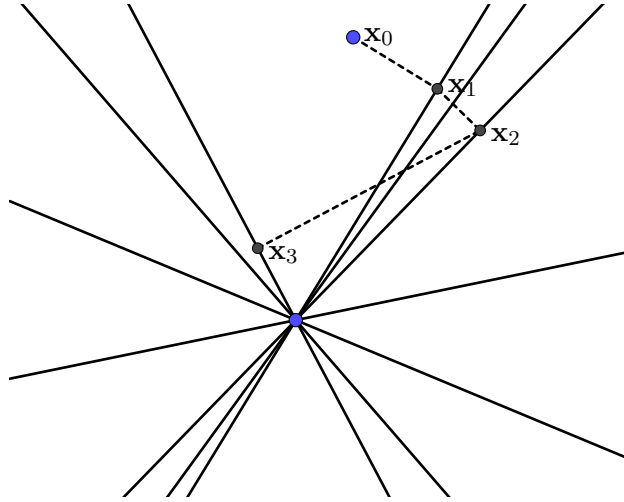


Figure 1.1: Several iterations of a Kaczmarz method. The iterate \mathbf{x}_{j+1} is the orthogonal projection of \mathbf{x}_j onto the solution set of the selected equation (represented by a line).

21 we let $\sigma_{\min}(A)$ be the smallest nonzero singular value of A and unless otherwise noted, we let
 22 $\|\cdot\|$ represent the Euclidean norm. We let $\|\cdot\|_F$ denote the Frobenius norm and $\|\cdot\|_\infty$ denote
 23 the ℓ^∞ norm. A visualization of several iterations of a Kaczmarz method are shown in Figure
 24 1.1.

25 The Kaczmarz method was originally proposed in the late 30s [36] and rediscovered in the
 26 1970's under the name *algebraic reconstruction technique (ART)* as an iterative method for
 27 reconstructing an image from a series of angular projections in computed tomography [27, 35].
 28 This method has seen popularity among practitioners and researchers alike since the beginning
 29 of the digital age [14, 33], but saw a renewed surge of interest after the elegant convergence
 30 analysis of the *Randomized Kaczmarz (RK) method* in [60]. In [60], the authors showed
 31 that for a consistent system with unique solution, RK (with specified sampling distribution)
 32 converges at least linearly in expectation with the guarantee

$$33 \quad (1.1) \quad \mathbb{E}\|\mathbf{e}_k\|^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{\|A\|_F^2}\right)^k \|\mathbf{e}_0\|^2.$$

34 Many variants and extensions followed, including convergence analyses for inconsistent and
 35 random linear systems [49, 15], connections to other popular iterative algorithms [44, 51, 56,
 36 57, 21], block approaches [52, 58], acceleration and parallelization strategies [22, 37, 47, 45],
 37 and techniques for reducing noise and corruption [68, 32].

38 Another popular Kaczmarz method extension is greedy (rather than randomized) row
 39 selection, which has been rediscovered several times in the literature as the “most violated
 40 constraint control” or the “maximal-residual control” [13, 54, 55]. This method was proposed
 41 in the 1950's as an iterative relaxation method for linear programming by Agmon, Motzkin,
 42 and Schoenberg under the name *Motzkin's relaxation method for linear inequalities (MM)*
 43 [48, 1]. In [1], the author showed that MM converges at least linearly (deterministically) with

44 the convergence rate of (1.1). The bodies of literature studying this greedy strategy have
 45 remained somewhat disjoint, with analyses for linear systems of equations in the numerical
 46 linear algebra community and analyses for linear systems of inequalities in the operations
 47 research and linear programming community [25, 26, 62, 4, 8, 9, 16]. There has been recent
 48 work in analyzing variants of this greedy strategy [20, 6, 7, 59]. In [59], the authors analyze
 49 MM on a system to which a Gaussian sketch has been applied. In [6, 7], the authors analyze
 50 variants of MM in which the equation selected in each iteration is chosen randomly amongst
 51 the set whose residual values are sufficiently near the maximal residual value. In [20], the
 52 authors provide a convergence analysis for a generalized version of MM in which the equation
 53 chosen in each iteration is that which has the maximal *weighted residual* value which are the
 54 residual values divided by the norm of the corresponding row of the measurement matrix.
 55 In [19], the authors illustrated the connection between MM and RK and proposed a family
 56 of algorithms that interpolate between the two, known as the *Sampling Kaczmarz-Motzkin*
 57 (*SKM*) *methods*.

58 The SKM methods operate by randomly sampling a subset of the system of equations,
 59 computing the residual of this subset, and projecting onto the equation corresponding to the
 60 largest magnitude entry of this sub-residual. The family of methods (parameterized by the size
 61 of the random sample of equations, β) interpolates between MM, which is SKM with $\beta = m$,
 62 and RK, which is SKM with $\beta = 1$. In [19], the authors prove that the SKM methods converge
 63 at least linearly in expectation with the convergence rate specified in (1.1). Meanwhile, the
 64 empirical convergence of this method is seen to depend upon β ; however, increasing β also
 65 increases the computational cost of each iteration so the per iteration gain from larger sample
 66 size may be outweighed by the in-iteration cost. This is reminiscent of other methods which
 67 use sampled subsets of data in each iteration, such as the block projection methods [2, 52, 53].

Like SKM, the randomized block Kaczmarz (RBK) methods use a subset of rows $\tau \subset [m]$
 to produce the next iterate; rather than forcing the next iterate to satisfy the single sampled
 equation as in RK, block iterates satisfy *all* the equations in the randomly sampled block.
 The $(k + 1)$ st RBK iteration is given by

$$\mathbf{x}_{k+1} = \mathbf{x}_k + (A_\tau)^\dagger(\mathbf{b}_\tau - A_\tau \mathbf{x}_k),$$

68 where A_τ and \mathbf{b}_τ represent the restriction onto the row indices in τ . In [52], the authors
 69 prove that on a system with a row-normalized measurement matrix and a well-conditioned
 70 row-paving RBK converges at least linearly in expectation with the guarantee

$$71 \quad (1.2) \quad \mathbb{E}\|\mathbf{e}_k\|^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{C\|A\|^2 \log(m+1)}\right)^k \|\mathbf{e}_0\|^2.$$

72 where C is an absolute constant and $\|A\|$ denotes the operator norm of the matrix. This can
 73 be a significant improvement over the convergence rate of (1.1) when $\|A\|_F^2 \gg \|A\|^2 \log(m+1)$.
 74 However, the cost per iteration scales with the size of the blocks. In [53], the authors generalize
 75 this result to inconsistent systems and show that, up to a convergence horizon, RBK converges
 76 to the least-squares solution.

77 In [29, 30], the authors introduce a framework of iterative methods known as the *sketch-*
 78 *and-project* methods. The sketch-and-project framework of methods produce each new it-
 79 erate by projecting the previous iterate onto a sketch of the linear system; i.e., $\mathbf{x}_{k+1} =$

80 $\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{x} - \mathbf{x}_k\|_B^2$ s.t. $S_k^\top A \mathbf{x} = S_k^\top \mathbf{b}$. Subsequent works proposed and analyzed vari-
 81 ants with momentum [40], inexact variants [42], and adaptive variants [28]. This framework
 82 includes as special cases many forms of row- and column-action methods and second-order
 83 iterative least-squares methods [64]. Kaczmarz methods which iteratively project onto the
 84 solution spaces of subsets of rows in each iteration (like Block RK or SKM) can be interpreted
 85 and analyzed in this framework. Single row-action methods are recovered when the sketch-
 86 ing matrices select a single row of the system. The SKM methods are recovered when the
 87 sketching matrices select a single row of the system and the choice of which sketch to use in
 88 each iteration is made in the same way as SKM (a randomized sample then a greedy selection
 89 based upon sketched residual). The results recovered from [29] for this interpretation of SKM
 90 coincide with (1.1).

91 In [50], the authors, inspired by the sketching framework in [29], construct a block-type
 92 method which iterates by projecting onto a Gaussian sketch of the equations. They show that
 93 this method converges at least linearly in expectation with the guarantee

$$94 \quad (1.3) \quad \mathbb{E} \|\mathbf{e}_k\|^2 \leq \left(1 - \left[\frac{\sqrt{s} \sigma_{\min}(A)}{9\sqrt{s} \|A\| + C \|A\|_F} \right]^2 \right)^k \|\mathbf{e}_0\|^2.$$

95 where C is an absolute constant and s is the number of rows in the resulting sketched system.
 96 This result requires a Gaussian sketch which is a costly operation, however the authors suggest
 97 using a Gaussian sketch of only a subset of the equations. This result is most related to SKM
 98 and to our main result due to the presence of s , the size of the sketched system, in the bound.

99 **2. Previous Results.** This section focuses on the convergence behavior of the RK, MM,
 100 and SKM methods. Each of these projection methods is a special case of Algorithm 2.1
 101 with a different *selection rule* (Line 4). In iteration j , RK uses the randomized selection
 102 rule that chooses $t_j = i$ with probability $\|\mathbf{a}_i\|_2^2 / \|A\|_F^2$, MM uses the greedy selection rule
 103 $t_j = \arg \max_i |\mathbf{a}_i^\top \mathbf{x}_{j-1} - b_i|$, and SKM uses the hybrid selection rule that first samples a subset
 104 of β rows, τ_j , uniformly at random from all subsets of size β , $\tau_j \sim \operatorname{unif}(\binom{[m]}{\beta})$, and then
 105 chooses $t_j = \arg \max_{i \in \tau_j} |\mathbf{a}_i^\top \mathbf{x}_{j-1} - b_i|$. As previously mentioned, RK and MM are special
 106 cases of the SKM method when the sample size $\beta = 1$ and $\beta = m$, respectively. Each of the
 107 methods converge linearly when the system is consistent with unique solution (RK and SKM
 108 converge linearly in expectation, MM converges linearly deterministically). In Table 2.1, we
 109 present the selection rules and convergence rates for RK, MM, and SKM. Note that under
 110 the assumption that A has been normalized so that $\|\mathbf{a}_i\|^2 = 1$, each of these upper bounds
 111 on the convergence rate is the same since $\|A\|_F^2 = m$. Thus, these results do not reveal any
 112 advantage the more computationally expensive methods (MM, SKM with $\beta \gg 1$) enjoy over
 113 RK. There are, in fact, pathological examples on which RK, MM, and SKM exhibit nearly
 114 the same behavior (e.g., consider the system defining two lines that intersect at one point in
 115 \mathbb{R}^2), so it is not possible to prove significantly different convergence rates without leveraging
 116 additional properties of the system.

In [31], the authors demonstrate that MM can converge faster than RK or SKM and
 that the convergence rate depends on the structure of the residual terms of the iterations,

Algorithm 2.1 Generic Kaczmarz Method

```

1: procedure KACZ( $A, \mathbf{b}, \mathbf{x}_0$ )
2:    $k = 1$ 
3:   repeat
4:     Choose  $t_k \in [m]$  according to selection rule.
5:      $\mathbf{x}_k = \mathbf{x}_{k-1} - \frac{\mathbf{a}_{t_k}^\top \mathbf{x}_{k-1} - b_{t_k}}{\|\mathbf{a}_{t_k}\|_2^2} \mathbf{a}_{t_k}$ .
6:      $k = k + 1$ 
7:   until stopping criterion reached
8:   return  $\mathbf{x}_k$ 
9: end procedure

```

	Selection Rule	Convergence Rate
RK [60]	$\mathbb{P}(t_j = i) = \frac{\ \mathbf{a}_i\ _2^2}{\ A\ _F^2}$	$\mathbb{E}\ \mathbf{e}_k\ ^2 \leq (1 - \frac{\sigma_{\min}^2(A)}{\ A\ _F^2})^k \ \mathbf{e}_0\ ^2$
SKM [19]	$\tau_j \sim \text{unif}(\binom{[m]}{\beta})$ $t_j = \arg \max_{i \in \tau_j} \mathbf{a}_i^\top \mathbf{x}_{j-1} - b_i $	$\mathbb{E}\ \mathbf{e}_k\ ^2 \leq (1 - \frac{\sigma_{\min}^2(A)}{m})^k \ \mathbf{e}_0\ ^2$
MM [1]	$t_j = \arg \max_i \mathbf{a}_i^\top \mathbf{x}_{j-1} - b_i $	$\ \mathbf{e}_k\ ^2 \leq (1 - \frac{\sigma_{\min}^2(A)}{m})^k \ \mathbf{e}_0\ ^2$

Table 2.1: The selection rules and convergence rates of RK, SKM, and MM. The presented results for MM and SKM assume that A has been normalized so that $\|\mathbf{a}_i\|^2 = 1$.

$\mathbf{r}_k = A\mathbf{x}_k - \mathbf{b}$. In particular, they prove that

$$\|\mathbf{e}_k\|^2 \leq \prod_{j=0}^{k-1} \left(1 - \frac{\sigma_{\min}^2(A)}{4\gamma_j}\right) \|\mathbf{e}_0\|^2,$$

117 where γ_j is the *dynamic range* of the i th residual, $\gamma_j := \frac{\|\mathbf{r}_j\|^2}{\|\mathbf{r}_j\|_\infty^2}$. Our main contribution in this
118 paper is to prove that the SKM methods can exhibit a similarly accelerated convergence rate
119 and the advantage scales with the size of the sample, β . Again, this advantage depends upon
120 the structure of the residuals of the iterations. We define here a generalization of the dynamic
121 range used in [31]; our dynamic range is defined as

$$122 \quad (2.1) \quad \gamma_j = \frac{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_{j-1} - \mathbf{b}_\tau\|^2}{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_{j-1} - \mathbf{b}_\tau\|_\infty^2}.$$

123 Now, we let \mathbb{E}_{τ_j} denote expectation with respect to the random sample τ_j conditioned upon
124 the sampled τ_i for $i < j$, and \mathbb{E} denote expectation with respect to all random samples τ_i for
125 $1 \leq i \leq j$ where j is understood to be the last iteration in the context in which \mathbb{E} is applied.
126 We state our main result below in Corollary 2.1; this is a corollary of our generalized result
127 which will be discussed and proven later.

Corollary 2.1. *Let A be normalized so $\|\mathbf{a}_i\| = 1$ for all rows $i = 1, \dots, m$. Suppose the system of equations $A\mathbf{x} = \mathbf{b}$ is consistent, define $\mathbf{x}^* = A^\dagger \mathbf{b}$, and let $\mathbf{x}_0 \in \text{range}(A^\top)$. Then SKM*

converges at least linearly in expectation and the bound on the rate depends on the dynamic range, γ_k of the random sample of β rows of A , τ_k . Precisely, in the k th iteration of SKM, we have

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta \sigma_{\min}^2(A)}{\gamma_k m}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2,$$

128 so applying expectation with respect to all iterations, we have

$$129 \quad \mathbb{E} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \prod_{j=1}^k \left(1 - \frac{\beta \sigma_{\min}^2(A)}{\gamma_j m}\right) \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

130

Corollary 2.1 shows that SKM experiences at least linear convergence where the contraction term is a product of terms that are less than one and dependent on the sub-sample size β . When $\beta = 1$, as in RK, $\gamma_k = 1$, so Corollary 2.1 recovers the upper bound for RK shown in [60]. However, when $\beta = m$ for MM, Corollary 2.1 offers an improved upper bound on the error over [31]; specifically

$$\|\mathbf{e}_k\|^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{\gamma_k}\right) \|\mathbf{e}_{k-1}\|^2.$$

Our result illustrates that the progress made by an iteration of the SKM algorithm depends upon the dynamic range of the residual of that iteration. The dynamic range of each iteration, γ_j , satisfies

$$1 \leq \gamma_j \leq \beta.$$

Note that the upper bound, $\gamma_j = \beta$, is achieved by a constant residual where $|\mathbf{a}_i^T \mathbf{x}_j - b_i| = |\mathbf{a}_{i'}^T \mathbf{x}_j - b_{i'}|$ for all $i, i' \in [m]$, while the lower bound is achieved by the residual with one nonzero entry. As smaller γ_j provides a smaller upper bound on the new error \mathbf{e}_j , we consider the situation with one nonzero entry in the residual as the “best case” and the situation with a constant residual as the “worst case.” We now compare our single iteration result in the best and worst cases to the previously known single iteration results of [1, 19, 31, 60]. These are summarized in Table 2.2; we present only the contraction terms α such that

$$\mathbb{E}_{\tau_k} \|\mathbf{e}_k\|^2 \leq \alpha \|\mathbf{e}_{k-1}\|^2,$$

131 for each upper bound in the case that A is normalized so that $\|\mathbf{a}_i\|^2 = 1$ for $i \in [m]$. In
 132 particular, note that the worst case residual provides the same upper bound rate as those of
 133 [60, 19, 1].

134 **3. Main Results.** Corollary 2.1 is a specialization of our general result to SKM with a
 135 fixed sample size β and systems that are row-normalized. Our general result requires neither
 136 row-normalization nor a static sample size. However, we must additionally generalize the
 137 SKM sampling distribution for systems that are not row-normalized. We now consider the
 138 general SKM method which samples β_k many rows of A in the k th iteration (according to
 139 probability distribution $p_{\mathbf{x}_{k-1}}$ defined in (3.1)) and projects onto the hyperplane associated to
 140 the largest magnitude entry of the sampled sub-residual.

	Best Case	Worst Case	Previous Best Case	Previous Worst Case
MM	$1 - \sigma_{\min}^2(A)$	$1 - \frac{\sigma_{\min}^2(A)}{m}$	$1 - \frac{\sigma_{\min}^2(A)}{4}$ [31]	$1 - \frac{\sigma_{\min}^2(A)}{m}$ [1, 19, 60]
SKM	$1 - \frac{\beta \sigma_{\min}^2(A)}{m}$		$1 - \frac{\sigma_{\min}^2(A)}{m}$	
RK	$1 - \frac{\sigma_{\min}^2(A)}{m}$			

Table 2.2: Contraction terms α such that $\mathbb{E}_{\tau_k} \|\mathbf{e}_k\|^2 \leq \alpha \|\mathbf{e}_{k-1}\|^2$ for the best and worst case bounds of MM, SKM, and RK.

141 The generalized probability distribution over the subset of rows of A of size β_k is denoted
 142 $p_{\mathbf{x}} : \binom{[m]}{\beta_k} \rightarrow [0, 1)$. The sampled subset of rows of A , $\tau_k \sim p_{\mathbf{x}}$ where

$$143 \quad (3.1) \quad p_{\mathbf{x}}(\tau_k) = \frac{\|\mathbf{a}_{t(\tau_k, \mathbf{x})}\|^2}{\sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x})}\|^2},$$

144 and $t(\tau, \mathbf{x}) = \arg \max_{t \in \tau} (\mathbf{a}_t^\top \mathbf{x} - b_t)^2$. Thus, our generalized SKM method is Algorithm 2.1 with
 145 selection rule $\tau_j \sim p_{\mathbf{x}_{j-1}}$ and $t_j = t(\tau_j, \mathbf{x}_{j-1})$. Similar to the RK probability distribution of [60],
 146 the computation of (3.1) is utilized here simply to theoretically analyze the SKM algorithm
 147 without requiring normalized rows. This choice of sampling distribution conveniently simplifies
 148 the expected value computation in the proof of Theorem 3.1 by cancelling the numerator of the
 149 probability with the squared norm of the sampled row. We do not suggest that this probability
 150 distribution be implemented in a real world setting as it is computationally prohibitive.

151 One could instead implement a uniform distribution over rows or learn the distribution
 152 with probabilities proportional to the squared norms of the rows (as suggested in [60]). Neither
 153 of these is guaranteed to coincide with the distribution defined in (3.1), due to the dependence
 154 on the iterate \mathbf{x} . However, for many datasets where the row norms are (approximately) equal,
 155 the uniform distribution (approximately) coincides with (3.1). In particular, when the rows
 156 of A all have equal norm, as in the case of incidence matrices (see Section 4.2), then (3.1)
 157 reduces to the uniform distribution over samples of size β_k . Past works which analyze SKM
 158 [19, 47, 46] assume that the rows of A are normalized and that the probability distribution
 159 over the samples of size β is uniform. To the best of our knowledge, ours is the first work
 160 in this area to analyze an iterative projection method with an iteration dependent sampling
 161 distribution.

162 Our main result shows that the generalized SKM converges at least linearly in expectation
 163 with a bound that depends on the dynamic range of the sampled sub-residual, the size of the
 164 sample, and the minimum squared nonzero singular value of A , $\sigma_{\min}^2(A)$. In the event that
 165 there are multiple rows within the sub-residual which achieve $\max_{t \in \tau} (\mathbf{a}_t^\top \mathbf{x} - b_t)^2$, an arbitrary
 166 choice can be made amongst those rows and the main result will not be affected by this choice.

167 Theorem 3.1 provides theoretical convergence guarantees for the generalized SKM method.
 168 Whereas previous guarantees for SKM required normalized rows or fixed sample sizes β [19,
 169 31, 47, 46], the guarantees presented here do not require either assumption. In addition, the
 170 contraction term of the generalized SKM method shows dependence on the dynamic range,
 171 another feature lacking in previous works. Following the statement of the theorem, we use

172 standard techniques in the Kaczmarz literature to prove our main result.

173 We additionally describe a simple generalization of the main result in the case that the
 174 samples are not made according to the generalized SKM distribution (3.1), but instead ac-
 175 cording to a distribution $\tilde{p}(\tau)$ whose probabilities are at least a constant factor of those in
 176 (3.1). We also remark on the special case in which rows have equal norm (and thus subsets τ
 177 are selected uniformly at random), the case where β_k is fixed, and the case in which $\beta_k = 1$
 178 in order to make connections to previous results. Due to the dependence of the sampling dis-
 179 tribution upon the current iterate, our main result is not easily iterable to provide the usual
 180 form of a Kaczmarz type result (e.g., $\mathbb{E}\|\mathbf{e}_k\|^2 \leq \alpha^k \|\mathbf{e}_0\|^2$) so we present the bound for only a
 181 single iteration. However, in the special cases we describe in Remarks 2 and 3 we are able to
 182 iterate the simplified expression due to the simplicity of the sampling distribution.

183 **Theorem 3.1.** *Suppose the system of equations $A\mathbf{x} = \mathbf{b}$ is consistent, define $\mathbf{x}^* = A^\dagger \mathbf{b}$, and*
 184 *let $\mathbf{x}_0 \in \text{range}(A^\top)$. Then generalized SKM converges at least linearly in expectation and the*
 185 *bound on the rate depends on the dynamic range, γ_k of the random sample of β_k rows of A ,*
 186 *τ_k . Precisely, in the k th iteration of generalized SKM, we have*

$$187 \quad \mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \binom{m}{\beta_k} \sigma_{\min}^2(A)}{\gamma_k m \sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2} \right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2.$$

188

189 *Proof.* We begin by rewriting the generalized SKM iterate \mathbf{x}_k and simplifying the resulting
 190 expression which yields

$$191 \quad \begin{aligned} \|\mathbf{x}_k - \mathbf{x}^*\|^2 &= \left\| \mathbf{x}_{k-1} - \frac{\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}^\top \mathbf{x}_{k-1} - b_{t(\tau_k, \mathbf{x}_{k-1})}}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2} \mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})} - \mathbf{x}^* \right\|^2 \\ 192 \quad &= \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \frac{(\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}^\top \mathbf{x}_{k-1} - b_{t(\tau_k, \mathbf{x}_{k-1})})^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2} \\ 193 \quad &= \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \frac{\|A_{\tau_k} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_k}\|_\infty^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2}, \\ 194 \end{aligned}$$

where the first equation uses the definition of the generalized SKM iterate, and the second follows from the fact that $\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}^\top (\mathbf{x}_{k-1} - \mathbf{x}^*) = \mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}^\top \mathbf{x}_{k-1} - b_{t(\tau_k, \mathbf{x}_{k-1})}$. Note that

$$\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}^\top \mathbf{x}^* = (AA^\dagger \mathbf{b})_{t(\tau_k, \mathbf{x}_{k-1})} = b_{t(\tau_k, \mathbf{x}_{k-1})}$$

195 since $\mathbf{b} \in \text{range}(A)$ and AA^\dagger is the projector onto $\text{range}(A)$.

196 Now, we take expectation of both sides (with respect to the sampled τ_k according to the

197 distribution (3.1)). This gives

$$\begin{aligned}
198 \quad \mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 &= \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \mathbb{E}_{\tau_k} \frac{\|A_{\tau_k} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_k}\|_\infty^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2} \\
199 \quad &= \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \sum_{\tau \in \binom{[m]}{\beta_k}} p_{\mathbf{x}_{k-1}}(\tau) \frac{\|A_\tau \mathbf{x}_{k-1} - \mathbf{b}_\tau\|_\infty^2}{\|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2} \\
200 \quad (3.2) \quad &= \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \sum_{\tau \in \binom{[m]}{\beta_k}} \frac{\|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}{\sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \frac{\|A_\tau \mathbf{x}_{k-1} - \mathbf{b}_\tau\|_\infty^2}{\|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2} \\
201 \quad &= \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \frac{1}{\sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \sum_{\tau \in \binom{[m]}{\beta_k}} \|A_\tau \mathbf{x}_{k-1} - \mathbf{b}_\tau\|_\infty^2 \\
202 \quad &= \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \frac{1}{\gamma_k \sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \sum_{\tau \in \binom{[m]}{\beta_k}} \|A_\tau \mathbf{x}_{k-1} - \mathbf{b}_\tau\|_\infty^2 \\
203 \quad &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - \frac{\binom{m}{\beta_k} \beta_k}{\gamma_k m \sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \|A \mathbf{x}_{k-1} - \mathbf{b}\|^2 \\
204 \quad &\leq \left(1 - \frac{\binom{m}{\beta_k} \beta_k \sigma_{\min}^2(A)}{\gamma_k m \sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2, \\
205 \quad &
\end{aligned}$$

206 where the last line follows from standard properties of singular values and the fact that
207 $\mathbf{x}_{k-1} \in \text{range}(A^\top)$ (since $\mathbf{x}_0 \in \text{range}(A^\top)$ and the SKM update preserves membership in
208 $\text{range}(A^\top)$). This completes our proof. ■

209 Now, we provide a corollary of the previous result which provides a bound on the expected
210 error for the SKM algorithm which samples subsets of rows τ_k according to an alternate
211 probability distribution $\tilde{p}(\tau)$ satisfying $\tilde{p}(\tau) \geq \epsilon p_{\mathbf{x}_{k-1}}(\tau)$ for all $\tau \in \binom{[m]}{\beta_k}$. In this case, we can
212 exploit the relationship between probabilities to reuse the proof of Theorem 3.1. Provided
213 that the probability distribution $\tilde{p}(\tau)$ is fixed between iterations we can iterate the bound
214 unlike in Theorem 3.1. An application of Corollary 3.2 with the uniform distribution is given
215 in Remark 1.

216 **Corollary 3.2.** *Suppose the system of equations $A\mathbf{x} = \mathbf{b}$ is consistent, define $\mathbf{x}^* = A^\dagger \mathbf{b}$, and
217 let $\mathbf{x}_0 \in \text{range}(A^\top)$. Suppose one runs SKM with $\tau_k \sim \tilde{p}(\tau)$ and $t_k = t(\tau_k, \mathbf{x}_{k-1})$ and that the
218 probabilities used to sample, $\tilde{p}(\tau)$, are at least a constant factor of the probabilities (3.1); that
219 is $\tilde{p}(\tau) \geq \epsilon p_{\mathbf{x}_{k-1}}(\tau)$ for all $\tau \in \binom{[m]}{\beta_k}$. Then we have that*

$$220 \quad (3.3) \quad \tilde{\mathbb{E}}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\epsilon \binom{m}{\beta_k} \beta_k \sigma_{\min}^2(A)}{\gamma_k m \sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2$$

221 where $\tilde{\mathbb{E}}_{\tau_k}$ denotes expectation taken with respect to the sampling of τ_k according to $\tilde{p}(\tau)$ and
222 conditioned on the choices of τ_j for $j < k$. Furthermore, if $\tilde{p}(\tau)$ is constant between iterations

223 (so $\beta_j = \beta$ is constant) and independent of \mathbf{x}_{k-1} , we can iterate the previous result and have

$$224 \quad (3.4) \quad \tilde{\mathbb{E}} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \prod_{j=1}^k \left(1 - \frac{\epsilon \binom{m}{\beta} \beta \sigma_{\min}^2(A)}{\gamma_j m \sum_{\pi \in \binom{[m]}{\beta}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{j-1})}\|^2} \right) \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

225 where $\tilde{\mathbb{E}}$ denotes expectation taken with respect to all samples of τ_j for $j = 1, \dots, k$.

226 *Proof.* This proof is identical to that of Theorem 3.1 but where we first replace \mathbb{E}_{τ_k} with
227 $\tilde{\mathbb{E}}_{\tau_k}$ and $p_{\mathbf{x}_{k-1}}(\tau)$ with $\tilde{p}(\tau)$. We replace the equation in (3.2) with an inequality and must
228 add an ϵ to the numerator of the subtracted term in each line from (3.2) on. The iterated
229 bound follows from recursively applying the bounds on the conditional expectations of each
230 iteration. \blacksquare

Remark 1. (Uniform probability distribution) We consider the case of the uniform distribution over samples, i.e., $\tilde{p}(\tau) = 1/\binom{m}{\beta_k}$. We note that since

$$\frac{\min_{i \in [m]} \|\mathbf{a}_i\|^2}{\binom{m}{\beta_k} \max_{i \in [m]} \|\mathbf{a}_i\|^2} \leq p_{\mathbf{x}_{k-1}}(\tau) \leq \frac{\max_{i \in [m]} \|\mathbf{a}_i\|^2}{\binom{m}{\beta_k} \min_{i \in [m]} \|\mathbf{a}_i\|^2}$$

when $\min_{i \in [m]} \|\mathbf{a}_i\|^2 > 0$, we have

$$\tilde{p}(\tau) \geq \frac{\min_{i \in [m]} \|\mathbf{a}_i\|^2}{\max_{i \in [m]} \|\mathbf{a}_i\|^2} p_{\mathbf{x}_{k-1}}(\tau).$$

231 Thus, Corollary 3.2 holds for the uniform distribution with $\epsilon = \min_{i \in [m]} \|\mathbf{a}_i\|^2 / \max_{i \in [m]} \|\mathbf{a}_i\|^2$.
232 We note that this additionally provides a convergence analysis for RK of [60] in the case that
233 the matrix A has unnormalized rows and the uniform distribution over rows is employed in
234 sampling.

235 The next remarks make simplifying assumptions on the generalized SKM algorithm and
236 our main result to provide better context for comparison with previous works.

237 **Remark 2.** (Recovery of RK guarantees) If all of the rows of A have equal norm (not
238 necessarily unit norm), then our result specializes to

$$239 \quad (3.5) \quad \mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \sigma_{\min}^2(A)}{\gamma_k \|A\|_F^2} \right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2,$$

240 and we can iteratively apply this per-iteration guarantee to give a bound on the error in ex-
241 pectation with respect to all samples,

$$242 \quad (3.6) \quad \mathbb{E} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \prod_{j=1}^k \left(1 - \frac{\beta_j \sigma_{\min}^2(A)}{\gamma_j \|A\|_F^2} \right) \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

243 Additionally, when $\beta_k = 1$, the sampling distribution (3.1) and theoretical error upper bound (3.5) \blacksquare
244 simplifies to the probability distribution and error guarantees of [60].

245 **Remark 3.** (*Improvement of MM guarantees*) Corollary 2.1 is obtained from Theorem 3.1
 246 when rows of A have unit norm and $\beta_k = \beta$. When $\beta_k = m$, then an improved convergence
 247 rate of

$$248 \quad \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{\gamma_k}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 \leq \prod_{j=1}^k \left(1 - \frac{\sigma_{\min}^2(A)}{\gamma_j}\right) \|\mathbf{x}_0 - \mathbf{x}^*\|^2,$$

249 for MM over that shown in [31],

$$250 \quad \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{4\gamma_k}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 \leq \prod_{j=1}^k \left(1 - \frac{\sigma_{\min}^2(A)}{4\gamma_j}\right) \|\mathbf{x}_0 - \mathbf{x}^*\|^2,$$

251 is obtained.

Remark 4. (*Connection to Block RK*) Note that this bound on the convergence rate of SKM additionally provides a bound on the convergence rate of a block Kaczmarz variant. This variant is distinct from the block Kaczmarz method considered in [52]. The analysis of [52] requires a pre-partitioned row paving, while the variant considered here allows the blocks to be sampled randomly and not pre-partitioned. Consider the block Kaczmarz variant which in each iteration selects a block of β_k rows of A , τ_k , and projects the previous iterate into the solution space of the entire block of β_k equations. This variant necessarily converges faster than SKM as it makes more progress in each iteration. In particular, note that $\{\mathbf{x} | A_{\tau_k} \mathbf{x} = \mathbf{b}_{\tau_k}\} \subset \{\mathbf{x} | \mathbf{a}_{t(\tau_k, \mathbf{y})}^T \mathbf{x} = b_{t(\tau_k, \mathbf{y})}\}$. Given iterate \mathbf{x}_{k-1} and sample of rows τ_k , let \mathbf{x}_k^{SKM} denote the iterate produced by SKM and \mathbf{x}_k^{BK} denote the iterate produced by this block Kaczmarz variant. Note that \mathbf{x}_k^{SKM} is the closest point to \mathbf{x}_{k-1} on the hyperplane associated to equation $t(\tau_k, \mathbf{x}_{k-1})$ so, since \mathbf{x}_k^{BK} also lies on this hyperplane, we have

$$\|\mathbf{x}_k^{SKM} - \mathbf{x}_{k-1}\|^2 \leq \|\mathbf{x}_k^{BK} - \mathbf{x}_{k-1}\|^2.$$

Now, we note that by orthogonality of the projections, we have

$$\|\mathbf{x}_k^{SKM} - \mathbf{x}_{k-1}\|^2 + \|\mathbf{x}_k^{SKM} - \mathbf{x}^*\|^2 = \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 = \|\mathbf{x}_k^{BK} - \mathbf{x}_{k-1}\|^2 + \|\mathbf{x}_k^{BK} - \mathbf{x}^*\|^2$$

so by the above inequality, we have

$$\|\mathbf{x}_k^{BK} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_k^{SKM} - \mathbf{x}^*\|^2.$$

252 A visualization of this situation is presented in Figure 3.1. Thus, the progress made by BK
 253 in any fixed iteration is at least as large as the progress made by SKM, so it must converge at
 254 least as quickly.

255 One may be assured that the contraction term in Theorem (3.1) is always strictly positive.
 256 We prove this simple fact in Proposition 3.3.

Proposition 3.3. For any matrix A defining a consistent system with $\mathbf{x}^* = A^\dagger \mathbf{b}$ and $\mathbf{x}_{j-1} \in \text{range}(A^\top)$, we have

$$257 \quad \gamma_j \geq \frac{\beta_j \binom{m}{\beta_j} \sigma_{\min}^2(A)}{m \sum_{\tau \in \binom{[m]}{\beta_j}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{j-1})}\|^2}.$$

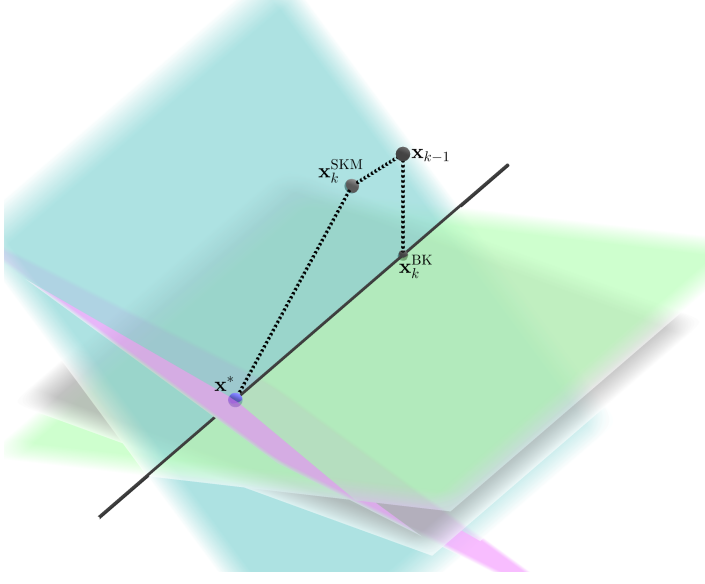


Figure 3.1: The SKM and BK iterates, $\mathbf{x}_k^{\text{SKM}}$ and \mathbf{x}_k^{BK} , generated by one iteration starting at \mathbf{x}_{k-1} satisfy $\|\mathbf{x}_k^{\text{BK}} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_k^{\text{SKM}} - \mathbf{x}^*\|^2$.

258 *Proof.* Beginning with the definition of γ_j , we have

$$\begin{aligned}
 259 \quad \frac{\sum_{\tau_j \in \binom{[m]}{\beta_j}} \|A_{\tau_j} \mathbf{x}_{j-1} - \mathbf{b}_{\tau_j}\|^2}{\sum_{\tau_j \in \binom{[m]}{\beta_j}} \|A_{\tau_j} \mathbf{x}_{j-1} - \mathbf{b}_{\tau_j}\|_\infty^2} &= \frac{\frac{\beta_j}{m} \binom{m}{\beta_j} \|A(\mathbf{x}_{j-1} - \mathbf{x}^*)\|^2}{\sum_{\tau_j \in \binom{[m]}{\beta_j}} |\mathbf{a}_{t(\tau_j, \mathbf{x}_{j-1})}^\top (\mathbf{x}_{j-1} - \mathbf{x}^*)|^2} \\
 260 \quad &\geq \frac{\frac{\beta_j}{m} \binom{m}{\beta_j} \sigma_{\min}^2(A) \|\mathbf{x}_{j-1} - \mathbf{x}^*\|^2}{\sum_{\tau_j \in \binom{[m]}{\beta_j}} \|\mathbf{a}_{t(\tau_j, \mathbf{x}_{j-1})}\|^2 \|\mathbf{x}_{j-1} - \mathbf{x}^*\|^2} \\
 261 \quad &= \frac{\beta_j \binom{m}{\beta_j} \sigma_{\min}^2(A)}{m \sum_{\tau \in \binom{[m]}{\beta_j}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{j-1})}\|^2}, \\
 262 \quad &
 \end{aligned}$$

263 where the inequality follows from properties of singular values and Cauchy-Schwartz. \blacksquare

264 Because Theorem 3.1 shows that the contraction coefficient for generalized SKM is depen-
 265 dent on the dynamic range, the following section discusses bounds on the dynamic range for
 266 special types of linear systems.

267 **4. Analysis of the Dynamic Range.** Since the dynamic range plays an integral part in the
 268 convergence behavior for generalized SKM, the dynamic range is analyzed here for different
 269 types of specialized linear systems. Note that the dynamic range has also appeared in other
 270 works, although not under the guise of “dynamic range”. For example, in [6] the authors
 271 proposed a Greedy Randomized Kaczmarz (GRK) algorithm that finds a subset of indices to
 272 randomly select the next row to project onto. The operation of finding this subset relies on a

273 ratio between the ℓ_∞ and ℓ_2 norms of the residual at the current iteration, essentially using
 274 a proxy of the dynamic range. In the next section, we analyze the dynamic range for random
 275 Gaussian linear systems and remark on the extension to other random linear systems. In the
 276 following section, we analyze the dynamic range for linear systems encoding average consensus
 277 problems on undirected graphs via the incidence matrix.

278 **4.1. Gaussian Matrices.** When entries of the measurement matrix A are drawn i.i.d. from
 279 a standard Gaussian distribution, it can be shown that the dynamic range is upper bounded
 280 by $\mathcal{O}(n\beta/\log \beta)$. The proof of the upper bound of γ_k is similar to Lemma 2 of [31], where the
 281 authors analyze the dynamic range for $\beta = m$. Here, we generalized the bound for varying
 282 samples sizes β_k .

283 **Proposition 4.1.** *Let $A \in \mathbb{R}^{m \times n}$ be a random Gaussian matrix with $a_{ij} \sim \mathcal{N}(0, \sigma^2)$. For
 284 each subset $\tau \in \binom{[m]}{\beta_k}$, let $I_\tau \subseteq \tau$ denote the set of rows in τ that are independent of \mathbf{x} and note
 285 $|I_\tau| \leq \beta_k$. Assuming there is at least m' rows in $[m]$ which are independent of \mathbf{x} , the dynamic
 286 range can be upper bounded as:*

$$287 \quad (4.1) \quad \gamma_j = \frac{\sum_{\tau \in \binom{[m]}{\beta}} \mathbb{E}_a \|A_\tau \mathbf{x}\|^2}{\sum_{\tau \in \binom{[m]}{\beta}} \mathbb{E}_a \|A_\tau \mathbf{x}\|_\infty^2} \leq \frac{\binom{m}{\beta_k} \left(\beta_k n + \sum_{i \in \tau \setminus I_\tau} \|\mathbf{a}_i\|^2 / \sigma^2 \right)}{\binom{m'}{\beta_k} \log(\beta_k)}.$$

288 **Remark 5.** *Note that the factor $\binom{m}{\beta_k} / \binom{m'}{\beta_k}$ is $\mathcal{O}(1)$ as $m \rightarrow \infty$ since*

$$289 \quad \frac{\binom{m}{\beta_k}}{\binom{m'}{\beta_k}} = \frac{m!}{\beta_k! (m - \beta_k)!} \frac{\beta_k! (m - j - \beta_k)!}{(m - j)!} = \prod_{i=0}^j \frac{m - i}{m - \beta_k - i}.$$

290 *Thus, we conclude that the expected dynamic range for any iteration k is $\mathcal{O}(n\beta_k/\log(\beta_k))$.*

291 **Proof.** Without loss of generality, we let the solution to the system $\mathbf{x}^* = 0$ so that $\mathbf{b} = 0$.
 292 We are then interested in finding an upper bound on the dynamic range (2.1) in expectation.
 293 Here, the expectation is taken with respect to the random i.i.d. draws of the entries of A .
 294 To that end, we derive upper bounds and lower bounds on the numerator and denominator

295 of (2.1). Starting with the upper bound on the numerator we have

$$\begin{aligned}
296 \quad & \sum_{\tau \in \binom{[m]}{\beta_k}} \mathbb{E}_a \|A_\tau \mathbf{x}\|^2 \leq \sum_{\tau \in \binom{[m]}{\beta}} \sum_{i \in \tau} \mathbb{E}_a (\|\mathbf{a}_i\|^2 \|\mathbf{x}\|^2) \\
297 \quad & = \sum_{\tau \in \binom{[m]}{\beta_k}} \left(\sum_{i \in I_\tau} \mathbb{E}_a \|\mathbf{a}_i\|^2 \|\mathbf{x}\|^2 + \sum_{i \in \tau \setminus I_\tau} \|\mathbf{a}_i\|^2 \|\mathbf{x}\|^2 \right) \\
298 \quad & = \sum_{\tau \in \binom{[m]}{\beta_k}} \left(\sum_{i \in I_\tau} n\sigma^2 \|\mathbf{x}\|^2 + \sum_{i \in \tau \setminus I_\tau} \|\mathbf{a}_i\|^2 \|\mathbf{x}\|^2 \right) \\
299 \quad & \leq \sum_{\tau \in \binom{[m]}{\beta_k}} \left(\beta_k n\sigma^2 + \sum_{i \in \tau \setminus I_\tau} \|\mathbf{a}_i\|^2 \right) \|\mathbf{x}\|^2 \\
300 \quad & = \binom{m}{\beta_k} \left(\beta_k n\sigma^2 + \sum_{i \in \tau \setminus I_\tau} \|\mathbf{a}_i\|^2 \right) \|\mathbf{x}\|^2 \\
301 \quad &
\end{aligned}$$

302 where the first inequality follows from the Cauchy-Schwartz inequality and remaining compu-
303 tation uses the fact that $\mathbb{E}_a \|\mathbf{a}_i\|^2 = n\sigma^2$ and simplifies the expression. The lower bound follows
304 from

$$\begin{aligned}
305 \quad & \sum_{\tau \in \binom{[m]}{\beta}} \mathbb{E}_a \|A_\tau \mathbf{x}\|_\infty^2 = \sum_{\tau \in \binom{[m]}{\beta}} \mathbb{E}_a \max_{i \in \tau} \langle \mathbf{a}_i, \mathbf{x} \rangle^2 \geq \sum_{\tau \in \binom{[m]}{\beta}} \mathbb{E}_a \max_{i \in I_\tau} \langle \mathbf{a}_i, \mathbf{x} \rangle^2 \\
306 \quad & \geq \sum_{\tau \in \binom{[m]}{\beta}} \left(\mathbb{E}_a \max_{i \in I_\tau} \langle \mathbf{a}_i, \mathbf{x} \rangle \right)^2 \geq \sum_{\substack{\tau \in \binom{[m]}{\beta} \\ |I_\tau| = \beta_k}} \left(\mathbb{E}_a \max_{i \in I_\tau} \langle \mathbf{a}_i, \mathbf{x} \rangle \right)^2 \\
307 \quad & \geq \sum_{\substack{\tau \in \binom{[m]}{\beta} \\ |I_\tau| = \beta_k}} \sigma^2 \|\mathbf{x}\|^2 \log(\beta_k) \\
308 \quad & \geq \binom{m'}{\beta_k} \sigma^2 \|\mathbf{x}\|^2 \log(\beta_k), \\
309 \quad &
\end{aligned}$$

310 where the second to last inequality uses the fact that for i.i.d. Gaussian random vari-
311 ables $g_1, g_2, \dots, g_N \sim \mathcal{N}(0, \sigma^2)$, we have that $\mathbb{E}(\max_{i \in [N]} g_i) \gtrsim \sigma \sqrt{\log N}$ and that $\langle \mathbf{a}_i, \mathbf{x} \rangle \sim$
312 $\mathcal{N}(0, \sigma^2 \|\mathbf{x}\|^2)$. Therefore, we have

$$\begin{aligned}
313 \quad & \gamma_k = \frac{\sum_{\tau \in \binom{[m]}{\beta}} \mathbb{E}_a \|A_\tau \mathbf{x}\|^2}{\sum_{\tau \in \binom{[m]}{\beta}} \mathbb{E}_a \|A_\tau \mathbf{x}\|_\infty^2} \\
314 \quad & \leq \frac{\binom{m}{\beta_k} \left(\beta_k n\sigma^2 + \sum_{i \in \tau \setminus I_\tau} \|\mathbf{a}_i\|^2 \right)}{\binom{m'}{\beta_k} \sigma^2 \log(\beta_k)}. \\
315 \quad &
\end{aligned}$$

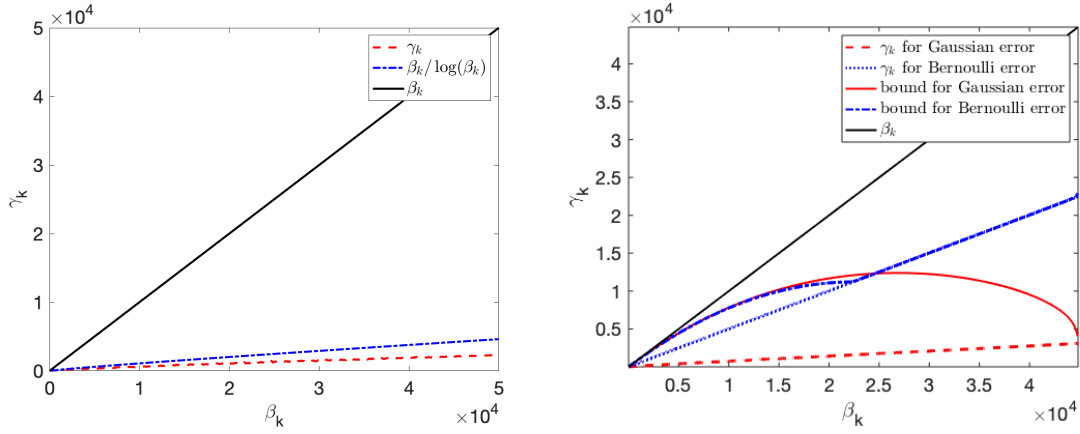


Figure 4.1: Dynamic ranges γ_k for various sample sizes β_k . Left: Gaussian matrix $A \in \mathbb{R}^{50000 \times 500}$ and Gaussian error $\mathbf{e}_k \in \mathbb{R}^{500}$ (red). Conjectured bound is plotted in blue. Right: Incidence matrix $Q \in \mathbb{R}^{44850 \times 300}$ for complete graph K_{300} with Gaussian error $\mathbf{e}_k \in \mathbb{R}^{300}$ (red) and a sparse error \mathbf{e}_k with Bernoulli random variable entries (blue). Bounds are plotted in the same colors with different line styles.

316 Dividing all terms by σ^2 attains the desired result (4.1). ■

317 We conjecture that the true bound is actually $\mathcal{O}(\beta_k / \log(\beta_k))$ and that the n is an artifact
 318 of our proof technique; throughout our experiments varying n (and for various m), we have not
 319 found any dependence of γ_k on n . For this reason, we have plotted γ_k and the corresponding
 320 conjectured bound in the left of Figure 4.1 for a Gaussian matrix of size 50000×500 .

321 **Remark 6.** To extend to other distributions, one can simply note that as the signal dimen-
 322 sion n gets large, the Law of Large numbers can be invoked and a similar computation can be
 323 used to show an upper bound on the dynamic range of the system.

With this estimate for γ_k , one can now ask for the optimal choice for β_k in terms of computational time. We can use Theorem 3.1 to estimate the expected number of SKM iterations required so that $\|\mathbf{x}_k - \mathbf{x}^*\|^2 / \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq \epsilon$ for $\epsilon < 1$; assuming that A is row normalized, we expect this relative error stopping threshold is reached for number of iterations

$$k \geq \frac{\log(\epsilon)}{\log\left(1 - \frac{\beta_k \sigma^2 \min(A)}{\gamma_k m}\right)}.$$

To estimate the optimal parameter β_k , we must have an accurate estimate for computational effort of each SKM iteration with sample size β . We first use an estimate for FLOPS required in each iteration of SKM as a proxy for computational time (which is notoriously difficult to estimate); the per-iteration flop requirement for each iteration of SKM is $\mathcal{O}(n\beta_k + n)$.

Therefore, the required FLOPS to reach relative error threshold ϵ is

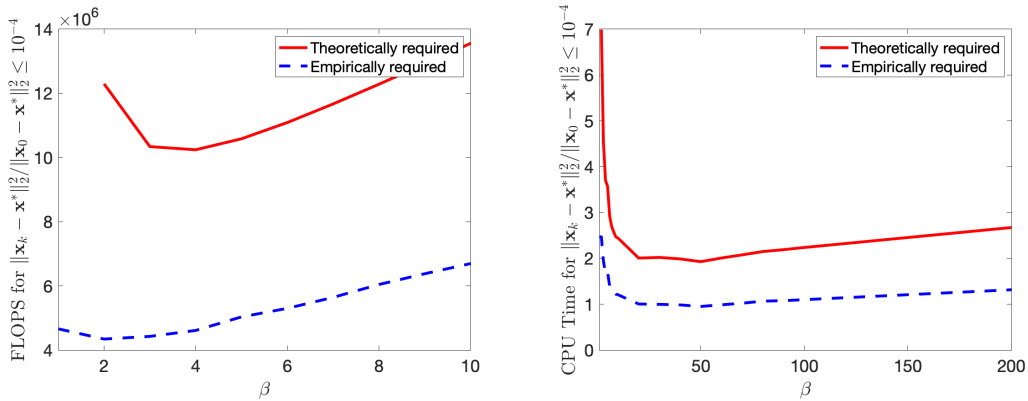
$$\frac{(n\beta_k + n) \log(\epsilon)}{\log\left(1 - \frac{\beta_k \sigma_{\min}^2(A)}{\gamma_k m}\right)} \approx \frac{(n\beta_k + n) \log(\epsilon)}{\log\left(1 - \frac{\log(\beta_k) \sigma_{\min}^2(A)}{nm}\right)},$$

324 where we have used the estimate $\gamma_k \approx n\beta_k/\log(\beta_k)$. This estimate is plotted for a variety of
 325 sample sizes β_k for a row-normalized Gaussian matrix $A \in \mathbb{R}^{50000 \times 500}$ with $\epsilon = 10^{-4}$ (solid red
 326 line) in Figure 4.2a ; we additionally plot $(n\beta_k + n)k_{\text{emp}}$ where k_{emp} is the number of iterations
 327 empirically required to reach relative error stopping threshold $\epsilon = 10^{-4}$ (dashed blue line).

We next estimate computational time required per iteration as the average CPU time for
 a single iteration of an empirical trial of SKM, t_{β_k} . We estimate the required CPU time to
 reach relative error threshold ϵ as

$$\frac{t_{\beta_k} \log(\epsilon)}{\log\left(1 - \frac{\beta_k \sigma_{\min}^2(A)}{\gamma_k m}\right)} \approx \frac{t_{\beta_k} \log(\epsilon)}{\log\left(1 - \frac{\log(\beta_k) \sigma_{\min}^2(A)}{nm}\right)}.$$

328 This estimate is plotted for a variety of sample sizes β_k for a row-normalized Gaussian matrix
 329 $A \in \mathbb{R}^{50000 \times 500}$ with $\epsilon = 10^{-4}$ (solid red line) in Figure 4.2b; we additionally plot the total
 330 CPU time empirically required to reach relative error stopping threshold $\epsilon = 10^{-4}$ (dashed
 331 blue line). We note that the estimated optimal choice for β_k differs given the estimates of
 332 per-iteration computational burden. This is to be expected as our estimate for FLOPS per
 333 SKM iteration ignores computational overhead and communication time that contribute to
 334 the CPU time.



(a) Estimate of theoretically required and empirically required FLOPS.

(b) Estimate of theoretically required and empirically required CPU time.

Figure 4.2: Comparison of theoretically and empirically required FLOPS and CPU time to reach given relative error stopping threshold and estimate of optimal β_k for Gaussian system defined by row-normalized $A \in \mathbb{R}^{50000 \times 500}$, averaged over five independent trials.

335 This estimate of the optimal choice for β_k for a given system was dependent upon having an
 336 estimate for γ_k that is fixed between iterations and does not depend upon the current iterate.
 337 Gaussian systems are the only type where we have such an estimate; in the next section,
 338 we provide a bound for systems defined by incidence matrices, but this bound is iterate-
 339 dependent. When there is no iteration-consistent estimate for γ_k , one can attempt to choose
 340 the optimal β_k using an estimate for γ_k based on the sampled portion of the residual. We
 341 take this naive, empirical approach in Section 5 and call this selection strategy ‘useDynRng’.

342 **4.2. Incidence Matrices.** In the previous subsection, we analyzed the dynamic range for
 343 systems with measurement matrices that are randomly generated. Deterministically generated
 344 measurement matrices are additionally of interest. In this subsection, we analyze the dynamic
 345 range associated to incidence matrices of undirected graphs, $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. The incidence matrix
 346 Q associated to an undirected graph is of size $|\mathcal{E}| \times |\mathcal{V}|$. For each edge, $(i, j) \in \mathcal{E}$ which connects
 347 vertex i to vertex j , the associated row of Q is all zeros with a one and negative one in the
 348 i th and j th entries. These types of matrices arise in one formulation of the *average consensus*
 349 problem as a system of linear equations.

350 The average consensus problem on a graph asks that all nodes on the graph learn the aver-
 351 age value of initial, secret values held by each node using only local information; that is, each
 352 node i initially knows c_i and at solution they should all know $\frac{1}{|\mathcal{V}|} \sum_{i \in \mathcal{V}} c_i$ with communication
 353 only across edges. This problem models computation in many real life applications such as
 354 clock synchronization [23], localization without GPS [67], distributed data fusion in sensor
 355 networks [66], and load balancing [17]. Many analyses of (asynchronous and synchronous)
 356 distributed methods for this problem exploit its formulation as a system of linear equations.
 357 The problem over a directed graph may be formulated as a linear system using either the
 358 incidence matrix (described above) or the Laplacian matrix, $L = D - A$ where D is the diag-
 359 onal matrix of node degrees and A is the adjacency matrix, or more generally as an *average*
 360 *consensus system* defined in [39].

361 The *gossip methods* that solve the average consensus problem are generalized by the Kacz-
 362 marz methods [43]. Early work making this connection focused on the formulation of the
 363 average consensus problem as a Laplacian system [69], but subsequent work generalized this
 364 connection to systems formulated more generally [39]. RK specializes to the *randomized gossip*
 365 method in which the pair of nodes which update are selected at random [65, 11]. The connec-
 366 tion between gossip methods and other Kaczmarz variants have been observed; the connection
 367 to block methods was noted in [39], extended methods in [69], and accelerated methods in
 368 [41, 38]. This connection has also spawned new gossip methods; in [34] the authors propose a
 369 privacy preserving gossip method, and in [5], the authors propose an accelerated, decentral-
 370 ized gossip method. In [43], the authors summarize many of these advances and connections
 371 between the Kaczmarz literature and gossip literature. They first noted and exploited the
 372 fact that SKM specializes to a variant of *greedy gossip with eavesdropping (GGE)* in which
 373 the nodes are selected from amongst a random sample to maximize the update [63]. By uti-
 374 lizing the connection noted between GGE and SKM in [43] and with some adjustments to
 375 the method and theoretical setup here (namely sampling according to the connectivity of the
 376 network and redefining the dynamic range accordingly), the approach for proving our main
 377 convergence rate can be employed in proving a similar convergence rate for GGE.

378 Now, we consider the dynamic range for an incidence matrix. We can derive a simple
 379 bound on the dynamic range in each iteration that depends only upon the entries of the
 380 current error vector, $\mathbf{e}_k := \mathbf{x}_k - \mathbf{x}^*$. In particular,

$$\begin{aligned}
 381 \quad (4.2) \quad \gamma_k &= \frac{\sum_{\tau \in \binom{[m]}{\beta_k}} \|Q_\tau \mathbf{e}_k\|^2}{\left\| \sum_{\tau \in \binom{[m]}{\beta_k}} Q_\tau \mathbf{e}_k \right\|_\infty^2} = \frac{\binom{m}{\beta_k} \frac{\beta_k}{m} \sum_{(i,j) \in \mathcal{E}} (e_k^{(i)} - e_k^{(j)})^2}{\sum_{\tau \in \binom{[m]}{\beta_k}} \max_{(i,j) \in \tau} (e_k^{(i)} - e_k^{(j)})^2} \\
 382 \quad &\leq \frac{\beta_k(m - \beta_k + 1) \sum_{(i,j) \in \mathcal{E}} (e_k^{(i)} - e_k^{(j)})^2}{m \sum_{n=\beta_k}^{|\mathcal{E}|} (e_k^{(n_i)} - e_k^{(n_j)})^2}, \\
 383 \quad & \\
 384 \quad &
 \end{aligned}$$

385 where n_i and n_j denote the vertices connected by the n th smallest magnitude difference
 386 across an edge. This bound improves for iterates with a sufficient amount of variation in the
 387 coordinates. We have plotted γ_k and the corresponding bounds in the right of Figure 4.1. We
 388 calculate these values for the incidence matrix $Q \in \mathbb{R}^{44850 \times 300}$ of the complete graph K_{300} in
 389 the cases when the error is a Gaussian vector (red) and a Bernoulli vector (blue).

Proposition 4.2. *If the right-hand-side vector associated to the system $Q\mathbf{x} = \mathbf{b}$ is $\mathbf{b} = \mathbf{0}$, as in the average consensus problem, then this bound on the dynamic range is easily computed from the current iterate,*

$$\gamma_k \leq \frac{\beta_k(m - \beta_k + 1) \sum_{(i,j) \in \mathcal{E}} (x_k^{(i)} - x_k^{(j)})^2}{m \sum_{n=\beta_k}^{|\mathcal{E}|} (x_k^{(n_i)} - x_k^{(n_j)})^2}.$$

390 **Remark 7.** *We note that this bound on the dynamic range holds for any incidence matrix*
 391 *Q , including those associated with directed graphs. In the case of directed graphs, however,*
 392 *additional assumptions must be made to ensure the well-posedness of the average consensus*
 393 *problem. Additionally, the Kaczmarz methods must be altered to ensure communication in*
 394 *only one direction along edges for directed graphs. Analyses of regular Kaczmarz methods,*
 395 *such as RK or SKM, do not apply to average consensus systems on directed graphs. We leave*
 396 *consideration of Kaczmarz type methods for this variant of the average consensus problem to*
 397 *future work.*

398 **5. Experiments.** In this section, we present simulated and real world experiments using
 399 SKM for varying sample sizes β . In the simulated experiments, we compare the theoretical
 400 convergence guarantees to the empirical performance of SKM, measured by approximation
 401 error $\|\mathbf{e}_j\|^2$, averaged over 20 random trials. The number of rows $m = 50000$ and number
 402 of columns $n = 500$ are fixed for all simulated experiments. The solution to the system is a
 403 vector $\mathbf{x}^* \in \mathbb{R}^n$ where each entry is drawn i.i.d. from a standard Gaussian distribution. In
 404 each experiment, the systems are consistent so that $\mathbf{b} = A\mathbf{x}^*$. The sample sizes considered for
 405 this experiment are $\beta = \{1, 100, 200, 500, 1000\}$. Unless otherwise stated, the rows of A are
 406 uniformly selected without replacement. The experiments presented in this section were
 407 performed in MATLAB 2017b on a MacBook Pro 2015 with a 2.7 GHz Dual-Core Intel Core
 408 i5 and 8GB RAM. For practical reasons, we normalize the rows of A and utilize the bound
 409 shown in Corollary 2.1.

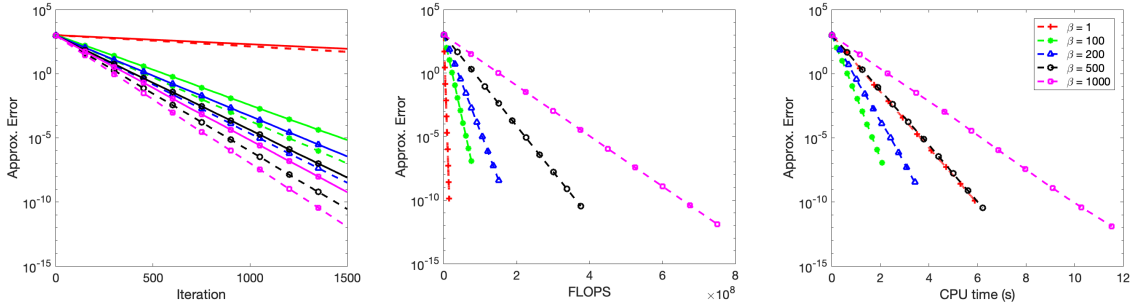


Figure 5.1: Comparison of SKM for various choices of fixed β values on linear system with entries of A drawn from i.i.d. from $\mathcal{N}(0,1)$. (left) Iteration vs Approximation Error with dashed lines representing average empirical performance of SKM and solid lines representing theoretical upper bounds for SKM. (middle) FLOPS vs Approximation Error. (right) CPU time vs Approximation Error

410 Figure 5.1 and Figure 5.2 show the results for Gaussian and Uniform random matrices A
 411 respectively. For Gaussian random matrices, each entry of A is drawn i.i.d. from a standard
 412 Gaussian distribution. For Uniform random matrices, entries of A are drawn i.i.d. uniformly
 413 from the interval $[0, 1]$. In each figure, we plot along the horizontal axis the (left subplot)
 414 iteration, (middle subplot) FLOPS or floating point operations, and (right subplot) CPU
 415 time in seconds. The vertical axis for all plots indicate the average approximation error across
 416 random trials. Note that the left most subplot for both figures also contains a solid line, which
 417 indicates the theoretical upper bound of the algorithm provided by Corollary 2.1.

418 For linear systems with Gaussian random matrices, we see in Figure 5.1 that the conver-
 419 gence upper bound proven in this work closely matches the behavior of the SKM algorithm
 420 regardless of the choice of sample size β . To compare this result to previous works, note that
 421 when $\beta = 1$, the upper bound provided in Corollary 2.1 simply recovers the previous known
 422 upper bound for SKM with normalized rows, a bound which was completely independent of
 423 β . In other words, the solid red line is the comparative previous known SKM upper bound
 424 for all β .

425 Of course, choices of large sample sizes β come at a cost, which are captured in the middle
 426 and right most subplots of Figure 5.1. When measuring efficiency, it seems that $\beta = 1$ makes
 427 the most progress with minimal FLOPS while $\beta = 100$ is optimal amongst the tested sample
 428 sizes with respect to CPU time. This difference is typically explained by the programming and
 429 computer architecture (e.g., it may be more efficient to work on batches of rows as opposed
 430 to single rows at a time).

431 Figure 5.2 uses a uniform random matrix A instead of a Gaussian random matrix. While
 432 the algorithm efficiency with respect to FLOPS and CPU time have similar conclusions to
 433 those in the Gaussian measurement matrix case (as one would expect), the iteration vs ap-
 434 proximation error plot now tells a different story. Unlike in the Gaussian case, the theoretical
 435 upper bound no longer closely tracks the approximation error of SKM. The looseness here
 436 comes from lower bounding the norm of $\|A\mathbf{x}\|_2^2$ with the magnitude of \mathbf{x} times the smallest

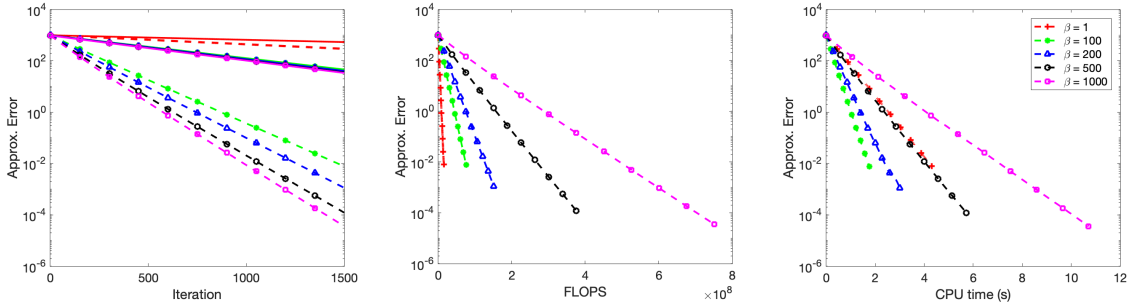


Figure 5.2: Comparison of SKM for various choices of fixed β values on linear system with entries of a A drawn from $\text{unif}([0, 1])$. (left) Iteration vs Approximation Error with dashed lines representing average empirical performance of SKM and solid lines representing theoretical upper bounds for SKM. (middle) FLOPS vs Approximation Error. (right) CPU time vs Approximation Error

437 nonzero singular value of A squared. Empirically, we have seen that this lower bound is
 438 tighter for Gaussian systems than Uniform systems. It should be noted that even though our
 439 theoretical bounds do not track the approximation error for SKM as tightly, they are still a
 440 slight improvement over the previous known bounds for SKM.

441 In addition to being an improvement over the previously known SKM bound, the con-
 442 vergence bound shown in this work enjoys the flexibility of being amenable to a dynamically
 443 selected sample size β_k . Figure 5.3 and Figure 5.4 show the empirical results from experiments
 444 where β_k is allowed to change at every iteration. In Figure 5.3 the measurement matrix A
 445 is again a random Gaussian matrix and in Figure 5.4 the measurement matrix entries are
 446 drawn i.i.d. from $\text{Unif}([0, 1])$. We consider three sampling regimes that change β_k at every
 447 iteration: ‘useDynRng’ which allocates β_k as a function of the dynamic range, ‘slowInc’ which
 448 increases β_k at every iteration until $\beta_k = m$, and finally ‘rand’ which uniformly at random
 449 selects a $\beta_k \in [m]$ at every iteration. More specifically, the ‘useDynRng’ uses the heuristic
 450 $\beta_k = \lceil \max(m, \frac{m \|A_{\tau_{k-1}} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_{k-1}}\|_{\infty}}{n \|A_{\tau_{k-1}} \mathbf{x}_k - \mathbf{b}_{\tau_{k-1}}\|_2}) \rceil$. Note that this choice of β_k relies directly on the inverse
 451 of an *approximation* of the dynamic range γ_k computed without incurring additional compu-
 452 tational cost, in a naive attempt to optimize the contraction term of the theoretical bound
 453 for SKM. Even though β_k changes at each iteration, we see that the theoretical guarantees
 454 proven in this work still track the progress of SKM. This indeed opens up new and interesting
 455 avenues of research including how one can compute an optimal β_k at every iteration. Since
 456 the focus of this work is the improvement of the convergence bound of SKM, we leave this for
 457 future work.

458 Figure 5.5 employs the upper bound on the dynamic range derived in Proposition 3.3 to
 459 approximate an upper bound for the error of SKM iterates when $\beta = 100$. Here, we compare
 460 the empirical performance of SKM with its previous known upper bound using the contraction
 461 term $1 - \frac{\sigma_{\min}^2}{m}$ and $1 - \frac{\log(\beta)\sigma_{\min}^2}{m}$. Note that we drop the factor of n apparent in Proposition 3.3
 462 as we suspect it to be an artifact of the proof technique used and conjecture that the true

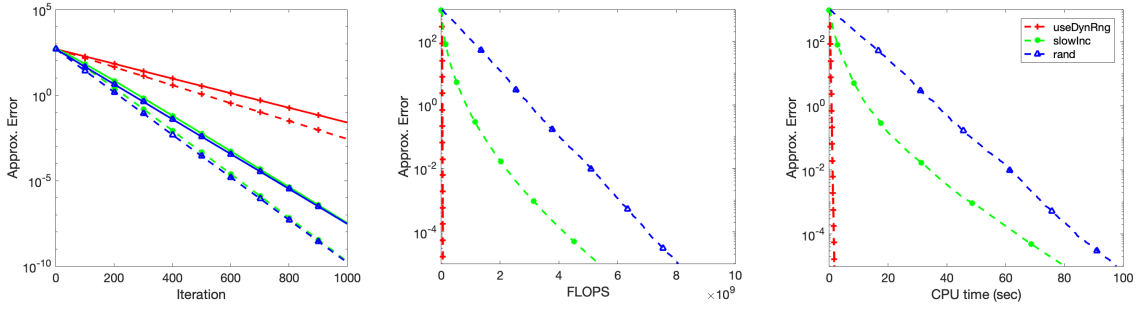


Figure 5.3: Comparison of SKM for various choices of dynamically selected β_k values on linear system with entries of A drawn from $\mathcal{N}(0, 1)$.

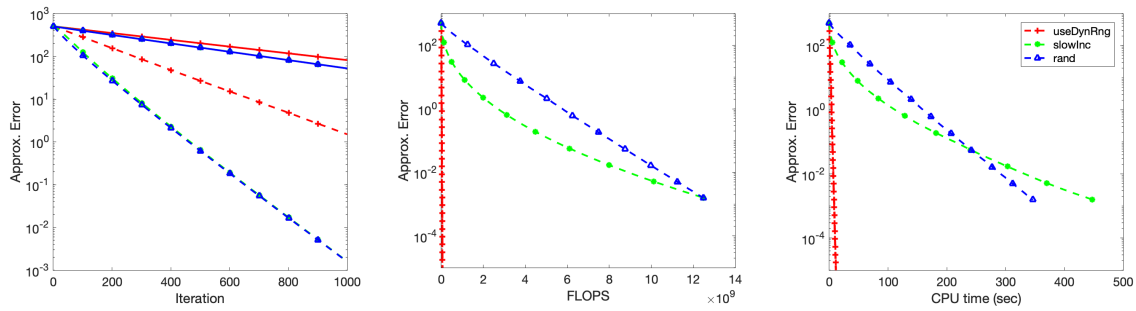


Figure 5.4: Comparison of SKM for various choices of dynamically selected β_k values on linear system with entries of A drawn from $\text{Unif}([0, 1])$.

463 upper bound of the dynamic range is actually $\mathcal{O}(\beta/\log(\beta))$.

464 In both the Gaussian and Uniform synthetic experiments, the row norms of A are of simi-
 465 lar magnitude on average and thus choosing β_k rows of the measurement matrix uniformly at
 466 random will behave similarly to the theoretically imposed probability distribution introduced
 467 in (3.1). In the next experiment, we consider a setting where the entries of the measurement
 468 matrix are $a_{ij} \sim \mathcal{N}(0, i/\sqrt{n})$ so that for each row, $\mathbb{E}\|\mathbf{a}_i\|^2 = i$. Since (3.1) is computationally
 469 impractical to implement, we will continue to select rows of A uniformly at random without
 470 replacement to evaluate the performance of SKM for various choices of β . The results of this
 471 experiment are presented in Figure 5.6. As in the previous synthetically generated experi-
 472 ments, $m = 50000$, $n = 500$, and the underlying signal x is a standard Gaussian random
 473 vector. Despite not sampling rows as imposed by (3.1), we see that SKM still converges with
 474 rates similar to those in Figure 5.1 and $\beta = 100$ outperforms the others with respect to CPU
 475 time.

476 Next we move on to evaluate the performance of SKM on incidence matrices of graphs. We
 477 start with the the AC systems discussed in Section 4.2. Here, the graph \mathcal{G} is a complete graph
 478 K_{100} with corresponding incidence matrix $Q \in \{-1, 0, 1\}^{4950 \times 100}$. The unknown underlying
 479 vector $\mathbf{x} \in \mathbb{R}^{100}$ is $\mathbf{x} = \hat{\mu}\mathbf{1}_{100}$ where $\mathbf{1}_{100}$ is a 100-dimensional vector of ones and $\hat{\mu}$ is the

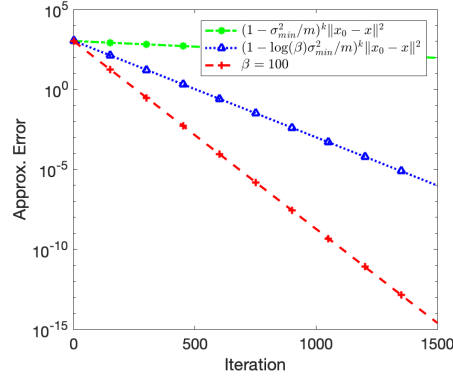


Figure 5.5: Comparison of SKM with $\beta = 100$ with previous known theoretical upper bounds and our upper bound with conjectured Gaussian system dynamic range γ_k .

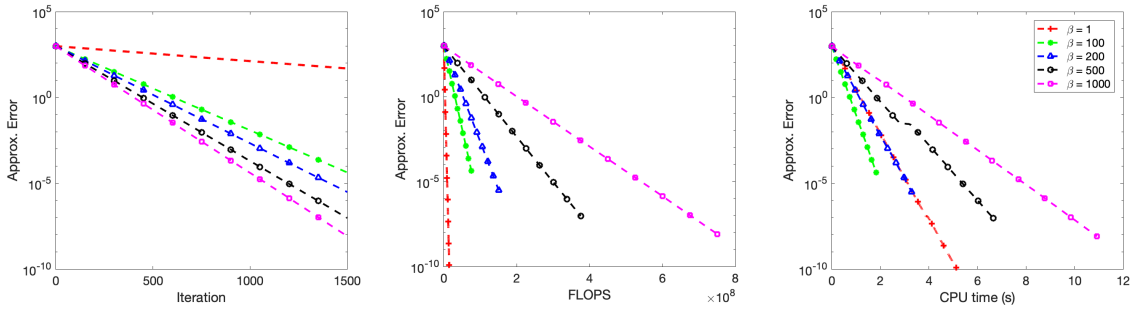


Figure 5.6: Comparison of SKM for various choices of fixed β values on linear system with entries a_{ij} drawn from $\mathcal{N}(0, i/\sqrt{n})$. (left) Iteration vs Approximation Error with dashed lines representing average empirical performance of SKM and solid lines representing theoretical upper bounds for SKM. (middle) FLOPS vs Approximation Error. (right) CPU time vs Approximation Error

480 empirical average of 100 random draws from a standard normal distribution. The results of
 481 this experiment are provided in Figure 5.7.

482 Figure 5.8 demonstrates the performance of SKM on a graph \mathcal{G} which reflects a scale-free
 483 network, i.e., a graph whose degree distribution follows the power law. To create the graph, we
 484 employ the implementation of the Barabási-Albert (BA) model [3] with an initial graph of five
 485 vertices and ending with a graph of 300 vertices [24]. For more details on scale-free networks,
 486 see [3]. In Figure 5.8 we again observe exponential convergence in the mean approximation
 487 error and optimal performance with respect to CPU time when $\beta = 10$.

488 In Figures 5.9a and 5.9b we present the performance of SKM for GGE problems (See Sec-
 489 tion 4.2). In such problems, instead of randomly selecting a subset of β rows of the incidence
 490 matrix of a graph uniformly, we randomly select a node (column), collect all rows correspond-
 491 ing to neighbors of said node (rows corresponding to nonzero entries in that column), and

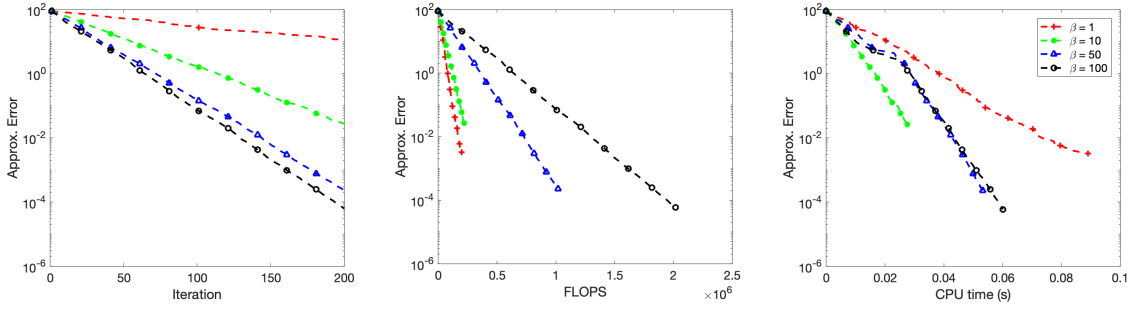


Figure 5.7: Comparison of SKM for various choices of fixed β on AC systems. (left) Iteration vs Approximation Error with dashed lines representing average empirical performance of SKM. (middle) FLOPS vs Approximation Error. (right) CPU time vs Approximation Error

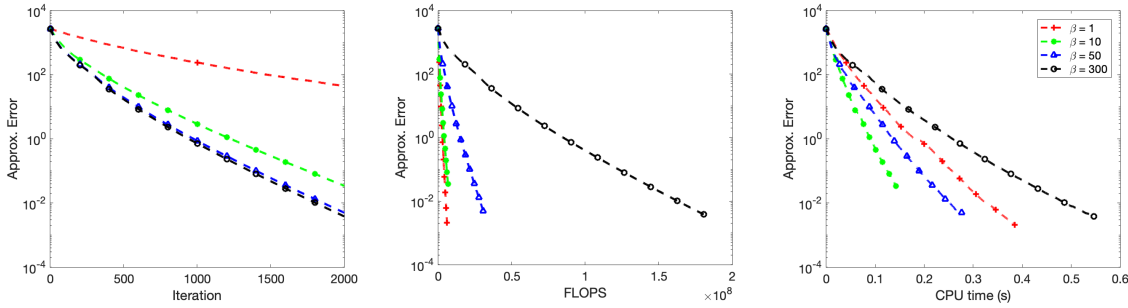
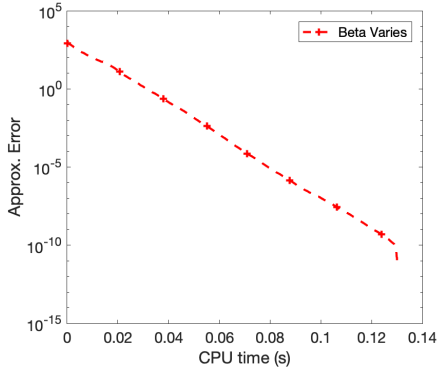


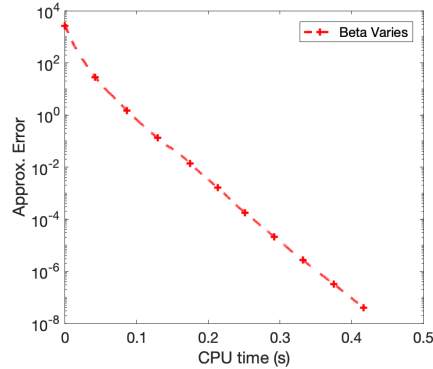
Figure 5.8: Comparison of SKM for various choices of fixed β on an incidence matrix of a scale-free network. (left) Iteration vs Approximation Error with dashed lines representing average empirical performance of SKM. (middle) FLOPS vs Approximation Error. (right) CPU time vs Approximation Error

492 then of those rows, select the row that returns the max absolute residual value to project
 493 on. Note that here, the size of the subset varies at each iteration, depending on the size of
 494 the neighborhood of the randomly selected node. Since this is significantly different from the
 495 standard SKM selection scheme, we refer to this method as GGE-SKM. Figure 5.9a presents
 496 the computational results for an incidence matrix of a complete graph (system set up as in
 497 Figure 5.7) while Figure 5.9b does the same for an incidence matrix of a BA model graph
 498 (system set up as in Figure 5.8).

499 Our last experiment demonstrates the performance of SKM with varying choices of β
 500 on real world data from the SuiteSparse Matrix Market [18]. For these experiments, we
 501 employ the ‘Maragal_4’, ‘well1850’, and ‘ash958’ matrices which are of dimension 1964×1034 ,
 502 1850×712 , and 958×292 respectively. It is useful to note that ‘Maragal_4’ and ‘well1850’
 503 have low rank data structures. These matrices are used as the measurement matrix A and
 504 the underlying signal is arbitrarily chosen to be in the range of A^T . Here, we allow SKM to
 505 run for a maximum of 10^6 iterations or terminate once the approximation error has reached



(a) Performance of CPU time vs Approximation Error of GGE-SKM on an incidence matrix of a complete graph.



(b) Performance of CPU time vs Approximation Error of GGE-SKM on an incidence matrix of a scale-free network.

Figure 5.9: Performance of GGE-SKM on GGE problem with two different graphs.

	SKM $\beta = 1$			SKM $\beta = 10$		
	Error	FLOPS	CPU time	Error	FLOPS	CPU time
Maragal_4	2.45	2.07e+09	176.25	0.164	1.13e+10	460.88
well1850	7.67	1.42e+09	112.75	0.064	7.83e+09	238.59
ash958	1.06e-06	4.91e+06	0.323	1.00e-06	5.33e+06	0.126

Table 5.1: Performance of SKM using $\beta = 1$ and $\beta = 10$ on real world data from SuiteSparse Matrix Market.

506 the allotted error tolerance of 10^{-6} . The results of using $\beta \in \{1, 10, 50\}$ are presented in
 507 Table 5.1 and Table 5.2. We also included the performance of Conjugate Gradient Least
 508 Squares (GCLS) [10, 61] for comparison. Note that we do not claim to have optimized either
 509 implementation.

510 Table 5.1 and Table 5.2 demonstrates that for a fixed number of iterations, increasing β
 511 results in a lower approximation error. It also highlights the trade off between the subset size,
 512 the FLOP cost, and the CPU time. As we increase β , in general, both the FLOP and CPU
 513 time increase as well. One interesting observation is that for the ‘ash958’ experiment, the
 514 optimal choice of subset size is $\beta = 10$ with respect to CPU time. This is further motivation
 515 for future work in optimal β selection. Finally, as expected, CGLS outperforms the three
 516 choices of β . However, SKM can be more naturally implemented in distributed computing
 517 settings. We leave that direction as an avenue for future work as well.

518 **6. Conclusion.** This work unifies the spectrum between the randomized Kaczmarz and a
 519 greedy variant of the Kaczmarz (Motzkin’s Method) algorithm by improving the convergence
 520 bound of SKM, a hybrid randomized-greedy algorithm. We show that the behavior of SKM

	SKM $\beta = 50$			CGLS		
	Error	FLOPS	CPU time	Error	FLOPS	CPU time
Maragal4	5.01e-02	5.27e+10	1593.4	3.97e-02	2.97e+07	5.41
well1850	2.49e-03	3.63e+10	1161	1.01e-06	4.06e+06	1.07
ash958	1.01e-06	1.44e+07	0.426	1.88e-06	41158	6.66e-3

Table 5.2: Performance of SKM using $\beta = 50$ and CGLS on real world data from SuiteSparse Matrix Market.

521 depends on the sample parameter β_k and the dynamic range of the linear system. This result
 522 improves upon previous work showing only the linear convergence of SKM. In presenting an
 523 improved convergence bound for SKM that highlights the impact of the sub-sample size β_k , we
 524 have opened up new and exciting avenues for SKM-type algorithms. Future directions of this
 525 work include finding optimal sample sizes for different types of linear systems and designing
 526 adaptive sample size selection schemes.

527 **Acknowledgements.** The authors would like to thank the manuscript referees for their
 528 thoughtful and detailed comments which significantly improved earlier versions of this work.
 529 The authors also thank Jacob Moorman, Liza Rebrova, Hanbaek Lyu, Deanna Needell, Jesús
 530 A. De Loera, and Roman Vershynin for useful conversations and suggestions.

531

REFERENCES

- 532 [1] S. AGMON, *The relaxation method for linear inequalities*, Canadian J. Math., 6 (1954), pp. 382–392.
 533 [2] R. AHARONI AND Y. CENSOR, *Block-iterative projection methods for parallel computation of solutions to*
 534 *convex feasibility problems*, Linear Algebra Appl., 120 (1989), pp. 165–175, [https://doi.org/10.1016/](https://doi.org/10.1016/0024-3795(89)90375-3)
 535 [http://dx.doi.org/10.1016/0024-3795\(89\)90375-3](http://dx.doi.org/10.1016/0024-3795(89)90375-3).
 536 [3] R. ALBERT AND A.-L. BARABÁSI, *Statistical mechanics of complex networks*, Rev. Mod. Phys., 74 (2002),
 537 p. 47.
 538 [4] E. AMALDI AND R. HAUSER, *Boundedness theorems for the relaxation method*, Math. Oper. Res., 30
 539 (2005), pp. 939–955, <https://doi.org/10.1287/moor.1050.0164>, [http://dx.doi.org/10.1287/moor.1050.](http://dx.doi.org/10.1287/moor.1050.0164)
 540 [0164](http://dx.doi.org/10.1287/moor.1050.0164).
 541 [5] N. S. AYBAT AND M. GÜRBÜZBALABAN, *Decentralized computation of effective resistances and accelera-*
 542 *tion of consensus algorithms*, in IEEE Glob. Conf. Sig., IEEE, 2017, pp. 538–542.
 543 [6] Z.-Z. BAI AND W.-T. WU, *On greedy randomized Kaczmarz method for solving large sparse linear systems*,
 544 SIAM J. Sci. Comput., 40 (2018), pp. A592–A606.
 545 [7] Z.-Z. BAI AND W.-T. WU, *On relaxed greedy randomized Kaczmarz methods for solving large sparse*
 546 *linear systems*, Appl. Math. Lett., 83 (2018), pp. 21–26.
 547 [8] U. BETKE, *Relaxation, new combinatorial and polynomial algorithms for the linear feasibility problem*,
 548 Discrete Comput. Geom., 32 (2004), pp. 317–338, <https://doi.org/10.1007/s00454-004-2878-4>, [http:](http://dx.doi.org/10.1007/s00454-004-2878-4)
 549 [//dx.doi.org/10.1007/s00454-004-2878-4](http://dx.doi.org/10.1007/s00454-004-2878-4).
 550 [9] U. BETKE AND P. GRITZMANN, *Projection algorithms for linear programming*, Eur. J. Oper. Res.,
 551 60 (1992), pp. 287 – 295, [https://doi.org/http://dx.doi.org/10.1016/0377-2217\(92\)90080-S](https://doi.org/http://dx.doi.org/10.1016/0377-2217(92)90080-S), [http:](http://www.sciencedirect.com/science/article/pii/037722179290080S)
 552 [//www.sciencedirect.com/science/article/pii/037722179290080S](http://www.sciencedirect.com/science/article/pii/037722179290080S).
 553 [10] Å. BJÖRCK, *Numerical methods for least squares problems*, SIAM, 1996.
 554 [11] S. BOYD, A. GHOSH, B. PRABHAKAR, AND D. SHAH, *Randomized gossip algorithms*, IEEE ACM T.
 555 Network., 14 (2006), pp. 2508–2530.
 556 [12] S. BOYD AND L. VANDENBERGHE, *Convex optimization*, Cambridge university press, 2004.

- 557 [13] Y. CENSOR, *Row-action methods for huge and sparse systems and their applications*, SIAM Rev., 23
558 (1981), pp. 444–466, <https://doi.org/10.1137/1023097>, <http://dx.doi.org/10.1137/1023097>.
- 559 [14] Y. CENSOR, P. P. B. EGGERMONT, AND D. GORDON, *Strong underrelaxation in Kaczmarz’s method for*
560 *inconsistent systems*, Numer. Math., 41 (1983), pp. 83–92.
- 561 [15] X. CHEN AND A. POWELL, *Almost sure convergence of the Kaczmarz algorithm with random mea-*
562 *surements*, J. Fourier Anal. Appl., (2012), pp. 1–20, <http://dx.doi.org/10.1007/s00041-012-9237-2>.
563 10.1007/s00041-012-9237-2.
- 564 [16] S. CHUBANOV, *A polynomial relaxation-type algorithm for linear programming*, Optimization Online,
565 February, (2011).
- 566 [17] G. CYBENKO, *Dynamic load balancing for distributed memory multiprocessors*, Journal of parallel and
567 distributed computing, 7 (1989), pp. 279–301.
- 568 [18] T. A. DAVIS AND Y. HU, *The University of Florida sparse matrix collection*, ACM T. Math. Software,
569 38 (2011), pp. 1–25.
- 570 [19] J. A. DE LOERA, J. HADDOCK, AND D. NEEDELL, *A sampling Kaczmarz-Motzkin algorithm for linear*
571 *feasibility*, SIAM J. Sci. Comput., 39 (2017), pp. S66–S87.
- 572 [20] K. DU AND H. GAO, *A new theoretical estimate for the convergence rate of the maximal weighted residual*
573 *Kaczmarz algorithm*, Numer. Math. - Theory Me., 12 (2019), pp. 627–639.
- 574 [21] B. DUMITRESCU, *On the relation between the randomized extended Kaczmarz algorithm and coordinate*
575 *descent*, BIT, (2014), pp. 1–11.
- 576 [22] Y. C. ELДАР AND D. NEEDELL, *Acceleration of randomized Kaczmarz method via the Johnson-*
577 *Lindenstrauss lemma*, Numer. Algorithms, 58 (2011), pp. 163–177, <https://doi.org/10.1007/s11075-011-9451-z>,
578 <http://dx.doi.org/10.1007/s11075-011-9451-z>.
- 579 [23] N. M. FRERIS AND A. ZOUZIAS, *Fast distributed smoothing of relative measurements*, in 2012 IEEE 51st
580 IEEE Conference on Decision and Control (CDC), IEEE, 2012, pp. 1411–1416.
- 581 [24] M. GEORGE, *B-A scale-free network generation and visualization*, 2020, [https://www.mathworks.com/](https://www.mathworks.com/matlabcentral/fileexchange/11947-b-a-scale-free-network-generation-and-visualization)
582 [matlabcentral/fileexchange/11947-b-a-scale-free-network-generation-and-visualization](https://www.mathworks.com/matlabcentral/fileexchange/11947-b-a-scale-free-network-generation-and-visualization).
- 583 [25] J.-L. GOFFIN, *The relaxation method for solving systems of linear inequalities*, Math. Oper. Res., 5 (1980),
584 pp. 388–414, <https://doi.org/10.1287/moor.5.3.388>, <http://dx.doi.org/10.1287/moor.5.3.388>.
- 585 [26] J.-L. GOFFIN, *On the nonpolynomiality of the relaxation method for systems of linear inequalities*, Math.
586 Program., 22 (1982), pp. 93–103, <https://doi.org/10.1007/BF01581028>, [http://dx.doi.org/10.1007/](http://dx.doi.org/10.1007/BF01581028)
587 [BF01581028](http://dx.doi.org/10.1007/BF01581028).
- 588 [27] R. GORDON, R. BENDER, AND G. T. HERMAN, *Algebraic reconstruction techniques (ART) for three-*
589 *dimensional electron microscopy and X-ray photography*, J. Theoret. Biol., 29 (1970), pp. 471–481.
- 590 [28] R. GOWER, D. MOLITOR, J. MOORMAN, AND D. NEEDELL, *Adaptive sketch-and-project methods for*
591 *solving linear systems*, arXiv preprint arXiv:1909.03604, (2019).
- 592 [29] R. M. GOWER AND P. RICHTÁRIK, *Randomized iterative methods for linear systems*, SIAM J. Matrix
593 Anal. A., 36 (2015), pp. 1660–1690.
- 594 [30] R. M. GOWER AND P. RICHTÁRIK, *Stochastic dual ascent for solving linear systems*, arXiv preprint
595 arXiv:1512.06890, (2015).
- 596 [31] J. HADDOCK AND D. NEEDELL, *On Motzkin’s method for inconsistent linear systems*, BIT, 59 (2019),
597 pp. 387–401.
- 598 [32] J. HADDOCK AND D. NEEDELL, *Randomized projection methods for linear systems with arbitrarily large*
599 *sparse corruptions*, SIAM J. Sci. Comput., 41 (2019), pp. S19–S36.
- 600 [33] M. HANKE AND W. NIETHAMMER, *On the acceleration of Kaczmarz’s method for inconsistent linear*
601 *systems*, Linear Algebra Appl., 130 (1990), pp. 83–98.
- 602 [34] F. HANZELY, J. KONEČNÝ, N. LOIZOU, P. RICHTÁRIK, AND D. GRISHCHENKO, *Privacy preserving ran-*
603 *domized gossip algorithms*, arXiv preprint arXiv:1706.07636, (2017).
- 604 [35] G. HERMAN AND L. MEYER, *Algebraic reconstruction techniques can be made computationally efficient*,
605 IEEE T. Med. Imaging, 12 (1993), pp. 600–609.
- 606 [36] S. KACZMARZ, *Angenäherte auflösung von systemen linearer gleichungen*, Bull. Int. Acad. Polon. Sci.
607 Lett. Ser. A, (1937), pp. 335–357.
- 608 [37] J. LIU, S. J. WRIGHT, AND S. SRIDHAR, *An asynchronous parallel randomized Kaczmarz algorithm*,
609 arXiv preprint arXiv:1401.4780, (2014).
- 610 [38] N. LOIZOU, M. RABBAT, AND P. RICHTÁRIK, *Provably accelerated randomized gossip algorithms*, in Int.

- 611 Conf. Acoust. Spee., IEEE, 2019, pp. 7505–7509.
- 612 [39] N. LOIZOU AND P. RICHTÁRIK, *A new perspective on randomized gossip algorithms*, in IEEE Glob. Conf.
613 Sig., IEEE, 2016, pp. 440–444.
- 614 [40] N. LOIZOU AND P. RICHTÁRIK, *Momentum and stochastic momentum for stochastic gradient, Newton,*
615 *proximal point and subspace descent methods*, arXiv preprint arXiv:1712.09677, (2017).
- 616 [41] N. LOIZOU AND P. RICHTÁRIK, *Accelerated gossip via stochastic heavy ball method*, in 2018 56th Annual
617 Allerton Conference on Communication, Control, and Computing (Allerton), IEEE, 2018, pp. 927–
618 934.
- 619 [42] N. LOIZOU AND P. RICHTÁRIK, *Convergence analysis of inexact randomized iterative methods*, arXiv
620 preprint arXiv:1903.07971, (2019).
- 621 [43] N. LOIZOU AND P. RICHTÁRIK, *Revisiting randomized gossip algorithms: General framework, convergence*
622 *rates and novel block and accelerated protocols*, arXiv preprint arXiv:1905.08645, (2019).
- 623 [44] A. MA, D. NEEDELL, AND A. RAMDAS, *Convergence properties of the randomized extended Gauss–Seidel*
624 *and Kaczmarz methods*, SIAM J. Matrix Anal. A., 36 (2015), pp. 1590–1604.
- 625 [45] J. D. MOORMAN, T. K. TU, D. MOLITOR, AND D. NEEDELL, *Randomized Kaczmarz with averaging*,
626 arXiv preprint arXiv:2002.04126, (2020).
- 627 [46] M. S. MORSHED, M. S. ISLAM, ET AL., *On generalization and acceleration of randomized projection*
628 *methods for linear feasibility problems*, arXiv preprint arXiv:2002.07321, (2020).
- 629 [47] M. S. MORSHED, M. S. ISLAM, AND M. NOOR-E-ALAM, *Accelerated sampling Kaczmarz Motzkin algo-*
630 *rithm for the linear feasibility problem*, J. Global Optim., (2019), pp. 1–22.
- 631 [48] T. S. MOTZKIN AND I. J. SCHOENBERG, *The relaxation method for linear inequalities*, Canadian J. Math.,
632 6 (1954), pp. 393–404.
- 633 [49] 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>.
- 634 [50] D. NEEDELL AND E. REBROVA, *On block Gaussian sketching for iterative projections*, arXiv preprint
635 arXiv:1905.08894, (2019).
- 636 [51] D. NEEDELL, N. SREBRO, AND R. WARD, *Stochastic gradient descent and the randomized Kaczmarz*
637 *algorithm*, Math. Program. A, 155 (2016), pp. 549–573.
- 638 [52] D. NEEDELL AND J. A. TROPP, *Paved with good intentions: Analysis of a randomized block Kaczmarz*
639 *method*, Linear Algebra Appl., (2013).
- 640 [53] D. NEEDELL, R. ZHAO, AND A. ZOUZIAS, *Randomized block Kaczmarz method with projection for solving*
641 *least squares*, Linear Algebra Appl., 484 (2015), pp. 322–343.
- 642 [54] J. NUTINI, B. SEPEHRY, A. VIRANI, I. LARADJI, M. SCHMIDT, AND H. KOEPKE, *Convergence Rates for*
643 *Greedy Kaczmarz Algorithms*, UAI, (2016).
- 644 [55] S. PETRA AND C. POPA, *Single projection Kaczmarz extended algorithms*, Numer. Algorithms, (2015),
645 pp. 1–16, <https://doi.org/10.1007/s11075-016-0118-7>, <https://arxiv.org/abs/1504.00231>.
- 646 [56] C. POPA, *A fast Kaczmarz-Kovarik algorithm for consistent least-squares problems*, Korean J. Comput.
647 Appl. Math., 8 (2001), pp. 9–26.
- 648 [57] C. POPA, *A Kaczmarz-Kovarik algorithm for symmetric ill-conditioned matrices*, An. Ştiinţ. Univ. Ovidius
649 Constanţa Ser. Mat., 12 (2004), pp. 135–146.
- 650 [58] C. POPA, T. PRECLIK, H. KÖSTLER, AND U. RÜDE, *On Kaczmarz’s projection iteration as a direct solver*
651 *for linear least squares problems*, Linear Algebra Appl., 436 (2012), pp. 389–404.
- 652 [59] E. REBROVA AND D. NEEDELL, *Sketching for Motzkin’s iterative method for linear systems*, Proc. 50th
653 Asilomar Conf. on Signals, Systems and Computers, (2019).
- 654 [60] T. STROHMER AND R. VERSHYNIN, *A randomized Kaczmarz algorithm with exponential convergence*, J.
655 Fourier Anal. Appl., 15 (2009), pp. 262–278.
- 656 [61] SYSTEMS OPTIMIZATION LABORATORY, *CGLS: CG method for $Ax = b$ and least squares*, <https://web.stanford.edu/group/SOL/software/cgls/>.
- 657 [62] J. TELGEN, *On relaxation methods for systems of linear inequalities*, Eur. J. Oper. Res., 9
658 (1982), pp. 184–189, [https://doi.org/10.1016/0377-2217\(82\)90071-6](https://doi.org/10.1016/0377-2217(82)90071-6), [http://dx.doi.org/10.1016/0377-2217\(82\)90071-6](http://dx.doi.org/10.1016/0377-2217(82)90071-6).
- 659 [63] D. USTEBAY, B. N. ORESHKIN, M. J. COATES, AND M. G. RABBAT, *Greedy gossip with eavesdropping*,
660 IEEE T. Signal Proces., 58 (2010), pp. 3765–3776.
- 661 [64] H. XIANG AND L. ZHANG, *Randomized iterative methods with alternating projections*, arXiv preprint
662

- 665 arXiv:1708.09845, (2017).
- 666 [65] L. XIAO AND S. BOYD, *Fast linear iterations for distributed averaging*, Syst. Control Lett., 53 (2004),
667 pp. 65–78.
- 668 [66] L. XIAO, S. BOYD, AND S. LALL, *A scheme for robust distributed sensor fusion based on average consen-*
669 *sus*, in IPSN 2005. Fourth International Symposium on Information Processing in Sensor Networks,
670 2005., IEEE, 2005, pp. 63–70.
- 671 [67] J. ZHANG, J. CUI, Z. WANG, Y. DING, AND Y. XIA, *Distributed joint cooperative self-localization and*
672 *target tracking algorithm for mobile networks*, Sensors, 19 (2019), p. 3829.
- 673 [68] A. ZOUZIAS AND N. M. FRERIS, *Randomized extended Kaczmarz for solving least squares*, SIAM J. Matrix
674 Anal. A., 34 (2013), pp. 773–793.
- 675 [69] A. ZOUZIAS AND N. M. FRERIS, *Randomized gossip algorithms for solving Laplacian systems*, in 2015
676 European Control Conference (ECC), IEEE, 2015, pp. 1920–1925.