

Mean Imputation and Stochastic Coordinate Descent for Linear Systems with Missing Data

Meha Patel^{*}, Samuel Rath[†], and Chupeng Zheng[‡]

Project Advisor: Anna Ma[§]

Abstract. As big data problems become more prevalent, the need to accurately approximate solutions to large-scale linear systems increases. Many real-world big data problems are also accompanied by the risk of missing or incomplete data, further complicating the linear models assigned to them. Current methods to address missing data involve deletion or zero imputation, which introduces bias to the model. We propose a model that adapts Stochastic Coordinate Descent (SCD) to handle missing data in linear systems and utilizes μ -imputation to retrieve a better approximation of the original data. We prove that in expectation, our proposed algorithm, μ -imputation mSCD utilizes an unbiased estimator of the gradient of the least-squares objective function when using mean imputation in the absence of data. Furthermore, we compare our algorithm's performance on synthetic data to closely related algorithms: zero-imputation mSCD and SCD. Finally, we apply μ -imputation mSCD on real-world data to demonstrate the usefulness and viability of our proposed algorithm.

1. Introduction. Large-scale data is crucial in training algorithms that are being implemented and utilized today. However, the data collection process to train such algorithms is often imperfect and can lead to noisy and incomplete data. For example, malfunctions in physical measurement devices can cause data to become corrupt or, in extreme instances, unavailable. As another example, a person may skip questions on a questionnaire to save time resulting in incomplete survey data. While a straightforward approach for dealing with missing data is to impute the missing data with zeros or ignore missing data altogether (throwing out any and all incomplete data points), this quickly becomes impractical and inefficient.

Other methods for solving missing data problems consider imputation methods in which a new data set is created using information from the available data set. For example, in [6], Mukhopadhyay and Mukherjee propose an algorithm for imputing incomplete streaming data using a constant factor times the possibly imputed data at the previous time instance. While this method eliminates any inefficiency of solving a matrix with missing entries and provides a more accurate estimation of the full matrix, it is not independent of data being missing at random. Mukhopadhyay and Mukherjee assume that there is a previous entry from which we can draw information to determine the following entry. When we face the problem of incomplete data sets, we often find that the data is missing independently at random. This work focuses on missing data and how it arises in linear systems. In particular, we consider the following linear system of equations $\mathbf{Ax} = \mathbf{y}$ where $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the measurement matrix that is not entirely known, i.e. only some of its entries are known, and we consider the over-determined case in which $m \geq n$, $\mathbf{x} \in \mathbb{R}^n$ is the unknown signal we wish to find, and $\mathbf{y} \in \mathbb{R}^m$ are the measurements. In most general cases, this system can be solved as $\mathbf{x} = \mathbf{A}^\dagger \mathbf{y}$,

^{*}California State University Long Beach, Long Beach, CA (meha.patel01@student.csulb.edu)

[†]San Diego State University, San Diego, CA (sarath@sdsu.edu)

[‡]University of Chicago, Chicago, IL (chupenz@uchicago.edu)

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

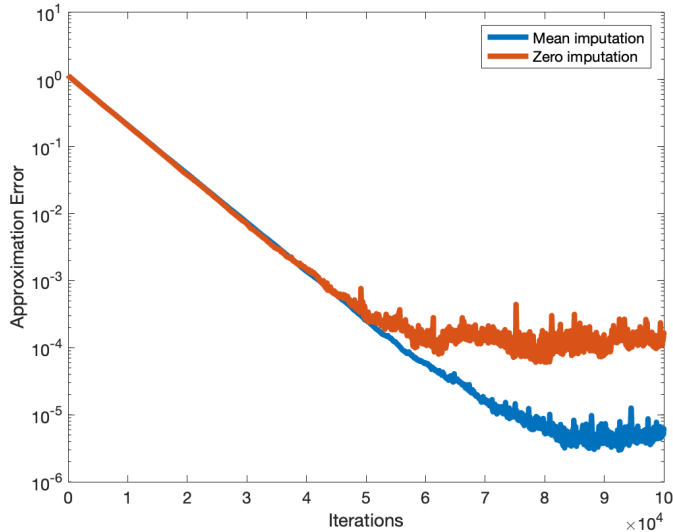


Figure 4. Performance of μ -imputed mSCD on Garment Productivity data set from the UCI Machine Learning Repo.

trials. We report the approximation error to the least squares solution for this data set using $\alpha = 5 \times 10^{-5}$ and $p = 0.85$. We observe that the approximation error decays linearly until the algorithm reaches a convergence horizon on the order of 10^{-5} for mean imputation and 10^{-4} for zero imputation, highlighting the benefit of using mean imputation.

5. Conclusion. We propose μ -imputed mSCD as a solution to big data problems in which missing or corrupt data issues arise. Our method utilizes an unbiased estimator of the gradient of the least-squares objective function as well as mean imputation to eliminate the issue of bias that may be introduced in the absence of data. The experimental results and theoretical proofs show that our method has less bias than current methods and, subsequently, performs better than these methods when solving large-scale linear systems with missing data. While we are able to show our algorithm’s convergence capabilities through theoretical and practical experiments, we hope to analyze convergence guarantees in future work. Furthermore, for future work we believe our method can be expanded to other types of imputation methods or patterns of missing data. For example, we believe our idea of using mean imputation can be utilized to improve the cumulative information method used in [5] to create a cumulative information mean-imputed algorithm.

Acknowledgments. We thank Dr. Anna Ma from the University of California, Irvine for advising us and this work. We also thank Chelsea Huynh and Michael Strand from the University of California, Irvine for their contributions in facilitating the conversations around our research topic.

Appendix. In this section we prove the major theorem of the paper, namely Theorem 3.3. To do this, we first introduce some convenient notation, π and Π , to simplify the calculation.

This notation will also be used for the proof of Theorem 3.4, and in fact one of our goals is that it conveys the essential parallels between these proofs. Finally, we also prove Theorem 3.5.

By the assumption of our model, we treat the entries of a given matrix $\tilde{\mathbf{A}}$ as i.i.d. Bernoulli random variables such that

$$\tilde{a}_{ij} = \begin{cases} a_{ij} & w.p. \ p \\ \mu & w.p. \ 1-p \end{cases},$$

for which the expected value is simply

$$(5.1) \quad \mathbb{E}_\delta [\tilde{a}_{ij}] = pa_{ij} + (1-p)\mu.$$

However, when computing expectations we often run into products of 2 entries. For these, we need to consider 2 cases:

$$\tilde{a}_{ij}\tilde{a}_{ik} = \begin{cases} a_{ij}a_{ik} & w.p. \ p^2 \\ \mu a_{ij} & w.p. \ p(1-p) \\ \mu a_{ik} & w.p. \ p(1-p) \\ \mu^2 & w.p. \ (1-p)^2 \end{cases},$$

when $j \neq k$, and

$$\tilde{a}_{ij}\tilde{a}_{ik} = \begin{cases} a_{ij}a_{ik} & w.p. \ p \\ \mu^2 & w.p. \ 1-p \end{cases},$$

when $j = k$. These come from the fact that when an entry is multiplied by itself, it still corresponds to one coin toss with 2 outcomes. On the other hand, when the entries are distinct, we now have 2 independent coin tosses and 4 outcomes total. Thus, we get the following expectation:

$$(5.2) \quad \mathbb{E}_\delta [\tilde{a}_{ij}\tilde{a}_{ik}] = \begin{cases} p^2 a_{ij}a_{ik} + p(1-p)(a_{ij}\mu + \mu a_{ik}) + (1-p)^2 \mu^2 & j \neq k \\ pa_{ij}a_{ik} + (1-p)\mu^2 & j = k \end{cases}.$$

The 2 cases in (5.2) occur so often in our proofs that we give them labels to condense the notation. These labels will also serve to indicate how the expected value operator $\mathbb{E}_\delta [\cdot]$ acts on different components in our calculations, depending on whether there are products of repeated or distinct entries. Let us define the functions $\Pi : \mathbb{R}^{m \times d} \times \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{m \times n}$ and $\pi : \mathbb{R}^{m \times d} \times \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{m \times n}$ such that

$$(5.3) \quad \Pi[\mathbf{A}, \mathbf{B}] := p^2 \mathbf{A}\mathbf{B} + p(1-p)(\mathbf{A}\mathbf{M}_B + \mathbf{M}_A\mathbf{B}) + (1-p)^2 \mathbf{M}_A\mathbf{M}_B$$

$$(5.4) \quad \pi[\mathbf{A}, \mathbf{B}] := p\mathbf{A}\mathbf{B} + (1-p)\mathbf{M}_A\mathbf{M}_B,$$

where $\mathbf{M}_A = \mu \mathbf{1}_{m \times d}$ and $\mathbf{M}_B = \mu \mathbf{1}_{d \times n}$, and whose dimensions are that of \mathbf{A} and \mathbf{B} respectively. Hence, using (5.3) and (5.4), we can now express (5.2) as

$$(5.5) \quad \mathbb{E}_\delta [\tilde{a}_{ij}\tilde{a}_{ik}] = \begin{cases} \Pi[a_{ij}, a_{ik}] & j \neq k \\ \pi[a_{ij}, a_{ik}] & j = k \end{cases}.$$

In addition to condensing the expectation, the functions $\Pi[\cdot, \cdot]$ and $\pi[\cdot, \cdot]$ have the added benefit of extending to higher dimensional expectations. For instance, suppose we need to take an expectation of the inner product of columns:

$$\begin{aligned}\mathbb{E}_\delta \left[\tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:k} \right] &= \mathbb{E}_\delta [\tilde{a}_{1j}\tilde{a}_{1k} + \tilde{a}_{2j}\tilde{a}_{2k} + \dots + \tilde{a}_{mj}\tilde{a}_{mk}] \\ &= \mathbb{E}_\delta [\tilde{a}_{1j}\tilde{a}_{1k}] + \mathbb{E}_\delta [\tilde{a}_{2j}\tilde{a}_{2k}] + \dots + \mathbb{E}_\delta [\tilde{a}_{mj}\tilde{a}_{mk}] \\ &= \sum_{i=1}^m \mathbb{E}_\delta [\tilde{a}_{ij}\tilde{a}_{ik}]\end{aligned}$$

which by (5.2) is

$$\begin{aligned}&= \sum_{i=1}^m [p^2 a_{ij} a_{ik} + p(1-p)(a_{ij}\mu + \mu a_{ik}) + (1-p)^2 \mu^2] \\ &= p^2 \sum_{i=1}^m a_{ij} a_{ik} + p(1-p) \left(\sum_{i=1}^m a_{ij} \mu + \sum_{i=1}^m \mu a_{ik} \right) + (1-p)^2 \sum_{i=1}^m \mu^2 \\ &= p^2 \mathbf{A}_{:j}^T \mathbf{A}_{:k} + p(1-p) (\mathbf{A}_{:j}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{:k}) + (1-p)^2 \boldsymbol{\mu}^T \boldsymbol{\mu}\end{aligned}$$

when $j \neq k$, and

$$\begin{aligned}&= \sum_{i=1}^m [p a_{ij} a_{ik} + (1-p) \mu^2] \\ &= p \sum_{i=1}^m a_{ij} a_{ik} + (1-p) \sum_{i=1}^m \mu^2 \\ &= p \mathbf{A}_{:j}^T \mathbf{A}_{:k} + (1-p) \boldsymbol{\mu}^T \boldsymbol{\mu}\end{aligned}$$

when $j = k$. But these are exactly $\Pi[\mathbf{A}_{:j}^T, \mathbf{A}_{:k}]$ and $\pi[\mathbf{A}_{:j}^T, \mathbf{A}_{:k}]$, respectively.

5.1. Proof of Theorem 3.3. Here we seek to show that $s_j(\mathbf{x})$ is an unbiased estimator for the gradient of our loss function, i.e. $\mathbb{E}[s_j(\mathbf{x})] = \nabla L(\mathbf{x})$. To do this we use the law of iterated expectation $\mathbb{E}[s_j(\mathbf{x})] = \mathbb{E}_j[\mathbb{E}_\delta[s_j(\mathbf{x})]]$, showing one at a time that for some function $\ell_j(\mathbf{x})$, $\mathbb{E}_\delta[s_j(\mathbf{x})] = \ell_j(\mathbf{x})$, while $\mathbb{E}_j[\ell_j(\mathbf{x})] = \nabla L(\mathbf{x})$.

Proof. Recall that $s_j(\mathbf{x}) = (c_j(\mathbf{x}) - d_j(\mathbf{x})) \mathbf{e}_j$ such that

$$c_j(\mathbf{x}) := \frac{1}{p^2} \tilde{\mathbf{A}}_{:j}^T (\tilde{\mathbf{A}} \mathbf{x} - p \mathbf{y}) - \frac{1-p}{p^2} \tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:j} x_j,$$

and

$$\begin{aligned}d_j(\mathbf{x}) &:= \frac{1-p}{p^2} \left[\left(\tilde{\mathbf{A}}_{:j}^T \mathbf{M} + \boldsymbol{\mu}^T \tilde{\mathbf{A}} - (1-p) \boldsymbol{\mu}^T \mathbf{M} \right) \mathbf{x} \right. \\ &\quad \left. - p \boldsymbol{\mu}^T \mathbf{y} - \left(\tilde{\mathbf{A}}_{:j}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \tilde{\mathbf{A}}_{:j} - \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j \right],\end{aligned}$$

where $c_j(\mathbf{x})$ addresses the missing entry without mean shift, while $d_j(\mathbf{x})$ accounts for mean

shift. Now, observe that

$$(5.6) \quad \mathbb{E}_\delta [s_j(\mathbf{x})] = (\mathbb{E}_\delta [c_j(\mathbf{x})] - \mathbb{E}_\delta [d_j(\mathbf{x})]) \mathbf{e}_j,$$

where the expectations of c_j and d_j will be functions of inner products of columns of \mathbf{A} . Using (5.3) and (5.4) that notation can be simplified. In particular, we can write

$$\begin{aligned} \mathbb{E}_\delta [\tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}] &= \mathbb{E}_\delta \left[\left(\tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:1}, \tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:2}, \dots, \tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:j}, \dots, \tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:n} \right) \right] \\ &= \left(\mathbb{E}_\delta [\tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:1}], \mathbb{E}_\delta [\tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:2}], \dots, \mathbb{E}_\delta [\tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:j}], \dots, \mathbb{E}_\delta [\tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:n}] \right) \\ &= \left(\Pi [\mathbf{A}_{:j}^T, \mathbf{A}_{:1}], \Pi [\mathbf{A}_{:j}^T, \mathbf{A}_{:2}], \dots, \pi [\mathbf{A}_{:j}^T, \mathbf{A}_{:j}], \dots, \Pi [\mathbf{A}_{:j}^T, \mathbf{A}_{:n}] \right) \\ (5.7) \quad &= \Pi [\mathbf{A}_{:j}^T, \mathbf{A}] - \Pi [\mathbf{A}_{:j}^T, \mathbf{A}_{:j}] \mathbf{e}_j^T + \pi [\mathbf{A}_{:j}^T, \mathbf{A}_{:j}] \mathbf{e}_j^T, \end{aligned}$$

Thus, using (5.7), we compute the expectation of c_j and d_j :

$$\begin{aligned} \mathbb{E}_\delta [c_j(\mathbf{x})] &= \frac{1}{p^2} \mathbb{E}_\delta [\tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}] \mathbf{x} - \frac{1}{p} \mathbb{E}_\delta [\tilde{\mathbf{A}}_{:j}^T] \mathbf{y} - \frac{1-p}{p^2} \mathbb{E}_\delta [\tilde{\mathbf{A}}_{:j}^T \tilde{\mathbf{A}}_{:j}] x_j \\ &= \frac{1}{p^2} \left(\Pi [\mathbf{A}_{:j}^T, \mathbf{A}] \mathbf{x} - \Pi [\mathbf{A}_{:j}^T, \mathbf{A}_{:j}] x_j + \pi [\mathbf{A}_{:j}^T, \mathbf{A}_{:j}] x_j \right) \\ &\quad - \frac{1}{p} (p \mathbf{A}_{:j}^T + (1-p) \boldsymbol{\mu}^T) \mathbf{y} - \frac{1-p}{p^2} (\pi [\mathbf{A}_{:j}^T, \mathbf{A}_{:j}]) x_j \\ &= \frac{1}{p^2} (p^2 \mathbf{A}_{:j}^T \mathbf{A} + p(1-p) (\mathbf{A}_{:j}^T \mathbf{M} + \boldsymbol{\mu}^T \mathbf{A}) + (1-p)^2 \boldsymbol{\mu}^T \mathbf{M}) \mathbf{x} \\ &\quad - \frac{1}{p^2} (p^2 \mathbf{A}_{:j}^T \mathbf{A}_{:j} + p(1-p) (\mathbf{A}_{:j}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{:j}) + (1-p)^2 \boldsymbol{\mu}^T \boldsymbol{\mu}) x_j \\ &\quad + \frac{1}{p^2} (p \mathbf{A}_{:j}^T \mathbf{A}_{:j} + (1-p) \boldsymbol{\mu}^T \boldsymbol{\mu}) x_j - \frac{1}{p} (p \mathbf{A}_{:j}^T + (1-p) \boldsymbol{\mu}^T) \mathbf{y} \\ &\quad - \frac{1-p}{p^2} (p \mathbf{A}_{:j}^T \mathbf{A}_{:j} + (1-p) \boldsymbol{\mu}^T \boldsymbol{\mu}) x_j \\ &= \left(\mathbf{A}_{:j}^T \mathbf{A} + \frac{1-p}{p} (\mathbf{A}_{:j}^T \mathbf{M} + \boldsymbol{\mu}^T \mathbf{A}) + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \mathbf{M} \right) \mathbf{x} \\ &\quad - \left(\mathbf{A}_{:j}^T \mathbf{A} + \frac{1-p}{p} (\mathbf{A}_{:j}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{:j}) + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j \\ &\quad + \left(\frac{1}{p} \mathbf{A}_{:j}^T \mathbf{A} + \frac{1-p}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j - \left(\mathbf{A}_{:j}^T + \frac{1-p}{p} \boldsymbol{\mu}^T \right) \mathbf{y} \\ &\quad - \left(\frac{1-p}{p} \mathbf{A}_{:j}^T \mathbf{A} + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j \\ (5.8) \quad &= \mathbf{A}_{:j}^T \mathbf{A} \mathbf{x} - \mathbf{A}_{:j}^T \mathbf{y} + \left(\frac{1-p}{p} (\mathbf{A}_{:j}^T \mathbf{M} + \boldsymbol{\mu}^T \mathbf{A}) + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \mathbf{M} \right) \mathbf{x} \end{aligned}$$

$$\begin{aligned}
& - \left(\frac{1-p}{p} (\mathbf{A}_{:j}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{:j}) + \frac{2(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j \\
& + \left(\frac{1-p}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j - \frac{1-p}{p} \boldsymbol{\mu}^T \mathbf{y} \\
\mathbb{E}_\delta [d_j(\mathbf{x})] &= \frac{1-p}{p^2} \left[\left(\mathbb{E}_\delta \left[\tilde{\mathbf{A}}_{:j}^T \mathbf{M} \right] + \mathbb{E}_\delta \left[\boldsymbol{\mu}^T \tilde{\mathbf{A}} \right] - (1-p) \boldsymbol{\mu}^T \mathbf{M} \right) \mathbf{x} \right. \\
& \quad \left. - p \boldsymbol{\mu}^T \mathbf{y} - \left(\mathbb{E}_\delta \left[\tilde{\mathbf{A}}_{:j}^T \boldsymbol{\mu} \right] + \mathbb{E}_\delta \left[\boldsymbol{\mu}^T \tilde{\mathbf{A}}_{:j} \right] - \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j \right] \\
&= \frac{1-p}{p^2} \left[\left(\pi \left[\mathbf{A}_{:j}^T, \mathbf{M} \right] + \pi \left[\boldsymbol{\mu}^T, \mathbf{A} \right] - (1-p) \boldsymbol{\mu}^T \mathbf{M} \right) \mathbf{x} \right. \\
& \quad \left. - p \boldsymbol{\mu}^T \mathbf{y} - \left(\pi \left[\mathbf{A}_{:j}^T, \boldsymbol{\mu} \right] + \pi \left[\boldsymbol{\mu}^T, \mathbf{A}_{:j} \right] - \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j \right] \\
&= \frac{1-p}{p^2} \left[\left(p \mathbf{A}_{:j}^T \mathbf{M} + p \boldsymbol{\mu}^T \mathbf{A} + 2(1-p) \boldsymbol{\mu}^T \mathbf{M} - (1-p) \boldsymbol{\mu}^T \mathbf{M} \right) \mathbf{x} \right. \\
& \quad \left. - p \boldsymbol{\mu}^T \mathbf{y} - \left(p \mathbf{A}_{:j}^T \boldsymbol{\mu} + p \boldsymbol{\mu}^T \mathbf{A}_{:j} + 2(1-p) \boldsymbol{\mu}^T \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j \right] \\
&= \left(\frac{1-p}{p} \mathbf{A}_{:j}^T \mathbf{M} + \frac{1-p}{p} \boldsymbol{\mu}^T \mathbf{A} + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \mathbf{M} \right) \mathbf{x} - \frac{1-p}{p} \boldsymbol{\mu}^T \mathbf{y} \\
& \quad - \left(\frac{1-p}{p} \mathbf{A}_{:j}^T \boldsymbol{\mu} - \frac{1-p}{p} \boldsymbol{\mu}^T \mathbf{A}_{:j} - \frac{(1-p)(1-2p)}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j \\
(5.9) \quad &= \left(\frac{1-p}{p} (\mathbf{A}_{:j}^T \mathbf{M} + \boldsymbol{\mu}^T \mathbf{A}) + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \mathbf{M} \right) \mathbf{x} \\
& \quad - \left(\frac{1-p}{p} (\mathbf{A}_{:j}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{:j}) + \frac{2(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j \\
& \quad + \left(\frac{1-p}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) x_j - \frac{1-p}{p} \boldsymbol{\mu}^T \mathbf{y}
\end{aligned}$$

Plugging (5.8) and (5.9) into (5.6), we get

$$(5.10) \quad \mathbb{E}_\delta [s_j(\mathbf{x})] = (\mathbb{E}_\delta [c_j(\mathbf{x})] - \mathbb{E}_\delta [d_j(\mathbf{x})]) \mathbf{e}_j = (\mathbf{A}_{:j}^T (\mathbf{A}\mathbf{x} - \mathbf{y})) \mathbf{e}_j,$$

for which we refer to the right-hand side expression as “ $\ell_j(\mathbf{x})$ ”. Recalling our loss function $L(\mathbf{x}) = \frac{1}{2n} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2$, it can be shown after some expansions that

$$\frac{\partial L(\mathbf{x})}{\partial x_j} = \frac{1}{n} \mathbf{A}_{:j}^T (\mathbf{A}\mathbf{x} - \mathbf{y}).$$

Thus, when we go to compute the column-wise expectation of $\ell_j(\mathbf{x})$ (assuming columns of \mathbf{A} are selected uniformly at random) we get

$$(5.11) \quad \mathbb{E}_j [\ell_j(\mathbf{x})] = \sum_{j=1}^n \frac{1}{n} \ell_j(\mathbf{x}) = \sum_{j=1}^n \frac{1}{n} (\mathbf{A}_{:j}^T (\mathbf{A}\mathbf{x} - \mathbf{y})) \mathbf{e}_j = \sum_{j=1}^n \frac{\partial L(\mathbf{x})}{\partial x_j} \mathbf{e}_j.$$

Finally, combining (5.10) and (5.11) gives us the desired gradient,

$$\mathbb{E} [s_j(\mathbf{x})] = \mathbb{E}_j [\mathbb{E}_\delta [s_j(\mathbf{x})]] = \mathbb{E}_j [\ell_j(\mathbf{x})] = \nabla L(\mathbf{x}). \quad \blacksquare$$

5.2. Proof of Theorem 3.4. Similar to the last proof, here we show that $t_i(\mathbf{x})$ is an unbiased estimator for the gradient of the loss function, $\mathbb{E}[t_i(\mathbf{x})] = \nabla F(\mathbf{x})$. Again, we use an iterated expectation $\mathbb{E}[t_i(\mathbf{x})] = \mathbb{E}_i[\mathbb{E}_\delta[t_i(\mathbf{x})]]$, showing first that for the function $f_i(\mathbf{x})$, $\mathbb{E}_\delta[t_i(\mathbf{x})] = \nabla f_i(\mathbf{x})$ holds, then finally that $\mathbb{E}_i[\nabla f_i(\mathbf{x})] = \nabla F(\mathbf{x})$.

Proof. Recall that $t_i(\mathbf{x}) = g_i(\mathbf{x}) - h_i(\mathbf{x})$ such that

$$g_i(\mathbf{x}) := \frac{1}{p^2} \tilde{\mathbf{A}}_{i:}^T \left(\tilde{\mathbf{A}}_{i:}^T \mathbf{x} - py_i \right) - \frac{1-p}{p^2} \text{diag} \left(\tilde{\mathbf{A}}_{i:}^T \tilde{\mathbf{A}}_{i:} \right) \mathbf{x},$$

and

$$h_i(\mathbf{x}) := \frac{1-p}{p^2} \left[\left(\tilde{\mathbf{A}}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \tilde{\mathbf{A}}_{i:} - (1-p) \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} - p \boldsymbol{\mu}^T y_i - \text{diag} \left(\tilde{\mathbf{A}}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \tilde{\mathbf{A}}_{i:} - \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \right],$$

where $g_i(\mathbf{x})$ plays the analogous role to $c_j(\mathbf{x})$ in the previous proof (missing entry without mean shift), and $h_i(\mathbf{x})$ the role of $d_j(\mathbf{x})$ (mean shift). Now, observe that

$$(5.12) \quad \mathbb{E}_\delta[t_i(\mathbf{x})] = \mathbb{E}_\delta[g_i(\mathbf{x})] - \mathbb{E}_\delta[h_i(\mathbf{x})],$$

where the expectations of g_i and h_i will be functions of inner products of columns of \mathbf{A} . Using (5.3) and (5.4) that notation can be simplified. In particular, we can write

$$(5.13) \quad \begin{aligned} \mathbb{E}_\delta \left[\tilde{\mathbf{A}}_{i:}^T \tilde{\mathbf{A}}_{i:} \right] &= \mathbb{E}_\delta \left[\begin{pmatrix} \tilde{a}_{i1} \tilde{a}_{i1}, & \tilde{a}_{i1} \tilde{a}_{i2}, & \dots & \tilde{a}_{i1} \tilde{a}_{in} \\ \tilde{a}_{i2} \tilde{a}_{i1}, & \tilde{a}_{i2} \tilde{a}_{i2}, & \dots & \tilde{a}_{i2} \tilde{a}_{in} \\ \vdots & \vdots & \ddots & \vdots \\ \tilde{a}_{in} \tilde{a}_{i1}, & \tilde{a}_{in} \tilde{a}_{i2}, & \dots & \tilde{a}_{in} \tilde{a}_{in} \end{pmatrix} \right] \\ &= \begin{pmatrix} \mathbb{E}_\delta [\tilde{a}_{i1} \tilde{a}_{i1}], & \mathbb{E}_\delta [\tilde{a}_{i1} \tilde{a}_{i2}], & \dots & \mathbb{E}_\delta [\tilde{a}_{i1} \tilde{a}_{in}] \\ \mathbb{E}_\delta [\tilde{a}_{i2} \tilde{a}_{i1}], & \mathbb{E}_\delta [\tilde{a}_{i2} \tilde{a}_{i2}], & \dots & \mathbb{E}_\delta [\tilde{a}_{i2} \tilde{a}_{in}] \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{E}_\delta [\tilde{a}_{in} \tilde{a}_{i1}], & \mathbb{E}_\delta [\tilde{a}_{in} \tilde{a}_{i2}], & \dots & \mathbb{E}_\delta [\tilde{a}_{in} \tilde{a}_{in}] \end{pmatrix} \\ &= \begin{pmatrix} \pi [a_{i1}, a_{i1}], & \Pi [a_{i1}, a_{i2}], & \dots & \Pi [a_{i1}, a_{in}] \\ \Pi [a_{i2}, a_{i1}], & \pi [a_{i2}, a_{i2}], & \dots & \Pi [a_{i2}, a_{in}] \\ \vdots & \vdots & \ddots & \vdots \\ \Pi [a_{in}, a_{i1}], & \Pi [a_{in}, a_{i2}], & \dots & \pi [a_{in}, a_{in}] \end{pmatrix} \\ &= \Pi [\mathbf{A}_{i:}^T, \mathbf{A}_{i:}] - \text{diag} (\Pi [\mathbf{A}_{i:}^T, \mathbf{A}_{i:}]) + \text{diag} (\pi [\mathbf{A}_{i:}^T, \mathbf{A}_{i:}]). \end{aligned}$$

Thus, taking each expectation separately and using (5.13), we get:

$$\begin{aligned}
\mathbb{E}_\delta [g_i(\mathbf{x})] &= \frac{1}{p^2} \mathbb{E}_\delta \left[\tilde{\mathbf{A}}_{i:}^T \tilde{\mathbf{A}}_{i:} \right] \mathbf{x} - \frac{1}{p} \mathbb{E}_\delta \left[\tilde{\mathbf{A}}_{i:}^T \right] y_i - \frac{1-p}{p^2} \text{diag} \left(\mathbb{E}_\delta \left[\tilde{\mathbf{A}}_{i:}^T \tilde{\mathbf{A}}_{i:} \right] \right) \mathbf{x} \\
&= \frac{1}{p^2} \left(\Pi \left[\mathbf{A}_{i:}^T, \mathbf{A}_{i:} \right] - \text{diag} \left(\Pi \left[\mathbf{A}_{i:}^T, \mathbf{A}_{i:} \right] \right) + \text{diag} \left(\pi \left[\mathbf{A}_{i:}^T, \mathbf{A}_{i:} \right] \right) \right) \mathbf{x} \\
&\quad - \frac{1}{p} \left(p \mathbf{A}_{i:}^T + (1-p) \boldsymbol{\mu}^T \right) y_i - \frac{1-p}{p^2} \left(\text{diag} \left(\pi \left[\mathbf{A}_{i:}^T, \mathbf{A}_{i:} \right] \right) \right) \mathbf{x} \\
&= \frac{1}{p^2} \left(p^2 \mathbf{A}_{i:}^T \mathbf{A}_{i:} + p(1-p) \left(\mathbf{A}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{i:} \right) + (1-p)^2 \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
&\quad - \frac{1}{p^2} \text{diag} \left(p^2 \mathbf{A}_{i:}^T \mathbf{A}_{i:} + p(1-p) \left(\mathbf{A}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{i:} \right) + (1-p)^2 \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
&\quad + \frac{1}{p^2} \text{diag} \left(p \mathbf{A}_{i:}^T \mathbf{A}_{i:} + (1-p) \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} - \frac{1}{p} \left(p \mathbf{A}_{i:}^T + (1-p) \boldsymbol{\mu}^T \right) y_i \\
&\quad - \frac{1-p}{p^2} \text{diag} \left(p \mathbf{A}_{i:}^T \mathbf{A}_{i:} + (1-p) \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
&= \left(\mathbf{A}_{i:}^T \mathbf{A}_{i:} + \frac{1-p}{p} \left(\mathbf{A}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{i:} \right) + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
&\quad - \text{diag} \left(\mathbf{A}_{i:}^T \mathbf{A}_{i:} + \frac{1-p}{p} \left(\mathbf{A}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{i:} \right) + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
&\quad + \text{diag} \left(\frac{1}{p} \mathbf{A}_{i:}^T \mathbf{A}_{i:} + \frac{1-p}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} - \left(\mathbf{A}_{i:}^T + \frac{1-p}{p} \boldsymbol{\mu}^T \right) y_i \\
&\quad - \text{diag} \left(\frac{1-p}{p} \mathbf{A}_{i:}^T \mathbf{A}_{i:} + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
(5.14) \quad &= \mathbf{A}_{i:}^T \mathbf{A}_{i:} \mathbf{x} - \mathbf{A}_{i:}^T y_i + \left(\frac{1-p}{p} \left(\mathbf{A}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{i:} \right) + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
&\quad - \text{diag} \left(\frac{1-p}{p} \left(\mathbf{A}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{i:} \right) + \frac{2(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
&\quad + \text{diag} \left(\frac{1-p}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} - \frac{1-p}{p} \boldsymbol{\mu}^T y_i,
\end{aligned}$$

and

$$\begin{aligned}
\mathbb{E}_\delta [h_i(\mathbf{x})] &= \frac{1-p}{p^2} \left[\left(\mathbb{E}_\delta \left[\tilde{\mathbf{A}}_{i:}^T \boldsymbol{\mu} \right] + \mathbb{E}_\delta \left[\boldsymbol{\mu}^T \tilde{\mathbf{A}}_{i:} \right] - (1-p) \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \right. \\
&\quad \left. - p \boldsymbol{\mu}^T y_i - \text{diag} \left(\mathbb{E}_\delta \left[\tilde{\mathbf{A}}_{i:}^T \boldsymbol{\mu} \right] + \mathbb{E}_\delta \left[\boldsymbol{\mu}^T \tilde{\mathbf{A}}_{i:} \right] - \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \right] \\
&= \frac{1-p}{p^2} \left[\left(\pi \left[\tilde{\mathbf{A}}_{i:}^T \boldsymbol{\mu} \right] + \pi \left[\boldsymbol{\mu}^T \tilde{\mathbf{A}}_{i:} \right] - (1-p) \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \right. \\
&\quad \left. - p \boldsymbol{\mu}^T y_i - \text{diag} \left(\pi \left[\tilde{\mathbf{A}}_{i:}^T \boldsymbol{\mu} \right] + \pi \left[\boldsymbol{\mu}^T \tilde{\mathbf{A}}_{i:} \right] - \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \right] \\
&= \frac{1-p}{p^2} \left[\left(p \mathbf{A}_{i:}^T \boldsymbol{\mu} + p \boldsymbol{\mu}^T \mathbf{A}_{i:} + 2(1-p) \boldsymbol{\mu}^T \boldsymbol{\mu} - (1-p) \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \right.
\end{aligned}$$

$$\begin{aligned}
& -p\boldsymbol{\mu}^T y_i - \text{diag} \left(p\mathbf{A}_{i:}^T \boldsymbol{\mu} + p\boldsymbol{\mu}^T \mathbf{A}_{i:} + 2(1-p)\boldsymbol{\mu}^T \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
= & \left(\frac{1-p}{p} \mathbf{A}_{i:}^T \boldsymbol{\mu} + \frac{1-p}{p} \boldsymbol{\mu}^T \mathbf{A}_{i:} + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} - \frac{1-p}{p} \boldsymbol{\mu}^T y_i \\
& - \text{diag} \left(\frac{1-p}{p} \mathbf{A}_{i:}^T \boldsymbol{\mu} - \frac{1-p}{p} \boldsymbol{\mu}^T \mathbf{A}_{i:} - \frac{(1-p)(1-2p)}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
(5.15) \quad = & \left(\frac{1-p}{p} (\mathbf{A}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{i:}) + \frac{(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
& - \text{diag} \left(\frac{1-p}{p} (\mathbf{A}_{i:}^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A}_{i:}) + \frac{2(1-p)^2}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} \\
& + \text{diag} \left(\frac{1-p}{p^2} \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \mathbf{x} - \frac{1-p}{p} \boldsymbol{\mu}^T y_i.
\end{aligned}$$

Plugging (5.14) and (5.15) into (5.12), we get

$$(5.16) \quad \mathbb{E}_\delta [t_i(\mathbf{x})] = \mathbb{E}_\delta [g_i(\mathbf{x})] - \mathbb{E}_\delta [h_i(\mathbf{x})] = \mathbf{A}_{i:}^T (\mathbf{A}_{i:} \mathbf{x} - y_i).$$

Now, recall our objective function $F(\mathbf{x}) = \frac{1}{2m} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2$, which can be rewritten as $F(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x})$, where $f_i(\mathbf{x}) = \frac{1}{2} (\mathbf{A}_{i:} \mathbf{x} - y_i)^2$. After some expansions, it is not too difficult to show that

$$(5.17) \quad \nabla f_i(\mathbf{x}) = \mathbf{A}_{i:}^T (\mathbf{A}_{i:} \mathbf{x} - y_i),$$

and notice this is exactly the right-hand side of (5.16). Thus, assuming rows i are uniformly selected from $[m]$, we take the row-wise expectation of $\nabla f_i(\mathbf{x})$ to get

$$(5.18) \quad \mathbb{E}_i [\nabla f_i(\mathbf{x})] = \sum_{i=1}^m \frac{1}{m} \nabla f_i(\mathbf{x}) = \nabla \left[\frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}) \right] = \nabla F(\mathbf{x}).$$

Finally, combining (5.16), (5.17), and (5.18), we get

$$\mathbb{E} [t_i(\mathbf{x})] = \mathbb{E}_i [\mathbb{E}_\delta [t_i(\mathbf{x})]] = \mathbb{E}_i [\nabla f_i(\mathbf{x})] = \nabla F(\mathbf{x}). \quad \blacksquare$$

the desired gradient.

5.3. Proof of Theorem 3.5. In this proof, we first define the μ , then we consider the different cases, such as 0-imputation and μ -imputation for the matrix \mathbf{A} . By calculating the expectation for both cases, $\mathbb{E}_\delta \|\mathbf{A} - \tilde{\mathbf{A}}_0\|_F^2 \geq \mathbb{E}_\delta \|\mathbf{A} - \tilde{\mathbf{A}}_\mu\|_F^2$ holds.

Proof. Given $\tilde{\mathbf{A}}_0$ is matrix \mathbf{A} with missing entries $\tilde{a}_{0,ij}$ in which 0 is imputed and $\tilde{\mathbf{A}}_\mu$ is matrix \mathbf{A} with entries $\tilde{a}_{\mu,ij}$ in which the mean value of \mathbf{A} , a fixed $\mu = \frac{1}{mn} \sum_{ij} a_{ij}$, is imputed, we want to find $\mathbb{E}_\delta [\|\mathbf{A} - \tilde{\mathbf{A}}_0\|_F^2]$ and $\mathbb{E}_\delta [\|\mathbf{A} - \tilde{\mathbf{A}}_\mu\|_F^2]$ and compare them.

We know,

$$\tilde{a}_{0,ij} = \begin{cases} 0 & w.p \ 1-p \\ a_{ij} & w.p \ p \end{cases} \Rightarrow \tilde{a}_{0,ij}^2 = \begin{cases} 0 & w.p \ 1-p \\ a_{ij}^2 & w.p \ p \end{cases}$$

$$\tilde{a}_{\mu,ij} = \begin{cases} \mu & w.p. 1-p \\ a_{ij} & w.p. p \end{cases} \Rightarrow \tilde{a}_{\mu,ij}^2 = \begin{cases} \mu^2 & w.p. 1-p \\ a_{ij}^2 & w.p. p \end{cases}$$

Therefore,

$$\begin{aligned} \mathbb{E}_\delta \left[\|\mathbf{A} - \tilde{\mathbf{A}}_0\|_F^2 \right] &= \mathbb{E}_\delta \left[\sum_{i=1}^m \sum_{j=1}^n (a_{ij} - \tilde{a}_{0,ij})^2 \right] \\ &= \sum_{i=1}^m \sum_{j=1}^n \mathbb{E}_\delta \left[(a_{ij}^2 - 2a_{ij}\tilde{a}_{0,ij} + \tilde{a}_{0,ij}^2) \right] \\ &= \sum_{i=1}^m \sum_{j=1}^n \mathbb{E}_\delta [a_{ij}^2] + \sum_{i=1}^m \sum_{j=1}^n \mathbb{E}_{\mathbf{A},\delta} [\tilde{a}_{0,ij}^2] - \sum_{i=1}^m \sum_{j=1}^n \mathbb{E}_\delta [2a_{ij}\tilde{a}_{0,ij}] \\ &= \|\mathbf{A}\|_F^2 + \sum_{i=1}^m \sum_{j=1}^n pa_{ij}^2 - \sum_{i=1}^m \sum_{j=1}^n 2pa_{ij}^2 \\ &= \|\mathbf{A}\|_F^2 - \sum_{i=1}^m \sum_{j=1}^n pa_{ij}^2 \\ &= \|\mathbf{A}\|_F^2 - p\|\mathbf{A}\|_F^2 = (1-p)\|\mathbf{A}\|_F^2 \end{aligned}$$

Similarly,

$$\begin{aligned} \mathbb{E}_\delta \left[\|\mathbf{A} - \tilde{\mathbf{A}}_\mu\|_F^2 \right] &= \mathbb{E}_\delta \left[\sum_{i=1}^m \sum_{j=1}^n (a_{ij} - \tilde{a}_{\mu,ij})^2 \right] \\ &= \sum_{i=1}^m \sum_{j=1}^n \mathbb{E}_\delta \left[(a_{ij}^2 - 2a_{ij}\tilde{a}_{\mu,ij} + \tilde{a}_{\mu,ij}^2) \right] \\ &= \sum_{i=1}^m \sum_{j=1}^n \mathbb{E}_\delta [a_{ij}^2] + \sum_{i=1}^m \sum_{j=1}^n \mathbb{E}_\delta [\tilde{a}_{\mu,ij}^2] - \sum_{i=1}^m \sum_{j=1}^n \mathbb{E}_\delta [2a_{ij}\tilde{a}_{\mu,ij}] \\ &= \|\mathbf{A}\|_F^2 + \sum_{i=1}^m \sum_{j=1}^n (pa_{ij}^2 + \mu^2(1-p)) - \sum_{i=1}^m \sum_{j=1}^n 2a_{ij}(pa_{ij} + \mu(1-p)) \\ &= \|\mathbf{A}\|_F^2 + \sum_{i=1}^m \sum_{j=1}^n pa_{ij}^2 + \sum_{i=1}^m \sum_{j=1}^n \mu^2(1-p) - 2 \sum_{i=1}^m \sum_{j=1}^n pa_{ij}^2 - 2 \sum_{i=1}^m \sum_{j=1}^n a_{ij}\mu(1-p) \\ &= \|\mathbf{A}\|_F^2 - \sum_{i=1}^m \sum_{j=1}^n pa_{ij}^2 + \sum_{i=1}^m \sum_{j=1}^n \mu^2(1-p) - 2 \sum_{i=1}^m \sum_{j=1}^n a_{ij}\mu(1-p) \\ &= \|\mathbf{A}\|_F^2 - \sum_{i=1}^m \sum_{j=1}^n pa_{ij}^2 + \sum_{i=1}^m \sum_{j=1}^n \mu(1-p)(\mu - 2a_{ij}) \\ &= (1-p)\|\mathbf{A}\|_F^2 + \sum_{i=1}^m \sum_{j=1}^n \mu(1-p)(\mu - 2a_{ij}) \end{aligned}$$

For the last term we have,

$$\begin{aligned}
\sum_{i=1}^m \sum_{j=1}^n \mu(1-p)(\mu - 2a_{ij}) &= \mu(1-p) \sum_{i=1}^m \sum_{j=1}^n (\mu - 2a_{ij}) \\
&= \mu(1-p) \left(mn\mu - \sum_{i=1}^m \sum_{j=1}^n 2a_{ij} \right) \\
&= \mu(1-p)(mn\mu - 2mn\mu) \\
&= -mn\mu^2(1-p),
\end{aligned}$$

which is non-positive. Thus, $\mathbb{E}_\delta[\|\mathbf{A} - \tilde{\mathbf{A}}_0\|_F^2] \geq \mathbb{E}_\delta[\|\mathbf{A} - \tilde{\mathbf{A}}_\mu\|_F^2]$ ■

REFERENCES

- [1] R. JIANG, T. SUN, D. SONG, AND J. J. LI, *Statistics or biology: the zero-inflation controversy about scrna-seq data*, *Genome biology*, 23 (2022), pp. 1–24.
- [2] D. LEVENTHAL AND A. S. LEWIS, *Randomized methods for linear constraints: Convergence rates and conditioning*, 2008.
- [3] A. MA AND D. NEEDELL, *Stochastic gradient descent for linear systems with missing data*.
- [4] A. MA, D. NEEDELL, AND A. RAMDAS, *Convergence properties of the randomized extended gauss–seidel and kaczmarz methods*, *SIAM Journal on Matrix Analysis and Applications*, 36 (2015), pp. 1590–1604.
- [5] S. MUKHOPADHYAY, *Stochastic gradient descent for linear systems with sequential matrix entry accumulation*, *Signal Processing*, 171 (2020), p. 107494.
- [6] S. MUKHOPADHYAY AND A. MUKHERJEE, *Imdlms: An imputation based lms algorithm for linear system identification with missing input data*, *IEEE Transactions on Signal Processing*, 68 (2020), pp. 2370–2385.
- [7] D. NEEDELL, *Randomized kaczmarz solver for noisy linear systems*, *BIT Numerical Mathematics*, 50 (2010), pp. 395–403.
- [8] Y. NESTEROV, *Efficiency of coordinate descent methods on huge-scale optimization problems*, *SIAM Journal on Optimization*, 22 (2012), pp. 341–362.
- [9] M. S. RAHIM, A. A. IMRAN, AND T. AHMED, *Mining the productivity data of garment industry*, *International Journal of Business Intelligence and Data Mining*, 1 (2021), p. 1.