# Stochastic Gradient Descent for Linear Systems with Missing Data

Anna Ma[1,*] and Deanna Needell [2]

[1] *Claremont Graduate University, Claremont, CA 91711, USA*
[2] *University of California, Los Angeles, Los Angeles CA 90095, USA*

**Abstract.** Traditional methods for solving linear systems have quickly become impractical due to an increase in the size of available data. Utilizing massive amounts of data is further complicated when the data is incomplete or has missing entries. In this work, we address the obstacles presented when working with large 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 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

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 [13, 16, 18, 23]. 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 [6, 7], matrix completion [3, 11, 12, 19], and maximum likelihood estimation [5, 15] can be wasteful, create biases, or be impractical for extremely large amounts of data. This work simultaneously addresses both issues of large-scale and missing data.

---

*Corresponding author. *Email addresses:* `a4ma@ucsd.edu` (A. Ma), `deanna@math.ucla.edu` (D. Needell)

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

Let $\tilde{A} = D \circ A$ where $A$ denotes the full matrix, and $\circ$ be the element-wise product, $D$ denotes a binary matrix (1 indicating the availability of an element and 0 indicating a missing entry). Formally, one wants to solve the following optimization program:

$$\text{Given } \tilde{A}, b \quad \text{s.t. } Ax = b \text{ and } \tilde{A} = D \circ A,$$

$$\text{Find } x_\star = \arg\min_{x \in \mathscr{W}} \frac{1}{2m} \|Ax - b\|^2, \tag{1.1}$$

where $\mathscr{W}$ is a convex domain containing the solution $x_\star$ (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 Randomized Kaczmarz (RK) 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 [20], SGD has proved to be particularly popular in machine learning [1, 2, 24]. SGD minimizes an objective function $F(x)$ over a convex domain $\mathscr{W}$ using unbiased estimates for the gradient of the objective, i.e., using $f_i(x)$ such that $\mathbb{E}[\nabla f_i(x)] = \nabla F(x)$. At each iteration, a random unbiased estimate, $\nabla f_i(x)$, is drawn and the minimizer of $F(x)$ is estimated with:

$$x_k = \mathscr{P}_{\mathscr{W}}\left(x_{k-1} - \alpha_k \nabla f_i(x_{k-1})\right), \tag{1.2}$$

where $\alpha_k$ is an appropriately chosen step size, or learning rate, at iteration $k$ and $\mathscr{P}_{\mathscr{W}}$ denotes the projection onto the convex set $\mathscr{W}$. To solve an overdetermined linear system $Ax = b$, one approach is to minimize the least-squares objective function $F(x) = \frac{1}{2m}\|Ax - b\|^2 = \frac{1}{m}\sum_{i=1}^{m} f_i(x)$, where $f_i(x) = \frac{1}{2}(A_i x - b_i)^2$, $A_i$ denotes the $i^{th}$ row of $A$, and $b_i$ denotes the $i^{th}$ entry of $b$. In this setting, a random *row* of the matrix $A$ is selected and

---

[1]The linear system is not assumed to be consistent; we will use the notation $Ax = b$ to denote a general linear system.