

# ADAPTED STOCHASTIC GRADIENT DESCENT FOR LINEAR SYSTEMS WITH MISSING DATA

ANNA MA, DEANNA NEEDELL

ABSTRACT. Stochastic iterative methods with low memory footprints are extremely useful when solving large linear systems. Utilizing large amounts of data is further complicated when the data is incomplete or has missing entries. In this work, we address the issues of having extremely large amounts of data and incomplete data simultaneously. In particular, we propose to adapt the Stochastic Gradient Descent method to address missing data in linear systems. Our proposed algorithm, the Adapted Stochastic Gradient Descent for Missing Data method (mSGD), is introduced and theoretical convergence guarantees are provided. In addition, we include numerical experiments on simulated and real world data that demonstrate the usefulness of our method.

## 1. INTRODUCTION

Traditional methods for solving linear systems have quickly become impractical due to an increase in the size of available data. When handling large amounts of data, it may not be possible to load the entire matrix (data set) into memory, as typically required by matrix inversions or matrix factorization. This has led to the study and advancement of stochastic iterative methods with low memory footprints such as Stochastic Gradient Descent, Randomized Kaczmarz, and Randomized Gauss-Seidel [MNR15, NWS14, SV09, LL10]. The need for algorithms that can process large amounts of information is further complicated by incomplete or missing data, which can arise due to, for example, attrition, errors in data recording, or cost of data acquisition. Standard methods for treating missing data, which include data imputation [Efr94, FC03], matrix completion [CCS10, KOM09, Rec11, KMO10], and maximum likelihood estimation [DLR77, LR14] can be wasteful, create biases, or be impractical for extremely large amounts of data. This work simultaneously addresses both issues of enormous data size and missing data.

Consider the system of linear equations  $\mathbf{Ax} = \mathbf{b}$ <sup>1</sup>, where  $\mathbf{A} \in \mathbb{C}^{m \times n}$  is a large, full-rank, overdetermined ( $m > n$ ) matrix. Suppose that  $\mathbf{A}$  is not known entirely, but instead only some of its entries are available. As a concrete example, suppose  $\mathbf{A}$  is the rating matrix from the survey of  $m$  users about  $n$  service questions, and  $\mathbf{b}$  contains the  $m$  “overall” ratings from each user (which is fully known). Each user may not answer all of the individual service questions, but the company wishes to understand how each question affects the overall rating of the user. That is, given partial knowledge of  $\mathbf{A}$ , one wishes to uncover  $\mathbf{x}_* = \arg \min_{\mathbf{x}} \frac{1}{2m} \|\mathbf{Ax} - \mathbf{b}\|^2$ .

---

<sup>1</sup> The linear system is not assumed to be consistent; we will use the notation  $\mathbf{Ax} = \mathbf{b}$  to denote a general linear system.

Let  $\mathbf{A}$  denote the complete matrix and  $\tilde{\mathbf{A}}$  be a matrix containing the known entries of  $\mathbf{A}$ . Formally, one wants to solve the following optimization program:

$$\begin{aligned} &\text{Given } \tilde{\mathbf{A}}, \mathbf{b} \text{ s.t. } \mathbf{Ax} = \mathbf{b} \\ &\text{Find } \mathbf{x}_\star = \arg \min_{\mathbf{x} \in \mathcal{W}} \frac{1}{2m} \|\mathbf{Ax} - \mathbf{b}\|^2, \end{aligned} \tag{1.1}$$

where  $\mathcal{W}$  is a convex domain containing the solution  $\mathbf{x}_\star$  (e.g. a ball with large enough radius).

**Contributions.** This work presents a stochastic iterative projection method for solving large-scale linear systems with missing data. We provide theoretical bounds for the proposed method's performance and demonstrate its usefulness on simulated and real world data sets.

**1.1. Stochastic Gradient Descent.** Stochastic iterative methods such as the Randomized Kaczmarz algorithm and Stochastic Gradient Descent (SGD) have gained interest in recent years due to their simplicity and ability to handle large-scale systems. Originally discussed in [RM51], SGD has gained recent attention for its usefulness in large-scale machine learning problems [Bot10, Zha04, Bot12]. SGD minimizes an objective function  $F(\mathbf{x})$  over a convex domain  $\mathcal{W}$  using unbiased estimates for the gradient of the objective, i.e., using  $f_i(\mathbf{x})$  such that  $\mathbb{E}[\nabla f_i(\mathbf{x})] = \nabla F(\mathbf{x})$ . At each iteration, a random unbiased estimate,  $\nabla f_i(\mathbf{x})$ , is drawn and the minimizer of  $F(\mathbf{x})$  is estimated with:

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \mathcal{P}_{\mathcal{W}}(\alpha_k \nabla f_i(\mathbf{x}_{k-1})), \tag{1.2}$$

where  $\alpha_k$  is the appropriately chosen step size at iteration  $k$  and  $\mathcal{P}_{\mathcal{W}}$  denotes the projection onto the convex set  $\mathcal{W}$ . To solve a linear system  $\mathbf{Ax} = \mathbf{b}$ , one can minimize the least-squares objective function  $F(\mathbf{x}) = \frac{1}{2m} \|\mathbf{Ax} - \mathbf{b}\|^2 = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x})$  where  $f_i(\mathbf{x}) = \frac{1}{2} (\mathbf{A}_i \mathbf{x} - \mathbf{b}_i)^2$  and  $\mathbf{A}_i$  denotes the  $i^{\text{th}}$  row of  $\mathbf{A}$ . SGD randomly chooses a row  $i \in \{1, 2, \dots, m\}$  of  $\mathbf{A}$  and computes (1.2) with  $\nabla f_i(\mathbf{x}_{k-1}) = \mathbf{A}_i^* (\mathbf{A}_i \mathbf{x}_{k-1} - \mathbf{b}_i)$  where  $*$  denotes the conjugate transpose.

The performance of SGD on linear systems depends on the choice of  $\alpha_k$  and the consistency of the system (i.e. whether a solution to the system exists). When the linear system is consistent, SGD achieves linear convergence with an appropriately chosen fixed step size [SR13]. For example, the Randomized Kaczmarz method, a special instance of SGD, has been shown to converge exponentially for consistent systems without decreasing step sizes [Kac37, SV09, NWS14]. Unfortunately, this is not the case when the system is inconsistent. When the linear system is inconsistent, or  $\mathbf{Ax} \approx \mathbf{b}$ , one must use decreasing step sizes to obtain the optimum (see e.g. [SR13, HN90, CEG83]). This phenomenon is explained by the norm of the unbiased estimates at the minimizer,  $\|\nabla f_i(\mathbf{x}_\star)\|^2$ . For consistent systems,  $\|\nabla f_i(\mathbf{x}_\star)\|^2 = \|\mathbf{A}_i^* (\mathbf{A}_i \mathbf{x}_\star - \mathbf{b}_i)\|^2 = 0$  since  $\mathbf{Ax}_\star = \mathbf{b}$ . Intuitively, as SGD progresses closer to the minimizer, the magnitude of the iterates get smaller and allow SGD to converge. When the system is inconsistent,  $\mathbf{Ax}_\star = \mathbf{b} + \mathbf{r}$  for some residual vector  $\mathbf{r}$  and  $\|\nabla f_i(\mathbf{x}_\star)\|^2 = \|\mathbf{A}_i^* (\mathbf{A}_i \mathbf{x}_\star - \mathbf{b}_i)\|^2 = \mathbf{r}_i^2 \|\mathbf{A}_i\|^2$ . As SGD gets closer to the solution, the magnitude of the iterates do not converge to 0. Using diminishing step sizes dampens the magnitude of the iterates over time, allowing SGD to converge. When SGD with fixed step size is applied to inconsistent systems, the iterates oscillate within a fixed distance from the solution [NWS14]. The fixed distance, also referred to as the convergence horizon, is proportional to the step size but inversely proportional to the rate of convergence. Therefore, there is a trade-off between the rate of convergence of SGD and the radius of convergence.

The proposed method, which we refer to as mSGD, is an SGD-type iterate with a correction term that takes into account the fact that not all entries of  $\mathbf{A}$  are available. We start with a discussion on the model under which mSGD operates and proceed to derive the iterate. After the introduction of the algorithm, the formal results are stated.

**Outline.** Section 2 introduces the proposed method, the Adapted Stochastic Gradient Descent for Missing Data method (mSGD) and the main theoretical results. The performance of mSGD on simulated data as well as real world data are shown in Section 3. Finally, we end with concluding remarks.

## 2. ADAPTED STOCHASTIC GRADIENT DESCENT FOR MISSING DATA

**2.1. Assumptions and Notation.** Suppose one wishes to solve  $\arg \min_x \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  where  $\mathbf{A} \in \mathbb{C}^{m \times n}$  ( $m > n$ ) but not all entries of  $\mathbf{A}$  are accessible. Instead, only a random subset of entries of  $\mathbf{A}$  are available. As a simplifying assumption, we model whether an entry of  $\mathbf{A}$  as missing by i.i.d. Bernoulli random variables that are equal to 1 with probability  $p$ . Practically, there are many applications in which this type of assumption holds. For example, suppose one does not have access to a matrix  $\mathbf{A}$  but instead new measurements with missing data are streaming in continuously, as in an online setting. In this case, if elements are missing at random, the simplifying assumption holds. As another example, consider an extremely large  $m \times n$  matrix  $\mathbf{A}$  where it is not possible to load entire rows of  $\mathbf{A}$  nor columns of  $\mathbf{A}$  due to memory constraints. Instead, one is restricted to only loading  $[pn]$  (random) elements of  $\mathbf{A}$  at a time. The proposed method solves  $\arg \min_x \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  in this example without making any additional assumptions on the structure of  $\mathbf{A}$  such as sparsity or low-rankness. In the case of having a fixed matrix  $\tilde{\mathbf{A}} \in \mathbb{C}^{m \times n}$  with missing entries, the theoretical results hold only if each row of the matrix is utilized once. If  $\tilde{\mathbf{A}}$  is an extremely overdetermined matrix (i.e.  $m \gg n$ ), then this is a reasonable assumption.

**Notation.** Here and throughout the remainder of this paper, the following notation is adopted. Let  $\mathbf{D}$  be an  $m \times n$  matrix where the entries of  $\mathbf{D}$ , denoted by  $\delta_{i,j}$  for  $i = 1, 2, \dots, m$  and  $j = 1, 2, \dots, n$ , are drawn i.i.d. from a Bernoulli distribution that are equal to 1 with probability  $p$  so that  $\mathbb{E}[\delta_{i,j}] = p$ . Let  $\mathbf{D}_i$  be the diagonal matrix whose diagonal is equal to the  $i$ th row of  $\mathbf{D}$ . The matrix  $\mathbf{D}$  is referred to as a *binary mask*. Let  $\tilde{\mathbf{A}}$  represent the matrix  $\mathbf{A}$  with missing elements filled in with zeros so that  $\tilde{\mathbf{A}} = \mathbf{D} \circ \mathbf{A}$  and  $\tilde{\mathbf{A}}_i = \mathbf{D}_i \mathbf{A}_i$ , where  $\circ$  denotes the element-wise product. Additionally, let  $\sigma_{\min}^2(\mathbf{A})$  be the smallest singular value of  $\mathbf{A}$  and  $\|\cdot\|$  denote the  $\ell_2$ -norm. Denote the expected value taken over the random selection of a of  $\tilde{\mathbf{A}}$  as  $\mathbb{E}_i[\cdot]$ , the expected value taken over all  $(2^{mn})$  possible binary masks  $\mathbf{D}$  as  $\mathbb{E}_\delta[\cdot]$ , and the full expected value as  $\mathbb{E}[\cdot]$ . Let  $\mathcal{W}$  be some convex domain containing  $\mathbf{x}_*$  and  $B := \max_{\mathbf{x} \in \mathcal{W}} \|\mathbf{x}\|^2$ .

Suppose one naively applies SGD to the system  $\tilde{\mathbf{A}}\mathbf{x} = \mathbf{b}$ . Since  $\mathbf{A}$  is not accessible, consider the objective  $\tilde{F}(\mathbf{x}) = \frac{1}{2m} \|\tilde{\mathbf{A}}\mathbf{x} - \mathbf{b}\|^2 = \frac{1}{m} \sum_{i=1}^m \tilde{f}_i(\mathbf{x})$  where  $\tilde{f}_i(\mathbf{x}) = \frac{1}{2} (\tilde{\mathbf{A}}_i \mathbf{x} - \mathbf{b}_i)^2$ . This objective function leads to the update:

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \alpha_k \left( \tilde{\mathbf{A}}_i^* (\tilde{\mathbf{A}}_i \mathbf{x}_{k-1} - \mathbf{b}_i) \right),$$

since  $\nabla \tilde{f}_i(\mathbf{x}) = \tilde{\mathbf{A}}_i^*(\tilde{\mathbf{A}}_i\mathbf{x} - \mathbf{b}_i)$ . Unfortunately, one computes that, taking the expectation with respect to the binary mask and gradient direction,

$$\begin{aligned} \mathbb{E}_i\mathbb{E}_\delta[\nabla \tilde{f}_i(\mathbf{x})] &= \mathbb{E}_i\mathbb{E}_\delta[\tilde{\mathbf{A}}_i^*(\tilde{\mathbf{A}}_i\mathbf{x} - \mathbf{b}_i)] \\ &= \frac{1}{m} \left( p^2 \mathbf{A}^* \mathbf{A} \mathbf{x} + (p - p^2) \text{diag}(\mathbf{A}^* \mathbf{A}) \mathbf{x} - p \sum_i \mathbf{A}_i^* \mathbf{b}_i \right) \\ &\neq \nabla F(\mathbf{x}). \end{aligned}$$

As a result, the iterates are not moving in the gradient descent direction toward the desired solution on average. Instead of using  $\nabla \tilde{f}_i(\mathbf{x})$  as the step direction, we propose to use  $\nabla \tilde{f}_i(\mathbf{x})$  to estimate  $\nabla F(\mathbf{x})$ . In other words, we want to represent  $\nabla F(\mathbf{x})$  in terms of  $\mathbb{E}[\nabla \tilde{f}_i(\mathbf{x})]$ . By doing so, iterates  $\mathbf{x}_k$  move in the gradient descent direction towards the true solution. Substituting  $\mathbf{b}_i$  with  $p\mathbf{b}_i$ ,

$$\nabla F(\mathbf{x}) = \frac{1}{p^2} \mathbb{E}[\nabla \tilde{f}_i(\mathbf{x})] - \frac{(1-p)}{p^2} \mathbb{E}[\text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i)] \mathbf{x}.$$

The detailed computation is available in the Appendix (Lemma 5.1). Therefore the appropriate update is

$$\mathbf{x}_k = \mathbf{x}_{k-1} - \alpha_k \left( \frac{1}{p^2} \left( \tilde{\mathbf{A}}_i^*(\tilde{\mathbf{A}}_i\mathbf{x}_{k-1} - p\mathbf{b}_i) \right) - \frac{1-p}{p^2} \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \mathbf{x}_{k-1} \right).$$

In classical SGD literature, the expected value is taken over the row choice  $i$ . However, in this setting there are two sources of randomness: the randomness from row selection and the randomness incurred from modeling whether or not data is missing. In this computation, the expected value is taken with respect to both sources of randomness. The method is outlined in Algorithm 1.

---

**Algorithm 1** Adapted Stochastic Gradient Descent for Missing Data (mSGD)

---

- 1: **procedure** ( $\tilde{\mathbf{A}}, \mathbf{b}, T, p, \{\alpha_k\}$ ) ▷ If using a fixed step size  $\alpha$ ,  $\alpha_k = \alpha$  for all  $k$ .
  - 2:   Initialize  $\mathbf{x}_0$
  - 3:   **for**  $k = 1, 2, \dots, T$  **do**
  - 4:     Choose row  $i$  of  $\tilde{\mathbf{A}}$  with probability  $\frac{1}{m}$
  - 5:      $g(\mathbf{x}_{k-1}) = \frac{1}{p^2} \left( \tilde{\mathbf{A}}_i^*(\tilde{\mathbf{A}}_i\mathbf{x}_{k-1} - p\mathbf{b}_i) \right) - \frac{1-p}{p^2} \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \mathbf{x}_{k-1}$
  - 6:      $\mathbf{x}_k = \mathcal{P}_{\mathcal{W}}(\mathbf{x}_{k-1} - \alpha g(\mathbf{x}_{k-1}))$  ▷  $\mathcal{P}_{\mathcal{W}}$  is the projection onto the set  $\mathcal{W}$ .
  - 7:   **end for**
  - 8:   Output  $\mathbf{x}_k$
  - 9: **end procedure**
- 

**2.2. Main Results.** Before the main results are presented, note the following properties of the objective function,

$$F(\mathbf{x}) = \frac{1}{2m} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2, \quad (2.1)$$

and the update function in Algorithm 1 (Line 5),

$$g(\mathbf{x}) = \frac{1}{p^2} \left( \tilde{\mathbf{A}}_i^*(\tilde{\mathbf{A}}_i\mathbf{x} - p\mathbf{b}_i) \right) - \frac{(1-p)}{p^2} \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \mathbf{x}_{k-1}, \quad (2.2)$$

as they play an important role in the convergence analysis of mSGD.

- The objective function (2.1) is  $\mu$ -strongly convex. For all  $\mathbf{x}, \mathbf{y} \in \mathcal{W}$ ,

$$(\mathbf{x} - \mathbf{y})^*(\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})) \geq \mu \|\mathbf{x} - \mathbf{y}\|^2,$$

where

$$\mu = \sigma_{\min}^2(\mathbf{A}). \quad (2.3)$$

- The update function  $g(\mathbf{x})$  is Lipschitz continuous, has Lipschitz constant  $L_{i,D}$  (for a fixed instance of  $i$  and  $D$ ), and supremum Lipschitz constant  $L_g$ . In other words, for all  $\mathbf{x}, \mathbf{y} \in \mathcal{W}$ ,

$$\|g(\mathbf{x}) - g(\mathbf{y})\| \leq L_{i,D} \|\mathbf{x} - \mathbf{y}\| \quad (2.4)$$

$$L_g = \sup_{i,D} L_{i,D}. \quad (2.5)$$

The supremum is taken over all choices of rows and all possible binary masks (i.e. all  $2^{mn}$  possible binary masks).

- There exists a constant  $G$  that uniformly bounds the expected norm of  $\|g(\mathbf{x})\|^2$ ,

$$\mathbb{E}[\|g(\mathbf{x})\|^2] \leq G, \quad (2.6)$$

for all  $\mathbf{x} \in \mathcal{W}$ . The expected norm of  $g(\mathbf{x}_\star)$  plays an important role in size of the convergence horizon. For this reason, let  $G_\star$  denote the upper bound of  $\mathbb{E}[\|g(\mathbf{x}_\star)\|^2]$ :

$$\mathbb{E}[\|g(\mathbf{x}_\star)\|^2] \leq G_\star. \quad (2.7)$$

The computation of  $L_{i,D}$  and  $L_g$  are shown in Lemma 5.2. Lemma 5.3 shows the computation of  $G$  and  $G_\star$ . The statements and proofs of both lemmas are provided in the Appendix so that we may proceed to the presentation of the main results.

Theorem 2.1 below shows that, in expectation, Algorithm 1 converges to the least squares solution of the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  with properly chosen step size. This theorem is an application of the previously proven result stated in Lemma 2.2. Because there is a trade-off between convergence rate and convergence to the optimal solution, the fixed step size setting is also discussed. Theorem 2.3 shows, using a fixed step size, Algorithm 1 converges to some convergence horizon. In addition, we provide an optimal step size choice based on a desired error tolerance,  $\epsilon$ , and a bound on the number of iterations required to obtain said tolerance in Corollary 2.4. Lastly, we remark on the recovery of classical SGD when  $p = 1$  both algorithmically and with respect to the proven error bounds.

**Theorem 2.1.** *Consider (1.1) with  $\tilde{\mathbf{A}} = \mathbf{D} \circ \mathbf{A}$  where entries of  $\mathbf{D}$  are drawn i.i.d. from a Bernoulli distribution and are equal to 1 with probability  $p$ . Let  $\mu$  be as defined in (2.3). Choosing  $\alpha_k = \frac{1}{\mu k}$ , Algorithm 1 converges in expectation with error*

$$\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}_\star\|^2] \leq \frac{17G(1 + \log(k))}{\mu^2 k},$$

where  $G = \frac{2B(2+p)(1-p)}{mp^3} \sum_i \|\mathbf{A}_i\|^4 + \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 |\mathbf{b}_i|^2$  and  $B = \max_{\mathbf{x} \in \mathcal{W}} \|\mathbf{x}\|^2$ .

It is clear that the convergence behavior of Algorithm 1 depends on  $G$ , the uniform upper bound on the expected norm of  $g(\mathbf{x})$  and on  $\sigma_{\min}^2(\mathbf{A})$ . As one would expect, the more data that is missing, the larger the upper bound on expected error. In particular, assuming all other variables

are constant and  $p \in (0, 1]$ , as  $p$  decreases,  $G$  increases. Note that as the number of iterations increases, Algorithm 1 still converges in expectation. Theorem 2.1 is an application of the following previously proved lemma.

**Lemma 2.2.** ([SZ13] Theorem 1) *Let  $F(\mathbf{x})$  be a  $\mu$ -strongly convex objective function,  $g(\mathbf{x})$  be such that  $\mathbb{E}[g(\mathbf{x})] = \nabla F(\mathbf{x})$ , and  $\mathbb{E}[\|g(\mathbf{x})\|^2] \leq G$  for all  $\mathbf{x} \in \mathcal{W}$ . Using step size  $\alpha_k = \frac{1}{\mu k}$  and update  $\mathbf{x}_k = \mathbf{x}_{k-1} - \alpha_k \mathcal{P}_{\mathcal{W}}(g(\mathbf{x}_{k-1}))$ , it holds that*

$$\mathbb{E}[F(\mathbf{x}_k) - F(\mathbf{x}_*)] \leq \frac{17G(1 + \log(k))}{\mu k}.$$

The next theorem details the convergence behavior of Algorithm 1 when using a fixed step size. Theorem 2.3 shows that Algorithm 1 experiences a convergence horizon that depends on  $L_g$  and  $G_*$ . For  $p \in (0, 1]$ , as  $p$  decreases,  $G_*$  and  $L_g$  both increase. Intuitively this makes sense as a larger amount of missing data should increase the size of the convergence horizon. Additionally, the convergence rate  $r = (1 - 2\alpha\mu(1 - \alpha L_g))$  also increases as  $p$  decreases. In other words, more missing data causes a slower convergence rate.

**Theorem 2.3.** *Consider (1.1) with  $\tilde{\mathbf{A}} = \mathbf{D} \circ \mathbf{A}$  where entries of  $\mathbf{D}$  are drawn i.i.d. from a Bernoulli distribution and are equal to 1 with probability  $p$ . Let  $L_g$ ,  $G_*$ , and  $\mu$  be as defined in (2.5), (2.7), and (2.3) respectively. Additionally, let the fixed step size be  $\alpha < \frac{1}{L_g}$ . Algorithm 1 converges with expected error*

$$\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}_*\|^2] \leq r^k \|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \frac{\alpha G_*}{\mu(1 - \alpha L_g)}. \quad (2.8)$$

where  $r = (1 - 2\alpha\mu(1 - \alpha L_g))$ ,  $L_g = \frac{2-p}{p^2} \sup_i \|\mathbf{A}_i\|^2$ , and  $\mu = \sigma_{\min}^2(\mathbf{A})$ . If  $\mathbf{Ax} = \mathbf{b}$  is consistent,  $G_* = \frac{2(1-p)(2-p)}{mp^3} \|\mathbf{x}_*\|^2 \sum_i \|\mathbf{A}_i\|^4$ . If the linear system is inconsistent (i.e.  $\mathbf{Ax} = \mathbf{b} + \mathbf{r}$  for some residual vector  $\mathbf{r}$ ), then  $G_* = \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 r_i^2 + \frac{2(1-p)(2-p)}{mp^3} \|\mathbf{x}_*\|^2 \sum_i \|\mathbf{A}_i\|^4$ .

Corollary 2.4 and the subsequent remark comment on the number of iterations required by Algorithm 1 to obtain some desired error tolerance  $\epsilon$  using a particular fixed step size  $\alpha^*$ . The corollary itself details this information in terms of the variables in Theorem 2.3 while the remark translates and simplifies  $\alpha^*$  and  $k$  (the number of iterations) into terms relating to  $\mathbf{A}$ . Note that the number of iterations required to reach a specified tolerance is a function of the ratio between the log of the initial error and  $\epsilon$ . The number of iterations increase as  $\epsilon$  decreases. Additionally, the remark shows that as  $p$  decreases, or as less data becomes available, more iterations are required to obtain an expected error of  $\epsilon$ .

**Corollary 2.4.** *Given an initial error  $\epsilon_0$  and choosing the fixed step size*

$$\alpha^* = \frac{\epsilon\mu}{2G_* + 2\mu\epsilon L_g},$$

after

$$k = 2 \log \left( \frac{2\epsilon_0}{\epsilon} \right) \left( \frac{L_g}{\mu} + \frac{G_*}{\mu^2\epsilon} \right)$$

iterations of Algorithm 1,  $\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}_*\|^2] \leq \epsilon$  holds in expectation.

*Remark.* Let  $\mathbf{a}_{\max}^2 = \max_i \|\mathbf{A}_i\|^2$  be the maximum squared row norm of  $\mathbf{A}$ . Given an initial error  $\epsilon_0$  and a desired tolerance  $\epsilon$  to the true solution, choosing the fixed step size

$$\alpha^* = \frac{p^3 \epsilon \sigma_{\min}^2(\mathbf{A})}{4(2-p)(1-p)\|\mathbf{b}\|^4 + 2(2-p)p\epsilon \sigma_{\min}^2(\mathbf{A})\mathbf{a}_{\max}^2},$$

after

$$k = 2 \log \left( \frac{2\epsilon_0}{\epsilon} \right) \left( \frac{(2-p)\mathbf{a}_{\max}^2}{p^2 \sigma_{\min}^2(\mathbf{A})} + \frac{2(2-p)(1-p)\|\mathbf{b}\|^4}{p^3 \sigma_{\min}^4(\mathbf{A})\epsilon} \right)$$

iterations of Algorithm 1,  $\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}_*\|^2] \leq \epsilon$  holds in expectation.

**Recovering SGD.** When  $p = 1$ , Algorithm 1 behaves as classical SGD does on the full linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Additionally, mSGD experiences similar convergence bounds as classical SGD for fixed step sizes [NWS14]. In particular, when  $p = 1$  the updating function  $g(\mathbf{x})$  reduces to  $g(\mathbf{x}) = \mathbf{A}_i^T(\mathbf{A}_i\mathbf{x} - \mathbf{b}_i)$ .

### 3. EXPERIMENTS

This section demonstrates the usefulness of Algorithm 1 on synthetic and real world data. Although the full data set is available in every experiment, missing data is simulated by computing a binary mask that dictates which elements are available and which are not at every iteration. By doing so, the simplifying assumption is satisfied, and the ground truth is known, so error can be computed. In each experiment, the percentage of available data is varied and log  $\ell_2$ -error to the least squares solution,  $\|\mathbf{x}_k - \mathbf{x}_*\|^2$  is averaged over 20 trials. For the fixed step size in simulated data,  $\alpha = 10^{-4}$  and for real world data  $\alpha = 10^{-5}$ . For updating step size,  $\alpha_k = \frac{c}{\sigma_{\min}^2(\mathbf{A})k}$  with  $c = 10^{-2}$ . Using  $\alpha_k = \frac{1}{\sigma_{\min}^2(\mathbf{A})k}$  (as described in Theorem 2.1) creates an initial increase in error followed by a decrease in error. This behavior is attributed to the step sizes being too large initially. It seems that the factor  $c$  can be optimized but we do not attempt to optimize such parameters here.

First, the performance of mSGD on standard Gaussian matrices is investigated. The results can be seen in Figure 1 and Figure 2. Here, elements of  $\mathbf{A} \in \mathbb{R}^{m \times n}$  are drawn i.i.d. from a standard Gaussian distribution where  $m = 1000$  and  $n = 200$ . Figure 1A and Figure 2A show the results of Algorithm 1 using a fixed step size ( $\alpha = 10^{-4}$ ) while Figure 1B and Figure 2B show results using updating step sizes. For inconsistent systems, we use  $\mathbf{b} + \mathbf{r}$  as the right hand side vector where  $\mathbf{r}$  is computed such that  $\mathbf{r} \in \text{null}(\mathbf{A}^*)$ .

The first real world data set was obtained from the UCI Machine Learning Repository [Lic13] and contains data from a bike rental service. Rows of  $\mathbf{A}$  contain hourly information from a bike share rental system and columns contain information such as weather, total number of rented bikes, time, and day of the week. In this experiment,  $m = 17379$  and  $n = 9$ . Figure 3 displays the performance of mSGD on this data set for fixed and updating step sizes.

The performance of Algorithm 1 on 2009 Adolescent California Health Interview Survey data [fHPR09] is shown in Figure 4. This data set contains information gathered via telephone interviews on health status, conditions, and insurance information from various adolescents. In this data set,  $m = 19516$  and  $n = 25$ . Figure 4 shows the results of mSGD on this data set using both fixed and updating step sizes. Notice that in Figure 4B for  $p = .3$  there is an initial increase in the

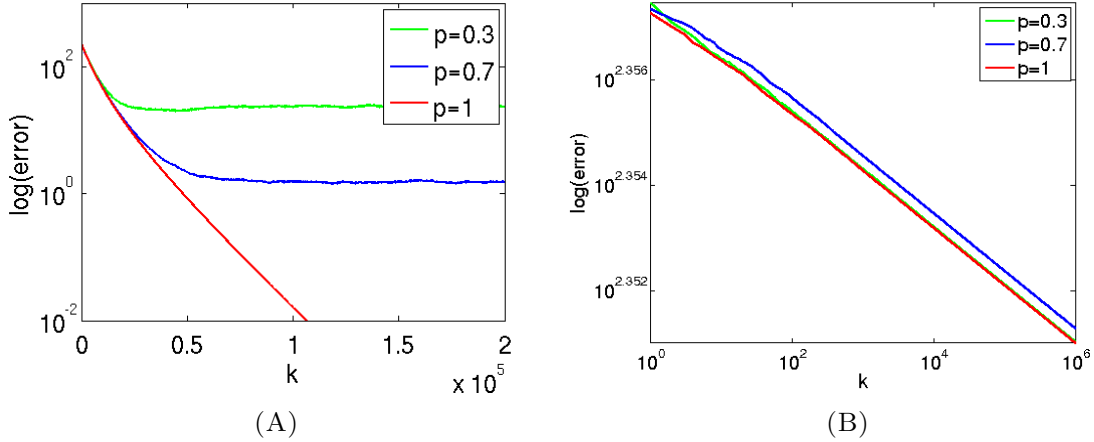


FIGURE 1. This figure compares the performance of Algorithm 1 on linear systems drawn from a standard Gaussian distribution. The percentage of data that is missing is varied. The x-axis is  $\log(\text{iteration})$  and the y-axis is the  $\log(\ell_2\text{-error})$ . Note that using a fixed step size (left), allows mSGD to converge much faster but to some convergence horizon. Using updating step sizes (right), continual progress is made at the cost of slower convergence.

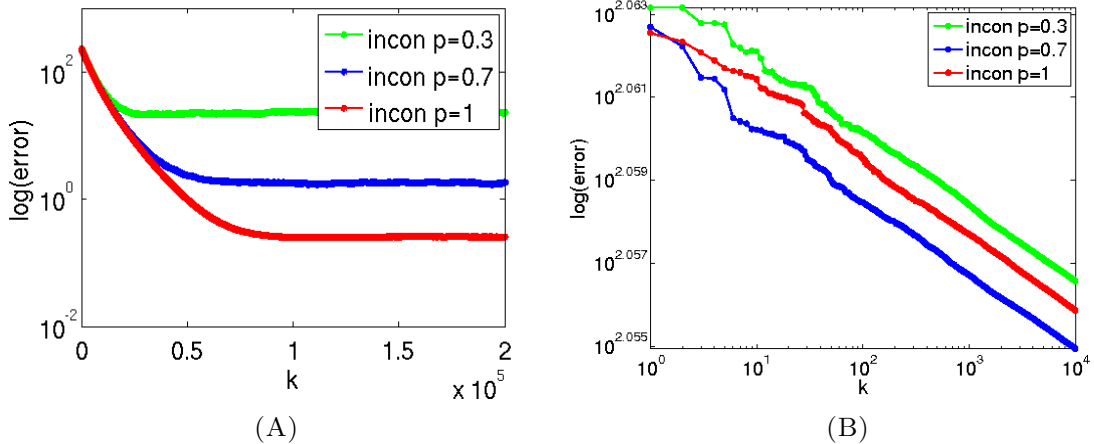


FIGURE 2. The performance of mSGD on inconsistent linear systems. For a fixed step size (left), mSGD converges to a convergence horizon and using updating step sizes (right) allows mSGD to continually progress at a slower rate.

$\ell_2$ -error. We attribute this to not choosing a small enough  $c$  where  $\alpha_k = \frac{c}{\sigma_{\min}^2(\mathbf{A})^k}$ . Since there is more missing data when  $p = 0.3$ , intuitively one would need to take much smaller steps. The initial



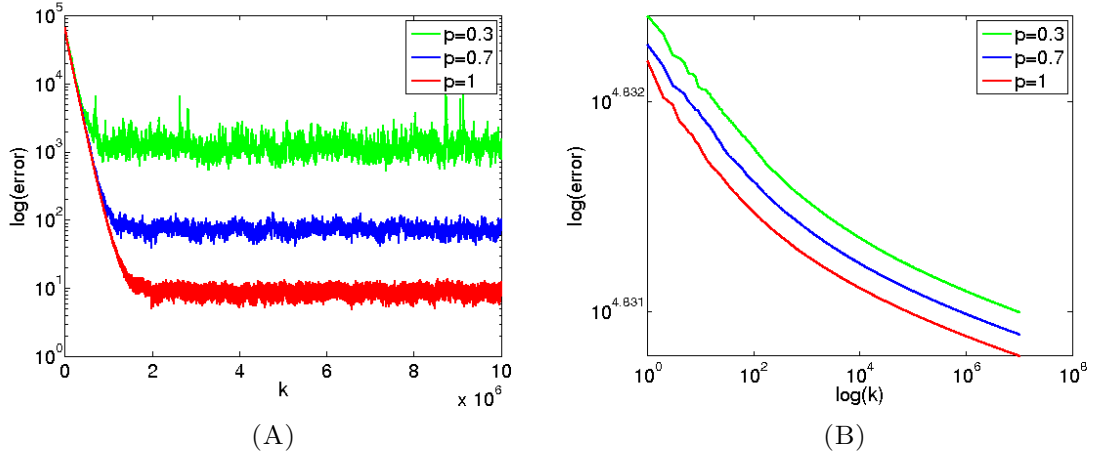


FIGURE 3. For bike data set, mSGD with a fixed step size (left) experiences a convergence horizon. Using updating step sizes (right), mSGD continues to progress toward the least squares solution.

step sizes in this case are too large, causing the error to increase initially. Regardless, we continue to use  $c = 10^{-2}$  here to be consistent with all other experiments.

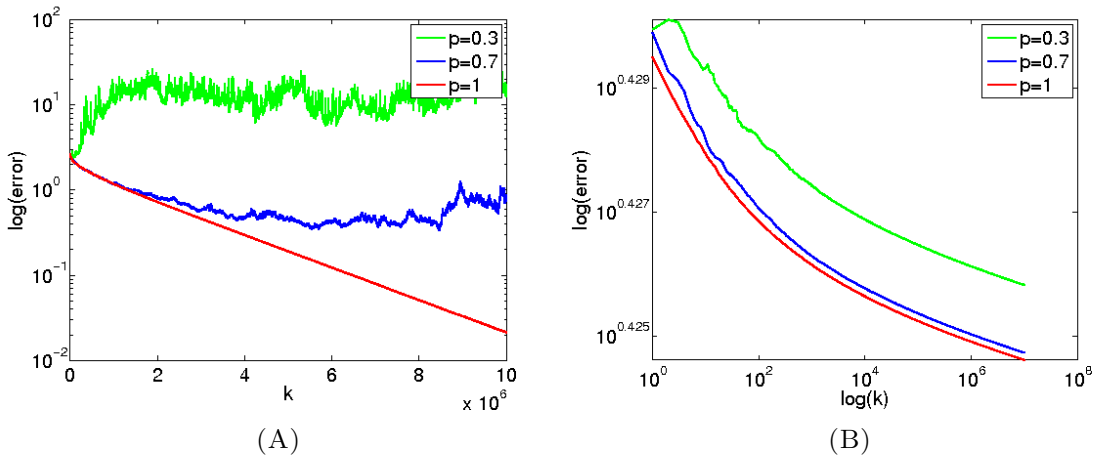


FIGURE 4. Using a fixed step size (left), the amount of missing data affects the convergence horizon. The smaller  $p$  is, or the more data that is missing, the larger the convergence radius. Using updating step sizes (right), there is a continual decrease in  $\ell_2$ -error for different percentages of available data.

Figure 5 compares mSGD and classical SGD applied to three different imputation treatments for missing data. For this experiment,  $\tilde{\mathbf{A}} \in \mathbb{R}^{100000 \times 25}$ . For each row  $\tilde{\mathbf{A}}_i$ , a random row  $\mathbf{A}_k$  of

the 2009 Adolescent California Health Interview Survey data is selected. For each element in  $\mathbf{A}_k$ , a Bernoulli random variable with is drawn to determine whether the element is missing or not. In this experiment, the probability an element is missing is  $p = 0.5$ . The right hand side vector  $\mathbf{y} \in \mathbb{R}^{100000 \times 1}$  is such that the  $i^{th}$  element is equal to  $\mathbf{y}_k$ , the  $k^{th}$  element of the right hand side vector in the 2009 Adolescent California Health Interview Survey linear system. Classical SGD is applied to  $\tilde{\mathbf{A}}$  in three ways: imputing 0 (if  $\tilde{\mathbf{A}}_{i,j}$  is missing,  $\tilde{\mathbf{A}}_{i,j} = 0$ ), imputing row means (if  $\tilde{\mathbf{A}}_{i,j}$  is missing,  $\tilde{\mathbf{A}}_{i,j}$  is the average over all non-missing elements in  $\tilde{\mathbf{A}}_i$ ), and imputing column means (if  $\tilde{\mathbf{A}}_{i,j}$  is missing,  $\tilde{\mathbf{A}}_{i,j}$  is the average over all non-missing elements in the  $j^{th}$  column of  $\tilde{\mathbf{A}}$ ). Notice that mSGD outperforms the imputation methods presented here, as expected.

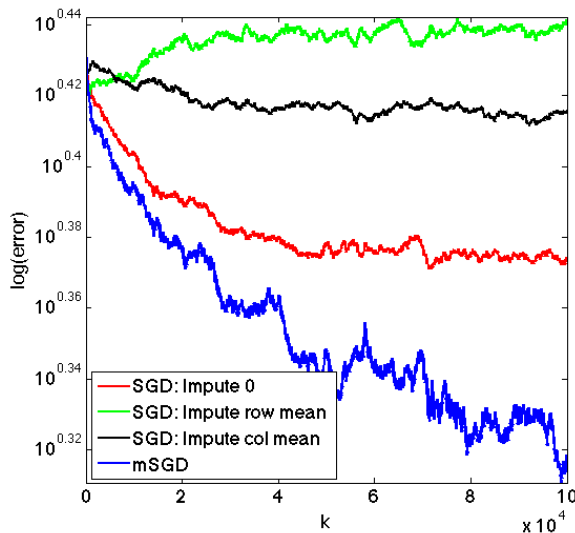


FIGURE 5

These experimental results support the theoretical findings presented in Section 2. Using a fixed step size, mSGD converges to some radius around the solution while using updating step size allows us to avoid the convergence horizon at the price of a slower convergence. For fixed step size, the amount of missing data affects the convergence horizon. In particular, as  $p$  decreases the size of the convergence horizon increases.

#### 4. CONCLUSION

In this work, we present a stochastic iterative projection method that solves linear systems with missing data. We prove that mSGD finds the least squares solution to the linear system with full data even though a system has missing data. Additionally, this work shows theoretical bounds for mSGD's performance using fixed and updating step sizes. The experiments show that mSGD is useful in real world settings when one wishes to solve a linear system with missing data without needing to impute missing values, which can be extremely costly.

## 5. APPENDIX

Consider the objective functions  $F(\mathbf{x}) = \frac{1}{2m}\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 = \frac{1}{m}\sum_{i=1}^m \frac{1}{2}(\mathbf{A}_i\mathbf{x} - \mathbf{b}_i)^2$  and  $\tilde{F}(\mathbf{x}) = \frac{1}{2m}\|\tilde{\mathbf{A}}\mathbf{x} - \mathbf{b}\|^2$ . Let  $\tilde{f}_i(\mathbf{x}) = \frac{1}{2}(\tilde{\mathbf{A}}_i\mathbf{x} - \mathbf{b}_i)^2$ . Let  $\mathbb{E}_\delta[\cdot]$  denote the expected value function with respect to the Bernoulli random variables of the binary mask  $\mathbf{D}$  and  $\mathbb{E}_i[\cdot]$  denote the expected value with respect to the choice of rows of  $\tilde{\mathbf{A}}$ . In addition, let  $\mu$  be the strong convexity parameter  $F(\mathbf{x})$  so that for any  $\mathbf{x}, \mathbf{y} \in \mathcal{W}$ ,  $(\mathbf{x} - \mathbf{y})^*(\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})) \geq \|\mathbf{x} - \mathbf{y}\|^2\mu$ .

First, we will show a few useful properties pertaining to the update function  $g(\mathbf{x})$ . In particular, Lemma 5.1 shows that in expectation  $g(\mathbf{x})$  allows us to make progress in the gradient direction of the objective  $F(\mathbf{x})$  (as opposed to the direction of  $\nabla \tilde{F}(\mathbf{x})$  in expectation). Next, Lemma 5.2 investigates the Lipschitz continuity of  $g(\mathbf{x})$  for a fixed row  $i$  and binary mask  $\mathbf{D}$  and its supremum Lipschitz constant of  $g(\mathbf{x})$  over all rows and binary masks. Lemma 5.3 shows that we can uniformly bound the expected norm of  $g(\mathbf{x})$  and provides said bound. Finally, we prove Theorem 2.3.

**Lemma 5.1.** *The expected value of the update function  $g(\mathbf{x})$  defined in (2.2) is the gradient of the objection function  $F(\mathbf{x})$ . In other words,  $\mathbb{E}g(\mathbf{x}) = \nabla F(\mathbf{x})$ .*

*Proof.* To prove this lemma, we will first take the expected value of  $\nabla \tilde{f}_i(\mathbf{x})$ . We take the expected value of  $g(\mathbf{x})$ , substitute  $\mathbb{E}[\nabla \tilde{f}_i(\mathbf{x})]$ , and simplify to complete the proof. Let's first check that

$$\mathbb{E}[\nabla \tilde{f}_i(\mathbf{x})] = \mathbb{E}[\tilde{\mathbf{A}}_i^*(\tilde{\mathbf{A}}_i\mathbf{x} - p\mathbf{b}_i)] = p^2\mathbf{A}^*\mathbf{A}\mathbf{x} + (p - p^2)\text{diag}(\mathbf{A}^*\mathbf{A})\mathbf{x} - p^2\sum_i \mathbf{A}_i\mathbf{b}_i. \quad (5.1)$$

Taking a simple derivative, we have that  $\nabla \tilde{f}_i(\mathbf{x}) = \tilde{\mathbf{A}}_i^*(\tilde{\mathbf{A}}_i\mathbf{x} - \mathbf{b}_i)$ . Recall that  $\mathbf{D}$  is a  $m \times n$  binary mask with entries  $\delta_{i,j} \stackrel{i.i.d.}{\sim} \text{Bern}(p)$ . Let  $\mathbf{D}_i = \text{diag}(\delta_{i,1}, \delta_{i,2}, \dots, \delta_{i,n})$  be a  $n \times n$  diagonal matrix so that  $\tilde{\mathbf{A}}_i = \mathbf{D}_i\mathbf{A}_i$ . Substituting  $\tilde{\mathbf{A}}_i = \mathbf{D}_i\mathbf{A}_i$  and taking the expectation with respect to the  $\delta_{i,j}$ 's,

$$\begin{aligned} \mathbb{E}_\delta[\nabla \tilde{f}_i(\mathbf{x})] &= \mathbb{E}_\delta[\tilde{\mathbf{A}}_i^*(\tilde{\mathbf{A}}_i\mathbf{x} - p\mathbf{b}_i)] \\ &= \mathbb{E}_\delta[\tilde{\mathbf{A}}_i^*\tilde{\mathbf{A}}_i\mathbf{x} - p\mathbb{E}_\delta[\tilde{\mathbf{A}}_i^*]\mathbf{b}_i] \\ &= \mathbb{E}_\delta[\tilde{\mathbf{A}}_i^*\tilde{\mathbf{A}}_i\mathbf{x} - p^2\mathbf{A}_i\mathbf{b}_i] \\ &\stackrel{(i)}{=} p^2\mathbf{A}_i^*\mathbf{A}_i\mathbf{x} + (p - p^2)\text{diag}(\mathbf{A}_i^*\mathbf{A}_i)\mathbf{x} - p^2\mathbf{A}_i\mathbf{b}_i. \end{aligned}$$

Letting  $[\mathbf{A}_i^*\mathbf{A}_i]_{jk}$  denote the  $(j, k)^{\text{th}}$  element of  $\mathbf{A}_i^*\mathbf{A}_i$ , step (i) uses the fact that,

$$\mathbb{E}_\delta[\tilde{\mathbf{A}}_i^*\tilde{\mathbf{A}}_i] = \begin{cases} p[\mathbf{A}_i^*\mathbf{A}_i]_{jk}, & j = k \\ p^2[\mathbf{A}_i^*\mathbf{A}_i]_{jk}, & j \neq k \end{cases}.$$

Now, we take the expectation with respect to the rows of  $\mathbf{A}$  to obtain:

$$\begin{aligned}
\mathbb{E}[g(\mathbf{x})] &\stackrel{(i)}{=} \frac{1}{p^2} \mathbb{E}[\tilde{\mathbf{A}}_i^* (\tilde{\mathbf{A}}_i \mathbf{x} - p \mathbf{b}_i)] - \frac{(1-p)}{p^2} \mathbb{E}[\text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i)] \mathbf{x} \\
&= \frac{1}{p^2} \mathbb{E}[\nabla \tilde{f}_i(\mathbf{x})] - \frac{(1-p)}{p^2} \mathbb{E}[\text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i)] \mathbf{x} \\
&\stackrel{(ii)}{=} \frac{1}{mp^2} \left( p^2 \mathbf{A}^* \mathbf{A} \mathbf{x} + (p-p^2) \text{diag}(\mathbf{A}^* \mathbf{A}) \mathbf{x} - p^2 \sum_i \mathbf{A}_i \mathbf{b}_i \right) - \frac{p(1-p)}{mp^2} \text{diag}(\mathbf{A}^* \mathbf{A}) \mathbf{x} \\
&= \frac{1}{m} \mathbf{A}^* \mathbf{A} \mathbf{x} + \frac{(p-p^2)}{mp^2} \text{diag}(\mathbf{A}^* \mathbf{A}) \mathbf{x} - \frac{1}{m} \sum_i \mathbf{A}_i \mathbf{b}_i - \frac{(p-p^2)}{mp^2} \text{diag}(\mathbf{A}^* \mathbf{A}) \mathbf{x} \\
&= \frac{1}{m} \left( \mathbf{A}^* \mathbf{A} \mathbf{x} - \sum_i \mathbf{A}_i \mathbf{b}_i \right) = \nabla F(\mathbf{x}).
\end{aligned}$$

Step (i) follows from the definition of  $g(\mathbf{x})$  and linearity of the expected value. Step (ii) utilizes (5.1) for the first expected value and evaluates the expectation of  $\mathbb{E}[\text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i)] = \mathbb{E}_i[\mathbb{E}_\delta[\text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i)]] = p \mathbb{E}_i[\text{diag}(\mathbf{A}_i^* \mathbf{A}_i)] = \frac{p}{m} \text{diag}(\mathbf{A}^* \mathbf{A})$ . The remaining steps follow by simplification.  $\square$

**Lemma 5.2.** *The update function  $g(\mathbf{x})$  of Algorithm 1 is Lipschitz continuous with Lipschitz constant  $L_{i,D}$ . In other words, for all  $\mathbf{x}, \mathbf{y} \in \mathcal{W}$ ,*

$$\|g(\mathbf{x}) - g(\mathbf{y})\| \leq L_{i,D} \|\mathbf{x} - \mathbf{y}\|.$$

In addition, we can bound the supremum Lipschitz constant,  $L_g$  by

$$L_g = \sup_{i,D} L_{g,i,D} \leq \frac{(2-p)}{p^2} a_{max}^2,$$

where  $a_{max}^2 = \max_i \|\mathbf{A}_i\|^2$ .

*Proof.* First we show that the Lipschitz constant  $L_{i,D}$  of  $g(\mathbf{x})$

$$\begin{aligned}
\|g(\mathbf{x}) - g(\mathbf{y})\| &= \left\| \left( \frac{1}{p^2} \tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i - \frac{(1-p)}{p^2} \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \right) (\mathbf{x} - \mathbf{y}) \right\| \\
&\leq \left\| \frac{1}{p^2} \tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i - \frac{(1-p)}{p^2} \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \right\| \|\mathbf{x} - \mathbf{y}\|.
\end{aligned}$$

Therefore we conclude that the Lipschitz constant of  $g(\mathbf{x})$  is  $L_{i,D} = \frac{1}{p^2} \left\| \tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i - (1-p) \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \right\|$ . Note that in determining supremum Lipschitz constant, we take the supremum over all possible rows and all possible binary masks. We bound the supremum Lipschitz constant,  $L_g$ , in the following

way:

$$\begin{aligned}
L_g &= \sup_{i,D} L_{i,D} = \sup_{i,D} \frac{1}{p^2} \|\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i - (1-p) \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i)\| \\
&\stackrel{(i)}{\leq} \sup_{i,D} \frac{1}{p^2} \left( \|\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i\| + (1-p) \|\text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i)\| \right) \\
&\leq \sup_{i,D} \frac{1}{p^2} \left( \|\tilde{\mathbf{A}}_i\|^2 + (1-p) \|\tilde{\mathbf{A}}_i\|^2 \right) \\
&\stackrel{(ii)}{\leq} \sup_i \frac{1}{p^2} (\|\mathbf{A}_i\|^2 + (1-p) \|\mathbf{A}_i\|^2) \\
&= \sup_i \frac{(2-p) \|\mathbf{A}_i\|^2}{p^2} \\
&\stackrel{(iii)}{\leq} \frac{2-p}{p^2} a_{max}^2.
\end{aligned}$$

Step (i) is an application of the triangle inequality. Step (ii) takes the supremum over all possible binary masks. Note since  $\mathbf{A}_i$  is a vector, the supremum is attained when  $\mathbf{D} = \mathbb{1}_{n \times n}$ , an  $n \times n$  matrix of ones. Finally, step (iii) takes the supremum over all possible rows.  $\square$

**Lemma 5.3.** *We can uniformly bound the expected value of the magnitude of the update function in the following way. We have that  $E\|g(\mathbf{x})\|^2 \leq G$ , where*

$$G = \frac{2B(2+p)(1-p)}{mp^3} \sum_i \|\mathbf{A}_i\|^4 + \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 |\mathbf{b}_i|^2,$$

where  $B = \max_{\mathbf{x} \in \mathcal{W}} \|\mathbf{x}\|^2$ . In addition, we have that

- if  $\mathbf{A}\mathbf{x}_* = \mathbf{b}$  (the linear system is consistent) then

$$G_* = \frac{2(1-p)(2-p)}{mp^3} \|\mathbf{x}_*\|^2 \sum_i \|\mathbf{A}_i\|^4.$$

- if  $\mathbf{A}\mathbf{x}_* = \mathbf{b} + \mathbf{r}$  (the linear system is inconsistent) then

$$G_* = \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 \mathbf{r}_i^2 + \frac{2(1-p)(2-p)}{mp^3} \|\mathbf{x}_*\|^2 \sum_i \|\mathbf{A}_i\|^4.$$

$G$  and  $G_*$  are also defined in (2.6) and (2.7) respectively.

*Proof.* We begin this proof by showing the upper bound of  $E[\|g(\mathbf{x})\|^2]$  for a general  $\mathbf{x}$ . From here, we obtain  $G_*$  by substituting  $\mathbf{x}$  with  $\mathbf{x}_*$  and making the appropriate assumptions on the consistency of the linear system. To get the uniform upper bound over all  $\mathbf{x}$ , we isolate  $\|\mathbf{x}\|^2$  and bound the norm by  $B = \max_{\mathbf{x} \in \mathcal{W}} \|\mathbf{x}\|^2$ . We have,

$$\begin{aligned}
\mathbb{E}[\|g(\mathbf{x})\|^2] &= \mathbb{E} \left[ \left\| \frac{1}{p^2} \left( \tilde{\mathbf{A}}_i^* (\tilde{\mathbf{A}}_i \mathbf{x} - p \mathbf{b}_i) \right) - \frac{(1-p)}{p^2} \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \mathbf{x} \right\|^2 \right] \\
&\stackrel{(i)}{\leq} \frac{2}{p^4} \mathbb{E} \left[ \left\| \tilde{\mathbf{A}}_i^* (\tilde{\mathbf{A}}_i \mathbf{x} - p \mathbf{b}_i) \right\|^2 \right] + \frac{2(1-p)^2}{p^4} \mathbb{E} \left[ \left\| \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \mathbf{x} \right\|^2 \right] \\
&\stackrel{(ii)}{=} \frac{2}{m^2 p^4} \mathbb{E} \left[ \|\tilde{\mathbf{A}}_i\|^2 (\tilde{\mathbf{A}}_i \mathbf{x} - p \mathbf{b}_i)^2 \right] + \frac{2(1-p)^2}{p^4} \mathbb{E} \left[ \left\| \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \mathbf{x} \right\|^2 \right] \\
&\stackrel{(iii)}{\leq} \frac{2}{p^4} \mathbb{E} \left[ \|\mathbf{A}_i\|^2 (\tilde{\mathbf{A}}_i \mathbf{x} - p \mathbf{b}_i)^2 \right] + \frac{2(1-p)^2}{p^4} \mathbb{E} \left[ \left\| \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \mathbf{x} \right\|^2 \right].
\end{aligned}$$

Step (i) follows by Jensen's inequality, step (ii) is simplification and uses the fact that  $(\tilde{\mathbf{A}}_i \mathbf{x} - p \mathbf{b}_i)$  is scalar. Lastly, step (iii) bounds the magnitude of a row of  $\mathbf{A}$  with missing data by the magnitude of a row of  $\mathbf{A}$  without missing data (i.e.  $\|\tilde{\mathbf{A}}_i\| = \|\mathbf{D}_i \mathbf{A}_i\| \leq \|\mathbf{A}_i\|$  for all  $\mathbf{D}_i$ ). From here, we use the fact that  $\mathbb{E} = \mathbb{E}_i \mathbb{E}_\delta$  to obtain the following:

$$\mathbb{E}[\|g(\mathbf{x})\|^2] \leq \frac{2}{p^4} \mathbb{E}_i \left[ \|\mathbf{A}_i\|^2 \underbrace{\mathbb{E}_\delta[(\tilde{\mathbf{A}}_i \mathbf{x} - p \mathbf{b}_i)^2]}_{(A)} \right] + \frac{2(1-p)^2}{p^4} \mathbb{E}_i \left[ \underbrace{\mathbb{E}_\delta \left[ \left\| \text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \mathbf{x} \right\|^2 \right]}_{(B)} \right]. \quad (5.2)$$

Now, we will focus on the computation of  $\mathbb{E}_\delta$ . First, we compute (A). We have that,

$$\begin{aligned}
\mathbb{E}_\delta \left[ (\tilde{\mathbf{A}}_i \mathbf{x} - p \mathbf{b}_i)^2 \right] &= \mathbb{E}_\delta \left[ (\tilde{\mathbf{A}}_i \mathbf{x})^2 \right] - 2p \mathbb{E}_\delta \left[ \tilde{\mathbf{A}}_i \right] \mathbf{x} \mathbf{b}_i + p^2 \mathbf{b}_i^2 \\
&= \mathbb{E}_\delta \left[ \left( \sum_{j=1}^n \tilde{\mathbf{A}}_{ij} \mathbf{x}_j \right)^2 \right] - 2p^2 \mathbf{A}_i \mathbf{x} \mathbf{b}_i + p^2 \mathbf{b}_i^2 \\
&= \mathbb{E}_\delta \left[ \sum_{j=1}^n \tilde{\mathbf{A}}_{ij}^2 \mathbf{x}_j^2 + 2 \sum_{j=1}^n \sum_{k=1}^{j-1} \tilde{\mathbf{A}}_{ij} \tilde{\mathbf{A}}_{ik} \mathbf{x}_j \mathbf{x}_k \right] - 2p^2 \mathbf{A}_i \mathbf{x} \mathbf{b}_i + p^2 \mathbf{b}_i^2 \\
&= \left( p \sum_{j=1}^n \mathbf{A}_{ij}^2 \mathbf{x}_j^2 + 2p^2 \sum_{j=1}^n \sum_{k=1}^{j-1} \mathbf{A}_{ij} \mathbf{A}_{ik} \mathbf{x}_j \mathbf{x}_k \right) - 2p^2 \mathbf{A}_i \mathbf{x} \mathbf{b}_i + p^2 \mathbf{b}_i^2 \\
&\stackrel{(i)}{=} \left( p^2 \sum_{j=1}^n \mathbf{A}_{ij}^2 \mathbf{x}_j^2 + (p-p^2) \sum_{j=1}^n \mathbf{A}_{ij}^2 \mathbf{x}_j^2 + 2p^2 \sum_{j=1}^n \sum_{k=1}^{j-1} \mathbf{A}_{i,j} \mathbf{A}_{i,k} \mathbf{x}_j \mathbf{x}_k \right) - 2p^2 \mathbf{A}_i \mathbf{x} \mathbf{b}_i + p^2 \mathbf{b}_i^2 \\
&= p^2 \left( \sum_{j=1}^n \mathbf{A}_{ij}^2 \mathbf{x}_j^2 + 2 \sum_{j=1}^n \sum_{k=1}^{j-1} \mathbf{A}_{ij} \mathbf{A}_{ik} \mathbf{x}_j \mathbf{x}_k - 2 \mathbf{A}_i \mathbf{x} \mathbf{b}_i + \mathbf{b}_i^2 \right) + (p-p^2) \left( \sum_{j=1}^n \mathbf{A}_{ij}^2 \mathbf{x}_j^2 \right) \\
&= p^2 \left( \left( \sum_{j=1}^n \mathbf{A}_{ij} \mathbf{x}_j \right)^2 - 2 \mathbf{A}_i \mathbf{x} \mathbf{b}_i + \mathbf{b}_i^2 \right) + (p-p^2) \mathbf{x}^* \text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x} \\
&= p^2 (\mathbf{A}_i \mathbf{x} - \mathbf{b}_i)^2 + p(1-p) \mathbf{x}^* \text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x}.
\end{aligned}$$

In step (i), we add and subtract the term  $p^2 \sum_{j=1}^n \mathbf{A}_{i,j}^2 \mathbf{x}_j^2$  so that we can combine terms. Other equalities follow by simplification and computation of expected value. Note that  $\mathbb{E}_\delta[\tilde{\mathbf{A}}_{i,j}] = p \mathbf{A}_{i,j}$  and  $\mathbb{E}_\delta[\tilde{\mathbf{A}}_{i,j} \tilde{\mathbf{A}}_{i,k}] = p^2 \mathbf{A}_{i,j} \mathbf{A}_{i,k}$  if  $j \neq k$ .

For term (B), we simply compute that

$$\begin{aligned}
\mathbb{E}_\delta \left[ \|\text{diag}(\tilde{\mathbf{A}}_i^* \tilde{\mathbf{A}}_i) \mathbf{x}\|^2 \right] &= \mathbb{E}_\delta \left[ \sum_{j=1}^n \tilde{\mathbf{A}}_{i,j}^2 \mathbf{x}_j^2 \right] \\
&= p \sum_{j=1}^n (\mathbf{A}_{i,j}^2 \mathbf{x}_j^2) \\
&= p \|\text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x}\|^2.
\end{aligned}$$

Now that we have (A) and (B), we can compute a general upper bound for  $\mathbb{E}[\|g(\mathbf{x})\|^2]$ . Starting with substituting (A) and (B) into (5.2),

$$\begin{aligned}
\mathbb{E} [\|g(\mathbf{x})\|^2] &\stackrel{(i)}{\leq} \frac{2}{p^2} \mathbb{E}_i \left[ \|\mathbf{A}_i\|^2 (\mathbf{A}_i \mathbf{x} - \mathbf{b}_i)^2 \right] + \frac{2p(1-p)}{p^4} \mathbb{E}_i \left[ \|\mathbf{A}_i\|^2 \mathbf{x}^* \text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x} \right] \\
&\quad + \frac{2p(1-p)^2}{p^4} \mathbb{E}_i \left[ \|\text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x}\|^2 \right] \\
&\stackrel{(ii)}{\leq} \frac{2}{p^2} \mathbb{E}_i \left[ \|\mathbf{A}_i\|^2 (\mathbf{A}_i \mathbf{x} - \mathbf{b}_i)^2 \right] + \left( \frac{2p(1-p)}{p^4} + \frac{2p(1-p)^2}{p^4} \right) \mathbb{E}_i \left[ \|\mathbf{A}_i\|^2 \mathbf{x}^* \text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x} \right] \\
&\leq \frac{2}{p^2} \mathbb{E}_i \left[ \|\mathbf{A}_i\|^2 (\mathbf{A}_i \mathbf{x} - \mathbf{b}_i)^2 \right] + \frac{2p(1-p)(2-p)}{p^4} \mathbb{E}_i \left[ \|\mathbf{A}_i\|^2 \mathbf{x}^* \text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x} \right] \\
&\stackrel{(iii)}{=} \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 (\mathbf{A}_i \mathbf{x} - \mathbf{b}_i)^2 + \frac{2p(1-p)(2-p)}{mp^4} \sum_i \|\mathbf{A}_i\|^2 \mathbf{x}^* \text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x} \\
&\stackrel{(iv)}{\leq} \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 (\mathbf{A}_i \mathbf{x} - \mathbf{b}_i)^2 + \frac{2p(1-p)(2-p)}{mp^4} \|\mathbf{x}\|^2 \sum_i \|\mathbf{A}_i\|^4
\end{aligned}$$

Step (i) substitutes (A) and (B) in (5.2). Step (ii) uses the fact that:

$$\begin{aligned}
\|\text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x}\|^2 &= \|\text{diag}(\mathbf{A}_i) \text{diag}(\mathbf{A}_i) \mathbf{x}\|^2 \\
&\leq \|\text{diag}(\mathbf{A}_i)\|^2 \|\text{diag}(\mathbf{A}_i) \mathbf{x}\|^2 \\
&\leq \|\text{diag}(\mathbf{A}_i)\|_F^2 \|\text{diag}(\mathbf{A}_i) \mathbf{x}\|^2 \\
&= \|\mathbf{A}_i\|^2 \|\text{diag}(\mathbf{A}_i) \mathbf{x}\|^2 \\
&= \|\mathbf{A}_i\|^2 \mathbf{x}^* \text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x}.
\end{aligned}$$

Here,  $\|\cdot\|_F$  denotes the Frobenious norm. In step (iii) we take the expected value over the choice of rows and then simplify. Finally, in step (iv) we simplify

$$\begin{aligned}
\sum_i \|\mathbf{A}_i\|^2 \mathbf{x}^* \text{diag}(\mathbf{A}_i^* \mathbf{A}_i) \mathbf{x} &\leq \sum_i \|\mathbf{A}_i\|^2 \|\mathbf{A}_i \mathbf{x}\|^2 \\
&\leq \sum_i \|\mathbf{A}_i\|^2 \|\mathbf{A}_i\|^2 \|\mathbf{x}\|^2 \\
&= \|\mathbf{x}\|^2 \sum_i \|\mathbf{A}_i\|^4.
\end{aligned}$$

From here, we substitute  $\mathbf{x}$  with  $\mathbf{x}_\star$  to compute  $G_\star$ . If  $\mathbf{A} \mathbf{x}_\star = \mathbf{b}$  (the linear system is consistent) then the terms  $(\mathbf{A}_i \mathbf{x} - \mathbf{b}_i)^2 = 0$  and we find that

$$G_\star = \frac{2(1-p)(2-p)}{mp^3} \|\mathbf{x}_\star\|^2 \sum_i \|\mathbf{A}_i\|^4.$$

Otherwise, if  $\mathbf{A} \mathbf{x}_\star = \mathbf{b} + \mathbf{r}$  for some residual vector  $\mathbf{r}$ , we have that

$$G_\star = \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 \mathbf{r}_i^2 + \frac{2(1-p)(2-p)}{mp^3} \|\mathbf{x}_\star\|^2 \sum_i \|\mathbf{A}_i\|^4,$$



where  $\mathbf{r}_i$  is the  $i^{\text{th}}$  element of the vector  $\mathbf{r}$ . To finish the proof of Lemma 5.3, we simplify starting from step (iv).

$$\begin{aligned}
\mathbb{E} [\|g(\mathbf{x})\|^2] &\leq \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 (\mathbf{A}_i \mathbf{x} - \mathbf{b}_i)^2 + \frac{2p(1-p)(2-p)}{mp^4} \|\mathbf{x}\|^2 \sum_i \|\mathbf{A}_i\|^4 \\
&\stackrel{(i)}{\leq} \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 (|\mathbf{A}_i \mathbf{x}|^2 + |\mathbf{b}_i|^2) + \frac{2p(1-p)(2-p)}{mp^4} \|\mathbf{x}\|^2 \sum_i \|\mathbf{A}_i\|^4 \\
&\stackrel{(ii)}{\leq} \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^4 \|\mathbf{x}\|^2 + \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 |\mathbf{b}_i|^2 + \frac{2p(1-p)(2-p)}{mp^4} \|\mathbf{x}\|^2 \sum_i \|\mathbf{A}_i\|^4 \\
&\stackrel{(iii)}{\leq} \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^4 B + \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 |\mathbf{b}_i|^2 + \frac{2p(1-p)(2-p)}{mp^4} B \sum_i \|\mathbf{A}_i\|^4 \\
&= \left( \frac{2B}{mp^2} + \frac{2p(1-p)(2-p)B}{mp^4} \right) \sum_i \|\mathbf{A}_i\|^4 + \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 |\mathbf{b}_i|^2 \\
&= \frac{2B}{mp^2} \left( 1 + \frac{p(1-p)(2-p)}{p^2} \right) \sum_i \|\mathbf{A}_i\|^4 + \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 |\mathbf{b}_i|^2 \\
&= \left( \frac{2B(2+p)(1-p)}{mp^3} \right) \sum_i \|\mathbf{A}_i\|^4 + \frac{2}{mp^2} \sum_i \|\mathbf{A}_i\|^2 |\mathbf{b}_i|^2
\end{aligned}$$

In step (i) we use Jensen's inequality. Note that  $\mathbf{A}_i \mathbf{x}$  and  $\mathbf{b}_i$  are both scalar values. In step (ii), we distribution the summation in the first term and use the fact that  $|\mathbf{A}_i \mathbf{x}|^2 \leq \|\mathbf{A}_i\|^2 \|\mathbf{x}\|^2$  by the Cauchy-Schwartz inequality. Step (iii) uses the definition of  $B = \max_{\mathbf{x} \in \mathcal{W}} \|\mathbf{x}\|^2$ . The remaining lines are simplification.  $\square$

Before we begin the proof of Theorem 2.3, we remind the reader that  $F(\mathbf{x})$  is strongly convex with strong convexity parameter  $\mu$ . In other words, for all  $\mathbf{x}, \mathbf{y} \in \mathcal{W}$  we have that

$$(\mathbf{x} - \mathbf{y})^* (\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})) \geq \mu \|\mathbf{x} - \mathbf{y}\|^2. \quad (5.3)$$

In addition, we define a new function  $G(\mathbf{x}) = \frac{1}{2p^2} \left( (\tilde{\mathbf{A}}_i \mathbf{x} - p\mathbf{b}_i)^2 - \frac{(1-p)}{2p^2} \|\text{diag}(\tilde{\mathbf{A}}_i) \mathbf{x}\|^2 \right)$  so that  $g(\mathbf{x}) = \nabla G(\mathbf{x})$ . The update function  $g(\mathbf{x})$  follows the Co-coercivity Lemma as stated in Lemma 5.4.

**Lemma 5.4.** ([NWS14] Lemma A.1) For  $G(\mathbf{x})$  a smooth function such that  $\nabla G(\mathbf{x}) = g(\mathbf{x})$ ,

$$\|g(\mathbf{x}) - g(\mathbf{y})\|^2 \leq L_{i,D} (\mathbf{x} - \mathbf{y})^* (g(\mathbf{x}) - g(\mathbf{y})),$$

where  $g(\mathbf{x})$  has Lipschitz constant  $L_{i,D}$ .

### 5.1. Proof of Theorem 2.3.

*Proof.* First, we bound expected error conditional on the previous  $k - 1$  iterations. Let  $\mathbb{E}_{k-1}[\cdot]$  denote the expected value conditional of the previous  $k - 1$  iterations and note that by the Law of

Iterated Expectation, we have that the full expected value over all iterations is  $\mathbb{E}[\cdot] = \mathbb{E}[\mathbb{E}_{k-1}[\cdot]]$ . Thus,

$$\begin{aligned} \mathbb{E}_{k-1} [\|\mathbf{x}_k - \mathbf{x}_\star\|^2] &= \mathbb{E}_{k-1} [\|\mathbf{x}_{k-1} - \alpha g(\mathbf{x}_{k-1}) - \mathbf{x}_\star\|^2] \\ &\stackrel{(i)}{=} \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* \mathbb{E}_{k-1}[g(\mathbf{x}_{k-1})] + \alpha^2 \mathbb{E}_{k-1} [\|g(\mathbf{x}_{k-1})\|^2] \\ &\stackrel{(ii)}{=} \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* \nabla F(\mathbf{x}_{k-1}) + \alpha^2 \mathbb{E}_{k-1} \|g(\mathbf{x}_{k-1})\|^2 \\ &\stackrel{(iii)}{=} \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* (\nabla F(\mathbf{x}_{k-1}) - \nabla F(\mathbf{x}_\star)) + \alpha^2 \mathbb{E}_{k-1} [\|g(\mathbf{x}_{k-1})\|^2] \end{aligned}$$

Step (i) follows by definition of the  $\ell_2$  norm, step (ii) by Lemma 5.1, and step (iii) since  $\nabla F(\mathbf{x}_\star) = 0$ . Continuing, we have,s

$$\begin{aligned} \mathbb{E}_{k-1} [\|\mathbf{x}_k - \mathbf{x}_\star\|^2] &\stackrel{(i)}{=} \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* \mathbb{E}_{k-1} [g(\mathbf{x}_{k-1})] + 2\alpha^2 \mathbb{E}_{k-1} [\|g(\mathbf{x}_{k-1})\|^2] \\ &\stackrel{(ii)}{=} \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* (\nabla F(\mathbf{x}_{k-1}) - \nabla F(\mathbf{x}_\star)) + 2\alpha^2 \mathbb{E}_{k-1} [\|g(\mathbf{x}_{k-1})\|^2] \\ &\stackrel{(iii)}{\leq} \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* (\nabla F(\mathbf{x}_{k-1}) - \nabla F(\mathbf{x}_\star)) \\ &\quad + 2\alpha^2 \mathbb{E}_{k-1} [\|g(\mathbf{x}_{k-1}) - g(\mathbf{x}_\star)\|^2] + 2\alpha^2 \mathbb{E}_{k-1} [\|g(\mathbf{x}_\star)\|^2] \\ &\stackrel{(iv)}{\leq} \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* (\nabla F(\mathbf{x}_{k-1}) - \nabla F(\mathbf{x}_\star)) \\ &\quad + 2\alpha^2 L_{i,g}(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* (\mathbb{E}_{k-1}[g(\mathbf{x}_{k-1})] - \mathbb{E}_{k-1}[g(\mathbf{x}_\star)]) + 2\alpha^2 G_\star \\ &\stackrel{(v)}{\leq} \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* (\nabla F(\mathbf{x}_{k-1}) - \nabla F(\mathbf{x}_\star)) \\ &\quad + 2\alpha^2 L_g(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* (\nabla F(\mathbf{x}_{k-1}) - \nabla F(\mathbf{x}_\star)) + 2\alpha^2 G_\star \\ &\leq \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(1 - \alpha L_g)(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* (\nabla F(\mathbf{x}_{k-1}) - \nabla F(\mathbf{x}_\star)) + 2\alpha^2 G_\star \\ &\stackrel{(vi)}{\leq} \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 - 2\alpha(1 - \alpha L_g)\mu \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 + 2\alpha^2 G_\star \\ &= (1 - 2\alpha\mu(1 - \alpha L_g)) \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 + 2\alpha^2 G_\star \\ &= r \|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 + 2\alpha^2 G_\star. \end{aligned}$$

Step (i) follows from the definition of the  $\ell_2$  norm. Step (ii) takes the expectation of  $g(\mathbf{x}_{k-1})$  using Lemma 5.1 and uses the fact that  $\nabla F(\mathbf{x}_\star) = 0$  to subtract  $2\alpha(\mathbf{x}_{k-1} - \mathbf{x}_\star)^* \nabla F(\mathbf{x}_\star)$ . In step (iii) we add and subtract the term  $\|g(\mathbf{x}_\star)\|^2$  then apply Jensen's inequality. Step (iv) is an application of the Lemma 5.4. Step (v) bounds  $L_{i,D}$  by  $L_g = \sup_{i,D} L_{i,D}$  and uses Lemma 5.1 to compute the expectation of  $\mathbb{E}_{k-1}[g(\mathbf{x})]$ . We use the strong convexity of  $F(\mathbf{x})$  in step (vi). The remaining lines are simplification. Now, by the Law of Iterated Expectation we recursively apply this bound to obtain the desired result,

$$\begin{aligned}
\mathbb{E}\|\mathbf{x}_k - \mathbf{x}_\star\|^2 &\leq r\|\mathbf{x}_{k-1} - \mathbf{x}_\star\|^2 + 2\alpha^2 G_\star \\
&\leq r^k\|\mathbf{x}_0 - \mathbf{x}_\star\|^2 + 2\alpha^2 G_\star \sum_{j=0}^{k-1} r^j \\
&\leq r^k\|\mathbf{x}_0 - \mathbf{x}_\star\|^2 + \frac{2\alpha^2 G_\star}{1-r}.
\end{aligned}$$

**A note on Inconsistent Linear Systems.** Theorem 2.3 also applies to inconsistent systems. Let  $\mathbf{A}\mathbf{x}_\star = \mathbf{b} + \mathbf{r}$  where  $\mathbf{r} \in \text{null}(\mathbf{A}^*)$ . In the proof of Theorem 2.3, we use the fact that  $\nabla F(\mathbf{x}_\star) = 0$  in step (ii). This is still true in the inconsistent setting as  $\nabla F(\mathbf{x}_\star) = \mathbf{A}^*(\mathbf{A}\mathbf{x}_\star - \mathbf{b}) = \mathbf{A}^*\mathbf{r} = 0$ . All other computations go through without issue.  $\square$

#### 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.

#### REFERENCES

- [Bot10] L. Bottou. Large-scale machine learning with stochastic gradient descent. In *Proc. of COMPSTAT'2010*, pages 177–186. Springer, 2010.
- [Bot12] L. Bottou. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade*, pages 421–436. Springer, 2012.
- [CCS10] J.-F. Cai, E. J. Candès, and Z. Shen. A singular value thresholding algorithm for matrix completion. *SIAM J. Optimz.*, 20(4):1956–1982, 2010.
- [CEG83] Y. Censor, P. P. Eggermont, and D. Gordon. Strong underrelaxation in kaczmarz’s method for inconsistent systems. *Numer. Math.*, 41(1):83–92, 1983.
- [DLR77] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *J. Roy. Stat. Soc. B Met.*, pages 1–38, 1977.
- [Efr94] B. Efron. Missing data, imputation, and the bootstrap. *J. Am. Stat. Assoc.*, 89(426):463–475, 1994.
- [FC03] M. Fichman and J. N. Cummings. Multiple imputation for missing data: Making the most of what you know. *Organ. Res. Methods*, 6(3):282–308, 2003.
- [fHPR09] U. C. for Health Policy Research. California health interview survey, 2009.
- [HN90] M. Hanke and W. Niethammer. On the acceleration of kaczmarz’s method for inconsistent linear systems. *Linear Algebra Appl.*, 130:83–98, 1990.
- [Kac37] S. Kaczmarz. Angenäherte auflösung von systemen linearer gleichungen. *Bull. Int. Acad. Polon. Sci. Lett. Ser. A*, 35:355–357, 1937.
- [KMO10] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from noisy entries. *J. Machine Learning Research*, 11(Jul):2057–2078, 2010.
- [KOM09] R. H. Keshavan, S. Oh, and A. Montanari. Matrix completion from a few entries. In *IEEE T. Inform. Theory*, pages 324–328. IEEE, 2009.
- [Lic13] M. Lichman. UCI machine learning repository, 2013.
- [LL10] D. Leventhal and A. S. Lewis. Randomized methods for linear constraints: convergence rates and conditioning. *Math. Oper. Res.*, 35(3):641–654, 2010.
- [LR14] R. J. Little and D. B. Rubin. *Statistical analysis with missing data*. John Wiley & Sons, 2014.

- [MNR15] A. Ma, D. Needell, and A. Ramdas. Convergence properties of the randomized extended gauss–seidel and kaczmarz methods. *SIAM J. Matrix Anal. and Appl.*, 36(4):1590–1604, 2015.
- [NWS14] D. Needell, R. Ward, and N. Srebro. Stochastic gradient descent, weighted sampling, and the randomized kaczmarz algorithm. In *Proc. Adv. in Neural Processing Systems (NIPS)*, pages 1017–1025, 2014.
- [Rec11] B. Recht. A simpler approach to matrix completion. *J. Machine Learning Research*, 12(Dec):3413–3430, 2011.
- [RM51] H. Robbins and S. Monro. A stochastic approximation method. *Ann. Math. Statist.*, pages 400–407, 1951.
- [SR13] M. Schmidt and N. L. Roux. Fast convergence of stochastic gradient descent under a strong growth condition. *arXiv preprint arXiv:1308.6370*, 2013.
- [SV09] T. Strohmer and R. Vershynin. A randomized kaczmarz algorithm with exponential convergence. *J. Fourier Anal. Appl.*, 15(2):262–278, 2009.
- [SZ13] O. Shamir and T. Zhang. Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. In *Proc. Int. Conf. Machine Learning*, pages 71–79, 2013.
- [Zha04] T. Zhang. Solving large scale linear prediction problems using stochastic gradient descent algorithms. In *Proc. Int. Conf. Machine Learning*, page 116. ACM, 2004.