# SAMPLING KACZMARZ MOTZKIN METHOD FOR LINEAR FEASIBILITY PROBLEMS: GENERALIZATION & ACCELERATION

**Md Sarowar Morshed**
Department of Mechanical & Industrial Engineering
Northeastern University
Boston, MA
morshed.m@northeastern.edu

**Md Saiful Islam**
Department of Mechanical & Industrial Engineering
Northeastern University
Boston, MA
islam.m@northeastern.edu

**Md. Noor-E-Alam**
Department of Mechanical & Industrial Engineering
Northeastern University
Boston, MA
mnalam@neu.edu

December 8, 2020

## ABSTRACT

*Randomized Kaczmarz* (RK), *Motzkin Method* (MM) and *Sampling Kaczmarz Motzkin* (SKM) algorithms are commonly used iterative techniques for solving a system of linear inequalities (i.e., $Ax \leq b$). As linear systems of equations represent a modeling paradigm for solving many optimization problems, these randomized and iterative techniques are gaining popularity among researchers in different domains. In this work, we propose a *Generalized Sampling Kaczmarz Motzkin* (GSKM) method that unifies the iterative methods into a single framework. In addition to the general framework, we propose a Nesterov type acceleration scheme in the SKM method called as *Probably Accelerated Sampling Kaczmarz Motzkin* (PASKM). We prove the convergence theorems for both GSKM and PASKM algorithms in the $L_2$ norm perspective with respect to the proposed sampling distribution. Furthermore, we prove sub-linear convergence for the Cesaro average of iterates for the proposed GSKM and PASKM algorithms. From the convergence theorem of the GSKM algorithm, we find the convergence results of several well-known algorithms like the Kaczmarz method, Motzkin method and SKM algorithm. We perform thorough numerical experiments using both randomly generated and real-world (classification with support vector machine and Netlib LP) test instances to demonstrate the efficiency of the proposed methods. We compare the proposed algorithms with SKM, *Interior Point Method* (IPM) and *Active Set Method* (ASM) in terms of computation time and solution quality. In the majority of the problem instances, the proposed generalized and accelerated algorithms significantly outperform the state-of-the-art methods.

***Keywords*** Kaczmarz Method · Randomized Projection · Sampling Kaczmarz Motzkin · Linear Feasibility · Nesterov's Acceleration · Iterative Methods

## 1 Introduction

We consider the following *Linear Feasibility* (LF) problem:

$$Ax \leq b, \ \ b \in \mathbb{R}^m, \ A \in \mathbb{R}^{m \times n}. \tag{1}$$

We confine the scope of our work in the regime of thin/tall coefficient matrix $A$ ($m \gg n$), as iterative methods are more competitive for such problems. Note that, while almost all of the classical methods are deterministic in nature, recent advances [1–12] suggest that randomized iterative methods can outperform existing deterministic methods for solving

many computational problems including linear feasibility, linear systems and convex optimization problems. From an algorithmic point of view, our work generalizes the SKM method and furthermore explores the possibility of faster variants of these methods. Before we delve into the contributions of this work, we give brief descriptions of some of the classical and modern techniques related to solving LF problems with iterative methods.

**Randomized Kaczmarz (RK)**   Kaczmarz method is one of the popular methods for solving linear systems due to its algorithmic simplicity [13]. Originally proposed in 1937 by Kaczmarz [13], the Kaczmarz method remained hidden to the research community until the early 1980s, when Gordon *et. al* proposed *Algebraic Reconstruction Techniques* (ART) in the area of image reconstruction [14]. Later, it has found applications in several areas like computer tomography [15, 16], digital signal processing [17], distributed computing [18, 19] and many other engineering and physics problems. It has been rediscovered several times as a family of methods including component solution, successive projection, row-action and cyclic projection methods (see [20]). Given a current point $x_k$, the Kaczmarz method generates new update $x_{k+1}$ based on the orthogonal projection of $x_k$ onto the hyper-plane $a_{i^*}^T x_k \leq b_{i^*}$,

$$x_{k+1} = x_k - \delta \frac{\left(a_{i^*}^T x_k - b_{i^*}\right)^+}{\|a_{i^*}\|_2^2} a_{i^*}. \tag{2}$$

The differences between the old and modern Kaczmarz schemes are the choice of projection hyper-planes in the update formula of equation (2) at each iteration and the choice of projection parameter $\delta$. The original Kaczmarz method chooses hyper-planes by $i^* \equiv k \mod m, k = 1, 2, 3, ..., m$ with parameter $\delta = 1$. Strohmer *et. al* [1] showed that instead of using cyclic rules, convergence can be improved by choosing $i^*$ from $\{1, 2, ..., m\}$ at random with probability proportional to $\|a_i^*\|_2^2$. This randomization scheme is very efficient for the linear system as well [2]. The projection parameter $\delta$ can be chosen any value in the range of $(0, 2]$ [11].

**Motzkin Method (MM)**   Another classical method for solving LF problems is the Motzkin method (MM) discovered by Motzkin *et. al* in the early 1950s [21, 22]. The work of Motzkin was rediscovered several times by other researchers in the field of *Machine Learning* (ML). For instance, the so-called perceptron algorithm in ML [23–25] can be classified as a member of Motzkin type methods. Furthermore, MM can be sought as the Kaczmarz method with "maximal-residual control" or with "most violated constraint control" [20, 26, 27]. The MM starts with an initial point $x_k$ and finds the next update $x_{k+1}$ as the projection of $x_k$ onto the most violated hyper-plane defined in the equation (1). Given the current point $x_k$, find the next projection hyper-plane $a_{i^*}$ as the maximum violated constraint (i.e., select $i^* = \arg\max_{i \in \{1, 2, ..., m\}} \{a_i^T x_k - b_i\}$) and then update $x_{k+1}$ as follows

$$x_{k+1} = (1 - \delta)x_k + \delta \, \mathcal{P}_{H_{i^*}}(x_k), \tag{3}$$

with the choice $0 \leq \delta < 2$, where $\mathcal{P}_{H_{i^*}}(x_k)$ denotes the orthogonal projection of $x_k$ onto the hyper-plane $H_{i^*} = a_{i^*}^T x_k \leq b_{i^*}$. The analysis of the MM depends on the so-called Hoffman constant (see Lemma 3.1 and Table 1). The main drawback of the standard MM is that it fails to terminate when the LF problem of (1) is infeasible. In the late 1980s, MM resurfaced for its connection to the ellipsoid method [28]. For rational data, it's proven that the system can detect infeasibility and for totally unimodular data, the scheme gives strong polynomial-time algorithms [29]. Recently, Chubanov [30, 31] developed a modified method compared to the traditional relaxation type methods [22], where instead of projecting on the original hyper-plane, one projects the new point to an induced hyper-plane.

In recent time, Kaczmarz type methods gained immense popularity in the research community. The work of Strohmer *et. al* [1] encouraged numerous extensions and variants of the RK method (see [2, 3, 5–8, 32]). For instance, in [5, 33], authors analyzed variants of the Kaczmarz method for a least square setup. A significant breakthrough came from the work of Gower *et. al* when they developed a generalized framework namely the *Gower-Richtarik* (GR) sketch. The authors showed that several well-known algorithms like *Randomized Kaczmarz* (RK), *Randomized Newton* (RN) and *Randomized Coordinate Descent* methods can be sought as special cases of the GR algorithm. For different choices of sampling distribution and a positive definite matrix, one can recover all of the above algorithms as special cases (see [8, 10, 34, 35] for a detailed discussion).

Another area of research spurred when Gower *et. al* provided the extension of the GR sketching method to combine several Quasi-Newton methods into one framework [36]. They showed that almost all of the available Quasi-Newton algorithms like *Bad Broyden* (BB), *Powell-Symmetric-Broyden* (PSB), *Good Broyden* (GB), *Broyden–Fletcher–Goldfarb–Shanno* (BFGS) and *Davidon–Fletcher–Powell* (DFP) can be derived as special cases of the GR sketch. In another work, they extended the GR method for finding the pseudo-inverse of a matrix [37]. Several variants of acceleration have been explored recently for the GR sketch [38, 39]. Special block variants of RK methods have been analyzed by Needell *et. al* [40–42]. From a linear programming perspective, Chubanov developed a polynomial-time algorithm for solving the $0 - 1$ linear system [30, 43, 44] and $0 - 1$ LF problem [31]. In recent time, other variants of both RK and SKM algorithms have been developed that deal with various types of sampling strategies [10, 45–48].

Moreover, a large number of scientific computing and machine learning tasks aim to solve the unconstrained mini-mization problem $x^* = \arg \min \Phi(x)$ with a differentiable function $\Phi : \mathbb{R}^n \to \mathbb{R}$ [49]. *Gradient Descent* (GD) and its variants have been the de facto choice in the artificial intelligence and machine learning community to solve such problems [50]. However, GD suffers from slow convergence as soon as the current solution approaches $x^*$. To achieve faster convergence, one of the major algorithmic development is the idea of momentum. The momentum method was first studied by Polyak [51] in the sense of rolling a heavy ball along with a well-defined cost function. However, despite its intuitiveness, Polyak's heavy ball momentum was difficult to analyze mathematically. Nesterov's acceleration method, proposed by Nesterov in his seminal work [52] for the GD provides the mathematical rigor that Polyak's method lacks and exhibits the worst-case convergence rate of $O(\frac{1}{k^2})$ for minimizing smooth convex functions compared to the original convergence rate of $O(\frac{1}{k})$. Since the introduction of Nesterov's work, numerous work has been done on algorithmic development of the first-order accelerated methods (for a detailed discussion see [53–56]). From then on, Nesterov and Polyak's work has been integrated into several well-known projection-based algorithms like *Coordinate Descent* [56], *Randomized Kaczmarz* [32], *Momentum Induced GR Sketching* [57], *Affine Scaling* [58], *Accelerated Quasi-Newton* [39], *Randomized Gossip* [59], *Sampling Kaczmarz Motzkin* [60] and the references therein. Particularly, Morshed *et. al* [60] investigated the acceleration scheme of Nesterov in the SKM algorithm for $\delta = 1$.

In this work, we develop a generalized framework namely the GSKM method that extends the SKM algorithm and proves the existence of a family of SKM type methods for solving LF problems. This general framework will provide an ideal platform for the researchers to experiment with a wide range of iterative projection methods and to design efficient algorithms for solving optimization problems in areas like artificial intelligence, machine learning, data mining, and engineering. In addition to the general framework, we propose a Nesterov type acceleration scheme in the SKM method ($0 < \delta < 2$) that outperforms state-of-the-art methods in terms of computation time and solution quality. With the convergence analysis of the GSKM algorithm, we synthesize the convergence analysis of SKM type methods into one convergence theorem from which one can derive convergence results of RK, MM and SKM methods. We also prove convergence of the average iterate (i.e., Cesaro average) generated by both GSKM and PASKM method. We prove sub-linear convergence rate for the Cesaro average under somewhat weaker conditions. We carry out thorough numerical experiments to show the effectiveness of the proposed methods in comparison with state-of-the-art methods for solving a wide range of linear feasibility test instances. Although the proposed methods deal with the case of linear feasibility problem with systems of inequalities, it can be noted that with some modification, like the one stated in the work of Lewis *et. al* [2], one can apply this method to linear systems with both equality and inequality constraints.

The remainder of the paper is organized as follows. The proposed algorithms are discussed in section 2, and the convergence analysis of the proposed algorithms is given in section 3. In section 4, we perform extensive numerical experiments on artificial and real test instances for a better understanding of the behavior of the proposed generalized and accelerated schemes. Besides, we compared the effectiveness of the proposed acceleration schemes with state-of-the-art techniques (i.e., SKM, IPM and ASM). And finally, the paper is concluded in section 5 with concluding remarks and future research directions.

## 2 Preliminaries & Contributions

In this section, we discuss the SKM algorithm and some preliminary technical tools to analyze the SKM type methods. We first discuss the notations and assumptions that will be used throughout the paper. We then briefly discuss the SKM method along with the expectation induced by the sampling distribution of the SKM method. To make the analysis easier and more formal, we introduce the function $f(x)$. Finally, we conclude the section with the proposed GSKM method and the PASKM method and their geometric interpretations.

### 2.1 Notation

We follow the standard linear algebra notation in this work. $\mathbb{R}^n$ denotes the $n$ dimensional real space, $\mathbb{R}^{m \times n}$ denotes the set of $m \times n$ real-valued matrices. For any matrix $A \in \mathbb{R}^{m \times n}$, $A^T$ denotes the transpose matrix $A$ and $a_i^T$ for $i = 1, 2, .., m$ denotes the rows of matrix $A$. Furthermore, $P = \{x \in \mathbb{R}^n | Ax \leq b\}$ denotes the feasible region of the feasibility problem and $\mathcal{P}(x)$ denotes the projection of $x \in \mathbb{R}^n$ onto the feasible region $P$. The notation $d(x, P)$ denotes the distance between $x \in \mathbb{R}^n$ and the feasible region $P$, i.e., $d(x, P) = \inf_{z \in P} \|x - z\| = \|x - \mathcal{P}(x)\|$. For any matrix $A$, the spectral norm and Frobenius norm are denoted by $\|A\|$ and $\|A\|_F$, respectively. For any function $f : X \mapsto Y$, we use $\nabla f$ to represent the gradient of $f$. Finally, $\langle x, y \rangle = x^T y$ denotes the standard inner product and $\|x\| = \sqrt{\langle x, x \rangle}$ as the euclidean ($L_2$) norm. The notation $x^+$ denotes the positive part of any real number (ie., $x^+ = \max\{x, 0\}$). For any two arbitrary matrices $M$, $N$, the notation $M \succ N$ implies the positive definiteness of the matrix $M - N$. The notation $\mathbb{E}_\mathbb{S}[\cdot]$ is used to denote the expectation with respect to the sampling distribution $\mathbb{S}$.

## 2.2 Assumptions

Throughout the paper, we assume that the system $Ax \leq b$ is consistent and the matrix $A$ has no zero rows. We also assumed that the rows of matrix $A$ are normalized (i.e., $\|a_i\|^2 = 1$ for all $i$). Note that, normalization simplifies the convergence analysis considerably. The normalization doesn't impact the computational time significantly (we could simply normalize each row for the first time it occurs during the computation). Moreover, normalization simplifies the convergence analysis considerably. In the description of algorithms, we do not enforce the assumption. Furthermore, it can be noted that the proposed algorithms generate the same iterates irrespective of normalization.

## 2.3 Sampling Kaczmarz Motzkin

The SKM method (Algorithm 1) for solving LF problems, proposed by De Loera *et. al* [11], combines the ideas of both Kaczmarz and Motzkin method. The authors provided a generalized convergence Theorem and a certificate of feasibility which synthesizes the convergence analysis of the Kaczmarz method and Motzkin method for solving LF problems. The proposed method requires only $O(n)$ memory storage and is much more efficient than the state-of-the-art techniques such as Kaczmarz type methods, IPMs and ASMs. The main advantage of SKM can be ascribed to its innovative way of projection plane selection. The hyper-plane selection goes as follows: at iteration $k$ the SKM algorithm selects a collection of $\beta$ rows namely $\tau_k$ uniformly at random out of $m$ rows of the constraint matrix $A$, then out of these $\beta$ rows the row with maximum positive residual is selected (i.e., choose row $i^*$ as $i^* = \arg\max_{i \in \tau_k} \{a_i^T x_k - b_i, 0\}$) and finally the next point $x_{k+1}$ is updated as follows

$$x_{k+1} = x_k - \delta \frac{\left(a_{i^*}^T x_k - b_{i^*}\right)^+}{\|a_{i^*}\|_2^2} a_{i^*}. \tag{4}$$

For ease of analysis, we denote the above sampling distribution as $\mathbb{S}_k$ at iteration $k$, i.e., at each iteration $k$ choose $\tau_k \sim \mathbb{S}_k$ and denote $i^*$ as $i^* = \arg\max_{i \in \tau_k \sim \mathbb{S}_k} \left(a_i^T x_k - b_i\right)^+$.

---

**Algorithm 1** SKM Algorithm: $x_{k+1} = \textbf{SKM}(A, b, x_0, K, \delta, \beta)$

---

Initialize $k \leftarrow 0$;
**while** $k \leq K$ **do**
    Choose a sample of $\beta$ constraints, $\tau_k$, uniformly at random from the rows of matrix $A$.
    From these $\beta$ constraints, choose $i^* = \arg\max_{i \in \tau_k} \{a_i^T x_k - b_i, 0\}$;
    Update $x_{k+1} = x_k - \delta \frac{\left(a_{i^*}^T x_k - b_{i^*}\right)^+}{\|a_{i^*}\|^2} a_{i^*}$;
    $k \leftarrow k + 1$;
**end while**
**return** $x$

---

The SKM method generalizes RK and MM, and it also combines their strength in choosing a constraint at each iteration. It has a cheaper per iteration cost compared to Motzkin's method and converges faster compared to the Kaczmarz method. Several extensions of the SKM method in terms of acceleration [60], improved rate [61] have been proposed recently.

## 2.4 Expectation

For the convergence analysis of Algorithm 1 and its variations (any algorithm that uses that specific type of sampling distribution), we need to discuss a specific expectation calculation. First of all, let us sort the residual vector $(Ax - b)^+$ from smallest to largest for any iterate $x$ and denote $(Ax - b)_{\underline{\mathbf{i_j}}}^+$ as the $(\beta + j)^{th}$ entry on the sorted list [1], i.e.,

$$\underbrace{(Ax - b)_{\underline{\mathbf{i_0}}}^+}_{\beta^{th}} \leq ... \leq \underbrace{(Ax - b)_{\underline{\mathbf{i_j}}}^+}_{(\beta+j)^{th}} \leq ... \leq \underbrace{(Ax - b)_{\underline{\mathbf{i_{m-\beta}}}}^+}_{m^{th}}. \tag{5}$$

Now, consider the list with all of the entries of the residual vector $(Ax - b)^+$, then we need to calculate the probability that particular entry of the residual vector is selected at any given iteration. Note that, the probability that any sample is selected is $\frac{1}{\binom{m}{\beta}}$ and each sample has an equal probability of selection. Another intersecting fact can be noted that the

---

[1]We use the notation $(Ax - b)_{\underline{\mathbf{i_j}}}^+$ throughout the paper to express the underlying expectation, where the indices $\underline{\mathbf{i_j}}$ represent the sampling process of 5.

size of the residual list controls the order and frequency that each entry of the residual vector will be expected to be selected. From now on, we will denote this specific choice of sampling distribution as $\mathbb{S}$ for any point $x \in \mathbb{R}^n$ [2]. To calculate the resulting expectation with respect to the above-mentioned sampling distribution, let us first denote, $\tau \sim \mathbb{S}$ as the set of sampled $\beta$ constraints and $i^*$ as [3]

$$i^* = \arg\max_{i \in \tau \sim \mathbb{S}}\{a_i^T x - b_i, 0\} \;=\; \arg\max_i (A_\tau x - b_\tau)_i^+, \tag{6}$$

where, $A_\tau$ denotes the collection of rows of $A$ restricted to the index set $\tau$ and $(A_\tau x - b_\tau)_i$ denotes the $i^{th}$ entry of $A_\tau x - b_\tau$. Using the above discussion with the list provided in equation (5), we have the following:

$$\mathbb{E}_{\mathbb{S}}\left[\left|(a_{i^*}^T x - b_{i^*})^+\right|^2\right] = \frac{1}{\binom{m}{\beta}} \sum_{j=0}^{m-\beta} \binom{\beta - 1 + j}{\beta - 1} \left|(Ax - b)_{\mathbf{i_j}}^+\right|^2, \tag{7}$$

where, $\mathbb{E}_{\mathbb{S}}$ denotes the required expectation corresponding to the sampling distribution $\mathbb{S}$. The above expectation calculation was first used by De Loera *et.al* in their work [11] where they first introduced the SKM method.

## 2.5 Function $f(x)$

In this section, we formalize the definition of function $f : \mathbb{R}^n \to \mathbb{R}$. Throughout section 3, we will use the properties of function $f(x)$ [4]. First, for any index $i$, let us define the following function

$$f_i(x) = \frac{1}{2}\left|(a_i^T x - b_i)^+\right|^2, \quad \nabla f_i(x) = (a_i^T x - b_i)^+ a_i. \tag{8}$$

Then to simplify the expectation expression of (7) further, we define the function $f$ and the gradient of $f$ as follows:

$$f(x) = \mathbb{E}_{\mathbb{S}}\left[f_{i^*}(x)\right], \quad \nabla f(x) = \mathbb{E}_{\mathbb{S}}\left[\nabla f_{i^*}(x)\right], \tag{9}$$

where, the index $i^*$ is selected by the rule provided in (6).

## 2.6 Contributions

***Generalized Sampling Kaczmarz Method (GSKM).*** For obtaining a generalized version of the SKM method, we suggest using history information in updating the current update. In particular, we take two random iterates $x_{k-1}$ and $x_k$ generated by successive SKM iteration and then update the next iterate $x_{k+1}$ as an affine combination of the previous two updates. Starting with $x_0 = x_1 \in \mathbb{R}^n$, for $k \geq 1$, we update

$$x_{k+1} = (1 - \xi)z_k + \xi z_{k-1},$$

where $z_k = x_k - \delta \frac{\left(a_{i^*}^T x_k - b_{i^*}\right)^+}{\|a_{i^*}\|^2} a_{i^*}$ is the $k^{th}$ update of the SKM algorithm. Note that, by taking $\xi = 0$, one can recover the original SKM algorithm. For simpler representation, we denote this method as a generalized SKM method or GSKM method. GSKM method is formally provided in Algorithm 2 and the convergence analysis is provided in subsection 3.2. Our convergence analysis suggests that for any $0 < \delta < 2$, one could choose any $\xi$ such that $\xi \in Q$ [5].

**Table 1:** Algorithms & convergence results obtained from GSKM.

| Parameters, $\beta$, $\delta$, $\xi$ | Row selection Rule, $(i^*)$ | Convergence Rate | Algorithm |
|---|---|---|---|
| $\beta = 1,\ \delta = 1,\ \xi = 0$ | $\mathbb{P}(i^*) = \frac{\|a_i\|^2}{\|A\|_F^2}$ | $\mathbb{E}\left[r_k^2\right] \leq \left(1 - \frac{\lambda_{\min}}{\|A\|_F^2}\right)^k r_0^2$ | RK [1] |
| $\beta = m,\ \delta = 1,\ \xi = 0$ | $i^* = \arg\max_j e_j(x_{k-1})$ | $r_k^2 \leq \left(1 - \frac{\lambda_{\min}}{m}\right)^k r_0^2$ | MM [22] |
| $0 < \delta < 2,\ \xi = 0$ | $\tau_k \sim \mathbb{S}_k$ <br> $i^* = \arg\max_{j \in \tau_k} e_j(x_{k-1})$ | $\mathbb{E}\left[r_k^2\right] \leq \left(1 - \frac{\eta}{mL^2}\right)^k r_0^2$ | SKM [11] |

---

[2] For ease of notation, throughout the paper, we will use $\mathbb{S}_k$ to denote the sampling distribution corresponding to any random iterate $x_k \in \mathbb{R}^n$

[3] Similarly, we will use $\tau_k \sim \mathbb{S}_k$ to denote the sampled set and $i^* = \arg\max_{i \in \tau_k \sim \mathbb{S}_k}\{a_i^T x_k - b_i, 0\} = \arg\max_{i \in \tau_k \sim \mathbb{S}_k}(A_{\tau_k} x_k - b_{\tau_k})_i^+$ for any iterate $x_k \in \mathbb{R}^n$.

[4] Similar type of functions with uniform sampling have been studied in [10] [12] in the context of stochastic gradient descent and alternating projection algorithms respectively.

[5] see (26).

---

**Algorithm 2** GSKM Algorithm: $x_{k+1} = \textbf{GSKM}(A, b, x_0, K, \delta, \beta, \xi)$

---

Choose $0 < \delta < 2$, $\xi \in Q$

Initialize $x_1 = x_0$, $z_1 = z_0$, $k = 0$;

**while** $1 \leq k \leq K$ **do**

Choose a sample of $\beta$ constraints, $\tau_k$, uniformly at random from the rows of matrix $A$. From these $\beta$ constraints, choose $i^* = \arg\max_{i \in \tau_k}\{a_i^T x_k - b_i, 0\}$ and update,

$$z_k = x_k - \delta\frac{\left(a_{i^*}^T x_k - b_{i^*}\right)^+}{\|a_{i^*}\|^2}a_{i^*}; \tag{10}$$

$$x_{k+1} = (1 - \xi)z_k + \xi z_{k-1}; \tag{11}$$

$k \leftarrow k + 1$;

**end while**

**return** $x$

---

In Table 1, we list the algorithms and their respective convergence Theorems recovered from the GSKM algorithm with different parameter choices. To simplify the notation, we denote, $r_k = d(x_k, P)$, $\eta = 2\delta - \delta^2$, $\lambda_{\min} = \lambda_{\min}^+(A^T A)$, $e_j(x) = a_j^T x - b_j$.

***Probably Accelerated Sampling Kaczmarz Method (PASKM)***. We propose an accelerated randomized projection method based on the SKM method and Nesterov accelerated gradient (NAG). Note that, NAG generates sequences $\{y_k\}$ and $\{v_k\}$ using the following update formulas:

$$y_k = \alpha_k v_k + (1 - \alpha_k)x_k, \quad x_{k+1} = y_k - \theta_k \nabla f(y_k),$$
$$v_{k+1} = \omega_k v_k + (1 - \omega_k)y_k - \gamma_k \nabla f(y_k). \tag{12}$$

In equation (12), $\nabla f$ is the gradient of the given function and $\alpha_k, \omega_k, \theta_k$ are the step sequences. Nesterov used updated values for the sequences $\alpha_k, \omega_k, \theta_k$ and obtained a better convergence rate for the acceleration of standard gradient descent. There are two available works directly involve applying Nesterov's acceleration in Kaczmarz type methods [6], first one is by Wright *et. al* [32] where the accelerated RK method is proposed for linear systems, the second one deals with applying acceleration in SKM for $\delta = 1$ [60].

---

**Algorithm 3** PASKM Algorithm: $x_{k+1} = \textbf{PASKM}(A, b, x_0, K, \delta, \beta)$

---

Initialize $v_0 \leftarrow x_0$, $k \leftarrow 0$;

**while** $k \leq K$ **do**

Choose $\gamma, \omega, \alpha$ considering either (28) or (32) and update

$$y_k = \alpha v_k + (1 - \alpha)x_k; \tag{13}$$

Choose a sample of $\beta$ constraints, $\tau_k$, uniformly at random from the rows of matrix $A$. From these $\beta$ constraints, choose $i^* = \arg\max_{i \in \tau_k}\{a_i^T y_k - b_i, 0\}$; Update

$$x_{k+1} = y_k - \delta\frac{\left(a_{i^*}^T y_k - b_{i^*}\right)^+}{\|a_{i^*}\|^2}a_{i^*}; \tag{14}$$

$$v_{k+1} = \omega v_k + (1 - \omega)y_k - \gamma\frac{\left(a_{i^*}^T y_k - b_{i^*}\right)^+}{\|a_{i^*}\|^2}a_{i^*}; \tag{15}$$

$k \leftarrow k + 1$;

**end while**

**return** $x$

---

In this work, we consider the general case $0 < \delta < 2$ and develop a probably accelerated scheme for the SKM algorithm. The main difference between the proposed PASKM algorithm and the above-mentioned method is the choice of step

---

[6]Recently, heavy ball momentum method has been proposed in the context of SKM method [62]
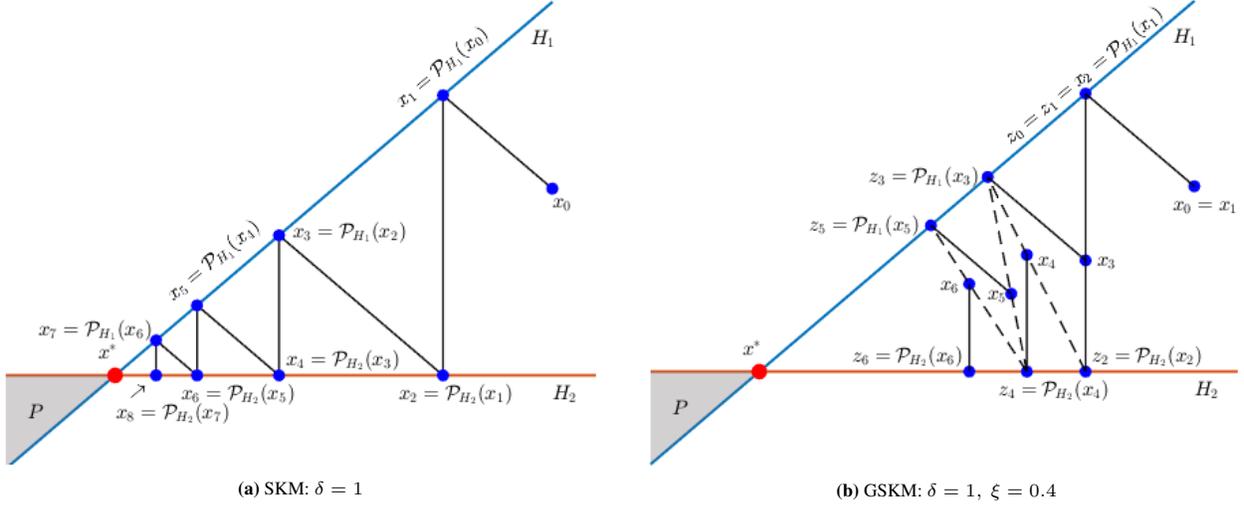
sequences. We propose to use precomputed values for the parameters $\omega, \gamma, \alpha$ for every iterate compared to the iterative parameter selection process in [32, 56, 60]. Now, using the definition of function $\hat{f}_i$ (see (8)) in (12), we derive the following scheme:

$$
\begin{aligned}
y_k &= \alpha v_k + (1 - \alpha)x_k, \quad x_{k+1} = y_k - \delta \nabla f_{i^*}(x), \\
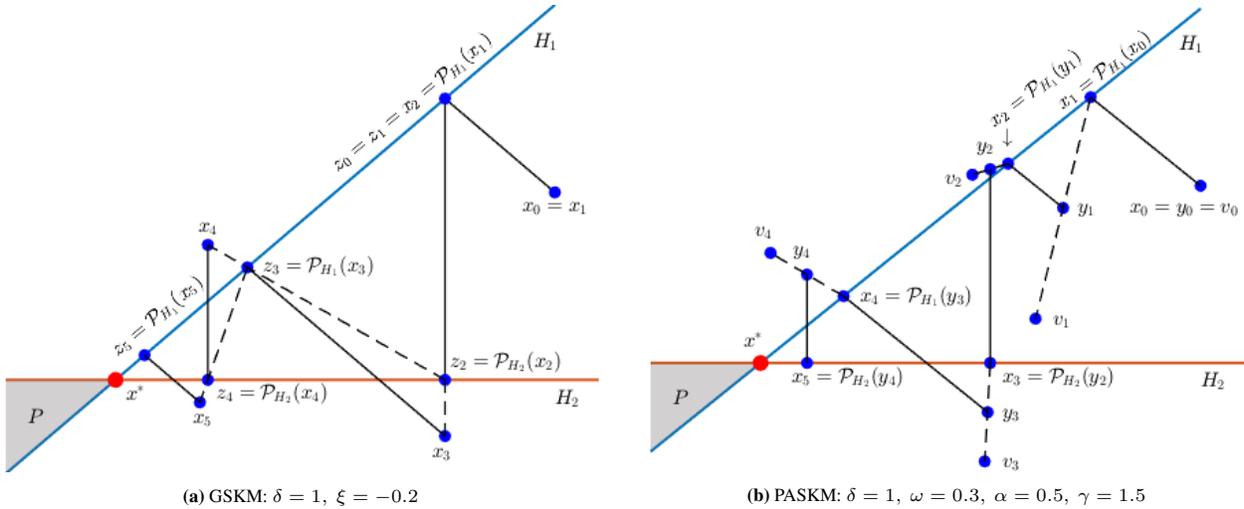v_{k+1} &= \omega v_k + (1 - \omega)y_k - \gamma \nabla f_{i^*}(x),
\end{aligned}
\tag{16}
$$

with $i^*$ chosen as $i^* = \arg\max_{j \in \tau_k} e_j(x_{k-1})$, where $\tau_k \sim \mathbb{S}_k$. The PASKM method is formalized as Algorithm 3 and the detailed convergence analysis of the method is provided in Section 3. This method generally outperforms both the SKM and GSKM algorithms for almost all of the test instances considered in this work (see Section 4).

## 2.7 Geometric Interpretation

The goal of this section is to provide a geometric interpretation of the proposed GSKM and PASKM methods. We shed more lights on how the proposed algorithms work in practice and the difference among SKM, GSKM and PASKM methods.



**(a)** SKM: $\delta = 1$

**(b)** GSKM: $\delta = 1$, $\xi = 0.4$

**Figure 1:** Graphical interpretation of the SKM method and the GSKM method with only two hyper-planes $H_j = \{x | a_j^T x \le b_j\}$



**(a)** GSKM: $\delta = 1$, $\xi = -0.2$

**(b)** PASKM: $\delta = 1$, $\omega = 0.3$, $\alpha = 0.5$, $\gamma = 1.5$

**Figure 2:** Graphical interpretation of the GSKM method and the PASKM method with only two hyper-planes $H_j = \{x | a_j^T x \le b_j\}$

In Figures 1 and 2, we illustrate the differences among SKM, GSKM and PASKM methods in an $\mathbb{R}^2$ plane. Our goal is to show how each of the proposed algorithms progress at each iteration. For illustration purposes, We performed the experiment with only two hyper-planes and the selection of hyper-planes is done in an alternative fashion. The

notation $\mathcal{P}_{H_1}(x)$ denotes the orthogonal projection of point $x$ onto the hyper-plane $H_1$. For comparison purposes, we started with the same starting point $x_0$ and drew the figures with the same scaling. For any given starting point $x_0$, each algorithm projects the point onto the most violated constraint from the sampled constraint set.

The projection step corresponds to the computation of the term $x_k - \delta \frac{(a_{i^*}^T x_k - b_{i^*})^+}{\|a_{i^*}\|^2} a_{i^*}$, which means that the current update $x_k$ is projected onto the violated hyper-plane. The projection parameter $\delta \in (0, 2]$ defines the type of projection. When $\delta = 1$, the projection is exact, that is the point $\mathcal{P}_H(x_k)$ belongs to the hyper-plane $H$. GSKM ($0 \leq \xi \leq 1$) can be seen as a kind of convex projection update which is slower compared to SKM. From Figure 2, it can be seen that the GSKM method with $-1 < \xi < 0$ proceeds faster compared to SKM and it requires an affine combination of the previous two successive projections (i.e., $z_{k-1}$ and $z_k$). Compared to SKM and GSKM, the PASKM method updates three different sequences $x_k, v_k, y_k$. From Figure 2, it can be noted that GSKM with negative $\xi$ and PASKM moves faster to the feasible region $P$ compared to the SKM method (later in the numerical section this comparison will become much more apparent for larger test instances).

## 2.8 Connection between GSKM and PASKM

Assume, $-1 < \xi \leq 0$. Then, we can simplify the update formula of the GSKM method as
$$x_{k+1} = (1 - \xi)x_k + \xi x_{k-1} - \delta(1 - \xi)\nabla f_{i^*}(x_k) - \delta\xi\nabla f_{j^*}(x_{k-1}), \tag{17}$$
where the indices $i^*$ and $j^*$ are selected following the rule of (8) for the iterate $x_k$ and $x_{k-1}$, respectively. Furthermore, take $\omega(1 - \alpha) = -\xi$ and $\gamma$ such that the condition $\alpha\gamma = \delta(1 - \xi)$ holds, then from the update formula of the PASKM method we get,
$$v_{k+1} \overset{(15)}{=} \omega v_k + (1 - \omega)y_k - \gamma\nabla f_{i^*}(y_k) \overset{(13)}{=} \left(1 - \frac{\xi}{\alpha}\right)y_k + \frac{\xi}{\alpha}x_k - \gamma\nabla f_{i^*}(y_k).$$

Similarly, from the definition of $y_{k+1}$, we have
$$\begin{aligned} y_{k+1} = \alpha v_{k+1} + (1 - \alpha)x_{k+1} &= (1 - \xi)y_k + \xi x_k - [\alpha\gamma + \delta(1 - \alpha)]\nabla f_{i^*}(y_k) \\ &= (1 - \xi)y_k + \xi [y_{k-1} - \delta\nabla f_{j^*}(y_{k-1})] - [\delta(1 - \xi) + \delta(1 - \alpha)]\nabla f_{i^*}(y_k) \\ &= (1 - \xi)y_k + \xi y_{k-1} - \delta(1 - \xi)\nabla f_{i^*}(y_k) - \delta\xi\nabla f_{j^*}(y_{k-1}), \end{aligned} \tag{18}$$
where the indices $i^*$ and $j^*$ are selected following the rule of (8) for the iterate $y_k$ and $y_{k-1}$, respectively. Considering update formulas (17) and (18), we can conclude that if the conditions $0 \leq \omega(1 - \alpha) = -\xi < 1$ and $\alpha\gamma = \delta(1 - \xi)$ hold, then the sequence $x_k$ generated by the GSKM algorithm and the sequence $y_k$ generated by the PASKM algorithm is the same sequence.

# 3 Main Results

In this section, we present the convergence analysis of the proposed algorithms. In the first subsection, we provided the necessary technical Lemmas & Theorems that will be used later for our convergence analysis. In the second subsection, we provided the convergence Theorems of the GASKM algorithm. Finally, the last subsection deals with the convergence analysis of the PASKM method.

## 3.1 Technical Tools

In this subsection, we will discuss two types of results. Most of the results derived are related to the properties of the function $f(x)$. Lemma 3.1 is the famous result of Hoffman regarding the linear system of inequalities. Lemmas 3.6-3.9 discuss the strong convexity and existence of Lipschitz constant along some restricted segment. Finally, Theorems 3.12 and 3.13 deal with developing decay bounds for some non-negative sequences. We will use Lemmas 3.6-3.9 frequently in our convergence analysis. Theorems 3.12 and 3.13 will be used to derive the proposed convergence bounds of the quantities $\mathbb{E}[d(x_k, P)]$ and $\mathbb{E}[d(x_k, P)^2]$.

**Lemma 3.1.** *(Hoffman [63], Theorem 4.4 in [2]) Let $x \in \mathbb{R}^n$ and $P$ be the feasible region, then there exists a constant $L > 0$ such that the following identity holds:*
$$d(x, P)^2 \leq L^2 \|(Ax - b)^+\|^2.$$

The constant $L$ is the so-called Hoffman constant. Note that, for a consistent system of equations (i.e., there exists a unique $x^*$ such that $Ax = b$), $L$ can be expressed in terms of the smallest singular value of matrix $A$, i.e.,
$$L^2 = \frac{1}{\|A^{-1}\|^2} = \frac{1}{\lambda_{min}^+(A^T A)}.$$

**Lemma 3.2.** *(Lemma 2.1 in [11]) Let $\{x_k\}$, $\{y_k\}$ be real non-negative sequences such that $x_{k+1} > x_k > 0$ and $y_{k+1} \geq y_k \geq 0$, then*

$$\sum_{k=1}^{n} x_k y_k \geq \sum_{k=1}^{n} \overline{x} y_k, \quad \text{where } \overline{x} = \frac{1}{n} \sum_{k=1}^{n} x_k.$$

**Lemma 3.3.** *For any $x \in \mathbb{R}^n$ and $\bar{x} \in P$, the following identity holds,*

$$d(x, P)^2 = \|x - \mathcal{P}(x)\|^2 \leq \|x - \bar{x}\|^2.$$

**Lemma 3.4.** *Let $\lambda_j$ be the $j^{th}$ eigenvalue of the matrix $W = \mathbb{E}_\mathbb{S}\left[a_{i*}a_{i*}^T\right]$, then for all $j$, the bound $0 \leq \lambda_j \leq 1$ holds.*

*Proof.* Since $W$ is positive semi-definite, we can write $\lambda_j \geq 0$ for all $j$. Also as the mapping $\lambda_{\max}(X)$ is convex, using Jensen's inequality we have,

$$\lambda_{\max}(W) = \lambda_{\max}\left[\mathbb{E}_\mathbb{S}\left[a_{i*}a_{i*}^T\right]\right] \leq \mathbb{E}_\mathbb{S}\left[\lambda_{\max}\left(a_{i*}a_{i*}^T\right)\right] \leq 1.$$

$\square$

**Lemma 3.5.** *For any $1 \leq \beta \leq m$, we have the following:*

$$\mathbb{E}_\mathbb{S}\left[a_{i*}a_{i*}^T\right] \preceq \frac{\beta}{m} A^T A.$$

*Proof.* See Appendix 1. $\square$

**Lemma 3.6.** *For any $x \in \mathbb{R}^n$ with $\lambda_{\max} = \lambda_{\max}(A^T A)$, we have the following:*

$$\frac{\mu_1}{2} d(x, P)^2 \leq f(x) \leq \frac{\mu_2}{2} d(x, P)^2,$$

*with $0 < \mu_1 = \frac{1}{mL^2} \leq \mu_2 = \min\left\{1, \frac{\beta}{m}\lambda_{\max}\right\} \leq 1$.*

*Proof.* See Appendix 1. $\square$

Lemma 3.6 states that the function $f$ is strongly convex with constant $\mu_1$ and has Lipschitz continuous gradient with constant $\mu_2$ when restricted along the segment $[x, \mathcal{P}(x)]$. Let, $f^* = \min_x f(x)$, then it can be easily checked that $f^* = f(x^*) = 0$. Here, $x^*$ is the optimal solution and it satisfies $Ax^* \leq b$. Moreover, the point $\mathcal{P}(x)$ satisfies the condition $\nabla f(\mathcal{P}(x)) = 0$. Then we rewrite the inequalities of Lemma 3.6 as follows

$$\frac{\mu_1}{2} \|x - \mathcal{P}(x)\|^2 + \langle \nabla f(\mathcal{P}(x)), x - \mathcal{P}(x) \rangle \leq f(x) - f^*, \tag{19}$$

$$f(x) - f^* \leq \langle \nabla f(\mathcal{P}(x)), x - \mathcal{P}(x) \rangle + \frac{\mu_2}{2} \|x - \mathcal{P}(x)\|^2. \tag{20}$$

Here, equation (19) and (20) represent the Lipschitz continuity condition and the strong convexity condition respectively along the line segment $[x, \mathcal{P}(x)]$. For our convergence analysis of Algorithm 2 and 3, we will need inequalities like (19) and (20) along the segment $[x, y]$ for any $x, y \in \mathbb{R}^n$. Following two Lemmas deal with the problem of finding such bounds.

**Lemma 3.7.** *For any $x, y \in \mathbb{R}^n$, we have the following:*

$$\langle x - y, \mathbb{E}_\mathbb{S}\left[(a_{i*}^T y - b_{i*})^+ a_{i*}\right] \rangle = \langle x - y, \nabla f(y) \rangle$$
$$\leq f(x) - f(y) \leq \frac{\mu_2}{2} d(x, P)^2 - \frac{\mu_1}{2} d(y, P)^2.$$

*Proof.* See Appendix 1. $\square$

**Remark 3.8.** *We note that the condition of Lemma 3.7 is weaker than the traditional strong convexity, and it is also weaker than the essentially strong convexity condition defined in [64]. For instance, the essentially strong convexity requires the following identity:*

$$f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle - \frac{\epsilon}{2} \|x - y\|^2, \quad \forall x, y, \text{ s.t. } \mathcal{P}(x) = \mathcal{P}(y),$$

*for some $\epsilon > 0$. The above condition clearly implies (33). Moreover, the restricted secant inequality condition defined in [64] can be written as*

$$\langle \nabla f(x), x - \mathcal{P}(x) \rangle \geq \epsilon \|x - \mathcal{P}(x)\|^2. \tag{21}$$

*Note that, with the choice $x = \mathcal{P}(y)$ in Lemma 3.7, we have the following:*

$$\langle \nabla f(y), y - \mathcal{P}(y) \rangle \geq \frac{\mu_2}{2} \|y - \mathcal{P}(y)\|^2.$$

*Here, we used the fact $d(\mathcal{P}(y), P)^2 = \|\mathcal{P}(y) - \mathcal{P}(y)\|^2 = 0$. This implies that the function $f(x)$ satisfies the restricted secant inequality condition of (21) with constant $\epsilon = \frac{\mu_1}{2}$. Indeed it can be shown that the constant $\epsilon = \frac{\mu_1}{2}$ can be improved further (see the following Lemma).*

**Lemma 3.9.** *For any $y \in \mathbb{R}^n$ and $\bar{y}$ such that $A\bar{y} \leq b$, we have the following:*

$$\langle \bar{y} - y, \mathbb{E}_{\mathbb{S}} \left[ a_{i^*} (a_{i^*}^T y - b_{i^*})^+ \right] \rangle = \langle \bar{y} - y, \nabla f(y) \rangle \leq -2f(y) \leq -\mu_1 d(y, P)^2.$$

*Proof.* See Appendix 1. □

**Remark 3.10.** *Substituting $\bar{y} = \mathcal{P}(y)$, in Lemma 3.9 we have,*

$$\langle \mathcal{P}(y) - y, \mathbb{E}_{\mathbb{S}} \left[ a_{i^*} (a_{i^*}^T y - b_{i^*})^+ \right] \rangle \leq -2f(y) \leq -\mu_1 d(y, P)^2.$$

*Note that, similar types of results can be found in the literature. For instance, in [12], authors obtained similar result with respect to a different expectation, they used $\mathbb{E}[x] = \frac{1}{n} \sum_i x_i$ for any $x \in \mathbb{R}^n$, which is commonly used to analyze randomized Kaczmarz type methods (see [1, 2]). Furthermore, we believe a better upper bound than the one obtained in Lemma 3.7 can be obtained considering some restrictions on the data matrix $A$. To that end, one needs to obtain a better version of equation (33), i.e., one needs to show that the function $f(x) - \frac{\epsilon}{2}\|x\|^2$ is convex along the line segment $[x, y]$.*

**Lemma 3.11.** *For any $x \in \mathbb{R}^n$ and $0 < \delta < 2$, we have the following:*

$$\mathbb{E}_{\mathbb{S}} \left[ d(z, P)^2 \right] = \mathbb{E}_{\mathbb{S}} \left[ \left\| x - \mathcal{P}(x) - \delta \left( a_{i^*}^T x - b_{i^*} \right)^+ a_{i^*} \right\|^2 \right] \leq h(\delta) \, d(x, P)^2,$$

*where, $z = x - \delta \left( a_{i^*}^T x - b_{i^*} \right)^+ a_{i^*}$, $\eta = 2\delta - \delta^2$ and $h(\delta) = 1 - \eta \mu_1 < 1$.*

*Proof.* See Appendix 1. □

Before we delved into the main Theorems, for any $\phi_1, \phi_2 \geq 0$, let us define the following parameters:

$$\phi = \frac{-\phi_1 + \sqrt{\phi_1^2 + 4\phi_2}}{2}, \quad \rho = \phi + \phi_1,$$

$$R_1 = \frac{1 + \phi}{\phi + \rho}, \; R_2 = \frac{1 - \rho}{\phi + \rho}, \; R_3 = \frac{\rho + \phi_2}{\phi + \rho}, \; R_4 = \frac{\phi - \phi_2}{\phi + \rho}. \tag{22}$$

The following two Theorems deal with the growth of non-negative real sequences. We will use these results in our main analysis of GSKM and PASKM method.

**Theorem 3.12.** *Let $\{G_k\}$ be a non-negative real sequence satisfying the following relation:*

$$G_{k+1} \leq \phi_1 G_k + \phi_2 G_{k-1}, \; \forall k \geq 1 \quad G_0 = G_1 \geq 0,$$

*if $\phi_1, \phi_2 \geq 0$ and $\phi_1 + \phi_2 < 1$ then the following bounds hold:*

  *1. (Lemma 9 in [57]) Let, $\phi$ be the largest root of $\phi^2 + \phi_1 \phi - \phi_2 = 0$, then*

$$G_{k+1} \leq (1 + \phi)(\phi + \phi_1)^k \, G_0, \; \forall k \geq 1.$$

2. *Define $\rho = \phi + \phi_1$, then we have the following:*

$$\begin{bmatrix} G_{k+1} \\ G_k \end{bmatrix} \leq \begin{cases} \begin{bmatrix} R_1\rho^{k+1} + R_2\phi^{k+1} \\ R_1\rho^k - R_2\phi^k \end{bmatrix} G_0 & k \text{ even}; \\ \begin{bmatrix} R_3\rho^k - R_4\phi^k \\ R_3\rho^{k-1} + R_4\phi^{k-1} \end{bmatrix} G_0 & k \text{ odd}, \end{cases}$$

*where, $0 \leq \phi < 1$ and $0 < \rho = \phi + \phi_1 < 1$.*

*Proof.* See Appendix 1. □

**Theorem 3.13.** *Let the real sequences $H_k \geq 0$ and $F_k \geq 0$ satisfy the following recurrence relation:*

$$\begin{bmatrix} H_{k+1} \\ F_{k+1} \end{bmatrix} \leq \begin{bmatrix} \Pi_1 & \Pi_2 \\ \Pi_3 & \Pi_4 \end{bmatrix} \begin{bmatrix} H_k \\ F_k \end{bmatrix}, \tag{23}$$

*where, $\Pi_1, \Pi_2, \Pi_3, \Pi_4 \geq 0$ such that the following relations*

$$\Pi_1\Pi_4 - \Pi_2\Pi_3 \geq 0, \qquad \Pi_1 + \Pi_4 < 1 + \min\{1, \Pi_1\Pi_4 - \Pi_2\Pi_3\}, \tag{24}$$

*hold. Then the sequence $\{H_k\}$ and $\{F_k\}$ converges and the following result holds:*

$$\begin{bmatrix} H_{k+1} \\ F_{k+1} \end{bmatrix} \leq \begin{bmatrix} \Pi_1 & \Pi_2 \\ \Pi_3 & \Pi_4 \end{bmatrix}^k \begin{bmatrix} H_1 \\ F_1 \end{bmatrix} = \begin{bmatrix} \Gamma_2\Gamma_3(\Gamma_1 - 1)\,\rho_1^k + \Gamma_1\Gamma_3(\Gamma_2 + 1)\,\rho_2^k \\ \Gamma_3(\Gamma_1 - 1)\,\rho_1^k + \Gamma_3(\Gamma_2 + 1)\,\rho_2^k \end{bmatrix} \begin{bmatrix} H_1 \\ F_1 \end{bmatrix}.$$

*where,*

$$\Gamma_1 = \frac{\Pi_1 - \Pi_4 + \sqrt{(\Pi_1 - \Pi_4)^2 + 4\Pi_2\Pi_3}}{2\Pi_3},$$

$$\Gamma_2 = \frac{\Pi_1 - \Pi_4 - \sqrt{(\Pi_1 - \Pi_4)^2 + 4\Pi_2\Pi_3}}{2\Pi_3}, \quad \Gamma_3 = \frac{\Pi_3}{\sqrt{(\Pi_1 - \Pi_4)^2 + 4\Pi_2\Pi_3}},$$

$$\rho_1 = \frac{1}{2}\left[\Pi_1 + \Pi_4 - \sqrt{(\Pi_1 - \Pi_4)^2 + 4\Pi_2\Pi_3}\right],$$

$$\rho_2 = \frac{1}{2}\left[\Pi_1 + \Pi_4 + \sqrt{(\Pi_1 - \Pi_4)^2 + 4\Pi_2\Pi_3}\right], \tag{25}$$

*and $\Gamma_1, \Gamma_3 \geq 0$ and $0 \leq \rho_1 \leq \rho_2 < 1$.*

*Proof.* See Appendix 1. □

### 3.2 Convergence Analysis of the GSKM Method

In this subsection, we study convergence properties of the proposed GSKM method, i.e., we study the convergence behavior of the quantities of $\mathbb{E}[\|x_k - \mathcal{P}(x_k)\|]$ and $\mathbb{E}[f(x_k)]$. For any $\xi \in \mathbb{R}$, let us define the sets $Q, Q_1, Q_2$ as

$$Q_1 = \{\xi \mid 0 \leq \xi \leq 1\}, \quad Q = Q_1 \cup Q_2,$$
$$Q_2 = \{-1 < \xi \leq 0 \mid (1 + \xi)\sqrt{h(\delta)} - \xi(1 + \delta\sqrt{\mu_2}) < 1\}. \tag{26}$$

We proved that whenever $\xi \in Q$ and $0 < \delta < 2$, the proposed GSKM method enjoys a global linear rate. We also provided convergence analysis of the function values (i.e., $f(x_k)$) with respect to the Cesaro average. Our results are global in nature and to the best of our knowledge, this is the first of its kind result for the SKM method.

**Theorem 3.14.** *Let $\{x_k\}$ be the sequence of random iterates generated by algorithm 2. With the choice of parameters, $0 < \delta < 2$ and $0 \leq \xi \leq 1$ ($\xi \in Q_1$), the sequence of iterates $\{x_k\}$ converges and the following results hold:*

1. Take $\phi_1 = (1 - \xi)h(\delta)$, $\phi_2 = \xi h(\delta)$ and $\rho, \phi$ as in equation (22), then

$$\mathbb{E}[d(x_{k+1}, P)^2] \leq \rho^k(1 + \phi)d(x_0, P)^2 \text{ and } \mathbb{E}[f(x_k)] \leq \frac{\mu_2(1 + \phi)}{2}\rho^k d(x_0, P)^2.$$

2. Take $\phi_1 = (1 - \xi)h(\delta)$ and $\phi_2 = \xi h(\delta)$, then

$$\mathbb{E}\begin{bmatrix} d(x_{k+1}, P)^2 \\ d(x_k, P)^2 \end{bmatrix} \leq \begin{cases} \begin{bmatrix} R_1\rho^{k+1} + R_2\phi^{k+1} \\ R_1\rho^k - R_2\phi^k \end{bmatrix} d(x_0, P)^2 & k \text{ even}; \\ \begin{bmatrix} R_3\rho^k - R_4\phi^k \\ R_3\rho^{k-1} + R_4\phi^{k-1} \end{bmatrix} d(x_0, P)^2 & k \text{ odd}, \end{cases}$$

where, the constants $R_1, R_2, R_3, R_4$ are defined in equation (22) and $0 \leq \phi, \phi_1, \phi_2 < 1$ and $0 < \rho = \phi + \phi_1 < 1$.

3. Also the average iterate $\tilde{x}_k = \sum_{l=1}^{k} x_l$ for all $k \geq 0$ satisfies the following

$$\mathbb{E}[d(\tilde{x}_k, P)^2] \leq \frac{(1 + \phi)\, d(x_0, P)^2}{k(1 - \rho)} \quad \text{and} \quad \mathbb{E}[f(\tilde{x}_k)] \leq \frac{(1 + \xi)d(x_0, P)^2}{2\delta k(2 - \delta)}.$$

*Proof.* See Appendix 2.

$\square$

In the above Theorem, we obtain a global linear rate for the GSKM method with $0 \leq \xi \leq 1$. Note that, when $0 \leq \xi \leq 1$, we have,

$$\rho = \phi + \phi_1 = \frac{(1 - \xi)h(\delta) + \sqrt{(1 - \xi)^2 h^2(\delta) + 4\xi h(\delta)}}{2}$$
$$\geq \frac{(1 - \xi)h(\delta) + \sqrt{(1 - \xi)^2 h^2(\delta)}}{2} = (1 - \xi)h(\delta).$$

Since the maximum value of $(1 - \xi)h(\delta)$ can be derived as $h(\delta)$, the above inequality attains equality when $\xi = 0$ (see the next Corollary). This gives us $1 > \rho = \phi + \phi_1 \geq h(\delta)$. Since the rate of the SKM algorithm is given by $h(\delta)$, we can say that the theoretical convergence rate of Algorithm 2 is always worse or equal compared to SKM whenever $0 \leq \xi \leq 1$.

**Corollary 3.14.1.** *(Theorem 1.3 in [11]) Let $\{x_k\}$ be the sequence of random iterates generated by the SKM method (algorithm 1) starting with $x_0 \in \mathbb{R}^n$. With $0 < \delta < 2$, the sequence of iterates $\{x_k\}$ converges and the following result holds:*

$$\mathbb{E}\left[d(x_{k+1}, P)^2\right] \leq [h(\delta)]^k\, d(x_0, P)^2.$$

*Proof.* Note that, if we let $\xi = 0$ in the GSKM method, then we have $x_{k+1} = z_k$, which is precisely the SKM method. Now, take $\xi = 0$ in Theorem 3.14, then considering the first part of the Theorem, we have $\rho = h(\delta)$. Furthermore, from the second part, we have $R_3\rho^k - R_4\phi^k = R_1\rho^{k+1} + R_2\phi^{k+1} = \rho^k = (h(\delta))^k$. This proves the result of Corollary 3.14.1 which is precisely the convergence rate obtained in [11] for the SKM method. $\square$

Our next Theorem, states that, for a range of negative values of the parameter $\xi$, the GSKM method enjoys a global linear rate.

**Theorem 3.15.** *Let $\{x_k\}$ be the sequence of random iterates generated by algorithm 2 and let $0 < \delta < 2$ and $\xi \in Q_2$. Define*

$$\Pi_1 = \sqrt{h(\delta)}, \ \Pi_2 = |\xi|, \ \Pi_3 = \delta\sqrt{\mu_2 h(\delta)}, \ \Pi_4 = |\xi|\left(1 + \delta\sqrt{\mu_2}\right), \tag{27}$$
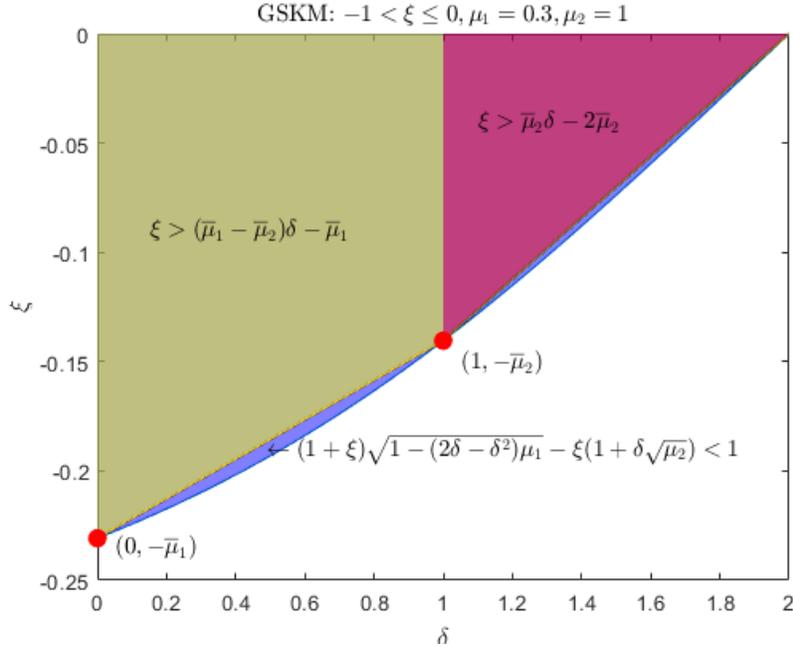
*and $\Gamma_1, \Gamma_2, \Gamma_3, \rho_1, \rho_2$ as in (24) with the parameter choice of (27). Then the sequence of iterates $\{x_k\}$ converges and the following result holds:*

$$\mathbb{E}\begin{bmatrix} d(x_{k+1}, P) \\ \|z_{k+1} - z_k\| \end{bmatrix} \leq \begin{bmatrix} -\Gamma_2\Gamma_3 \, \rho_1^k + \Gamma_1\Gamma_3 \, \rho_2^k \\ -\Gamma_3 \, \rho_1^k + \Gamma_3 \, \rho_2^k \end{bmatrix} d(x_0, P),$$

*where $\Gamma_1, \Gamma_3 \geq 0$ and $0 \leq \rho_1 \leq \rho_2 < 1$.*

*Proof.* See Appendix 2. □

**Parameter Choice for GSKM** Now we discuss allowable parameter selection for the GSKM algorithm based on Theorem 3.14 and 3.15.



**Figure 3:** Allowable parameter range

From Theorem 3.14, it can be noted that the GASKM method will converge for any $0 \leq \xi \leq 1$. Whenever $\xi$ is negative (i.e., $\xi \in Q_2$), the allowable range for $\xi$ can be shown in the following figure. In Figure 3, we plot the feasible region for allowable $\xi$ values for $\mu_1 = 0.3, \mu_2 = 1$. Denote, $\tilde{\mu}_1 = \frac{\mu_1}{\mu_1 + \sqrt{\mu_2}}$ and $\tilde{\mu}_2 = \frac{1 - \sqrt{1 - \mu_1}}{1 - \sqrt{1 - \mu_1} + \sqrt{\mu_2}}$. Then the feasible region of Figure 3 can be approximated piece-wise as $\xi > (\tilde{\mu}_1 - \tilde{\mu}_2)\delta - \tilde{\mu}_1$ for $0 < \delta \leq 1$ and $\xi > \tilde{\mu}_2\delta - 2\tilde{\mu}_2$ for $1 \leq \delta < 2$. Moreover, any $(\xi, \delta)$ pair that resides inside the region $\{0 < \delta < 2, \ -1 < \xi < 0, \ \xi \geq 0.5\tilde{\mu}_1(\delta - 2)\}$ also resides inside the feasible region of Theorem 3.15.

**Cesaro Average:** In the next Theorem, we propose the convergence analysis of the function values $f(x)$, with respect to the Cesaro average. Instead of bounding $\mathbb{E}[f(x_k)]$ in terms of initial function value $f(x_0)$, we bound the decay in terms of a larger quantity that results in a better convergence rate. To the best of our knowledge, this is the first result that shows $\mathcal{O}(\frac{1}{k})$ convergence of the Kaczmarz type methods for solving linear feasibility problems [7]. An interesting corollary of our method is the Cesaro average result for the SKM method. Furthermore, the result holds under weaker assumptions than the previous Theorems.

---

[7] Several works exits for the Kaczmarz type methods for solving linear systems [38, 57].

**Theorem 3.16.** *Let $\{x_k\}$ be the random sequence generated by Algorithm 2. Take, $-1 < \xi \leq 0$ and $0 < \delta <$ $\frac{2(1+\xi)}{1-2\xi}$. Define $\tilde{x}_k = \frac{1}{k}\sum_{l=1}^{k} x_l$ and $f(x)$ as in (9), then*

$$\mathbb{E}\left[f(\bar{x}_k)\right] \leq \frac{(1+\xi)(1+\xi-2\delta\xi\mu_2)\,d(x_0,P)^2 + 2\xi\delta(\delta\xi-\delta-1)f(x_0)}{2\delta k\,(2+2\xi+2\delta\xi-\delta)}.$$

*Proof.* See Appendix. □

**Corollary 3.16.1.** *Let $\{x_k\}$ be the random sequence generated by SKM method (algorithm 1). Define $\tilde{x}_k = \frac{1}{k}\sum_{l=1}^{k} x_l$ and $f(x)$ as in (9), then*

$$\mathbb{E}\left[f(\bar{x}_k)\right] \leq \frac{d(x_0,P)^2}{2\delta k\,(2-\delta)},$$

*holds for any $0 < \delta < 2$.*

*Proof.* Take $\xi = 0$ in Theorem 3.16, then the result follows. □

The next Theorem is an extension of the result obtained in [11] and to a certain extent, it can be taken as an extension of Telgen's result [28]. The Theorem gives one a certificate of feasibility after a finite number of GSKM iterations. Before delving into the Theorem, we will provide some known Lemmas for the SKM algorithm which holds for the GSKM algorithm too. We refer interested readers to the work of De-Loera *et. al* [11] for detailed proof of these Lemmas (Lemma 3.17 to Lemma 3.19).

**Lemma 3.17.** *(Lemma 1 in [11]) Define, $\theta(x) = \left[\max_i\{a_i^T x - b_i\}\right]^+$ as the maximum violation of point $x \in \mathbb{R}^n$ and the length of the binary encoding of a linear feasibility problem with rational data-points as*

$$\sigma = \sum_i \sum_j \ln\left(|a_{ij}| + 1\right) + \sum_i \ln\left(|b_i| + 1\right) + \ln\left(mn\right) + 2.$$

*Then if the rational system $Ax \leq b$ is infeasible, for any $x \in \mathbb{R}^n$, the maximum violation $\theta(x)$ satisfies the following lower bound:*

$$\theta(x) \geq \frac{2}{2^\sigma}.$$

**Lemma 3.18.** *(Lemma 4 in [11]) If $P$ is $n$-dimensional (full-dimensional) then the sequence of iterates $\{x_k\}$ generated by the GSKM method converges to a point $x \in P$.*

*Proof.* Since, by assumption, $P$ is full dimensional, then the rest of the proof follows the same argument as Lemma 4 in [11]. □

**Lemma 3.19.** *( [65]) If the rational system $Ax \leq b$ is feasible, then there is a feasible solution $x^*$ whose coordinates satisfy $|x_j^*| \leq \frac{2^\sigma}{2n}$ for $j = 1, ..., n$.*

**Certificate of feasibility:** To detect feasibility of the rational system $Ax \leq b$, one needs to find a point $x_k$ such that $\theta(x_k) < 2^{1-\sigma}$. Such a point if exists will be called a certificate of feasibility. When the system is feasible, one expects to find a certificate of feasibility after finitely many iterations, and that if one fails to find a certificate after finitely many iterations, one can obtain a lower bound on the probability that the system is infeasible. Moreover, as discussed in the next Theorem, if the system is feasible, one can bound the probability of finding a certificate of feasibility.

**Theorem 3.20.** *Suppose $A, b$ are rational matrices with binary encoding length, $\sigma$, and that we run the GSKM method ($0 < \delta < 2$, $\xi \in Q$) on the system $Ax \leq b$ ($\|a_i\| = 1, i = 1, 2, ..., m$) with $x_0 = 0$. Suppose the number*

*of iterations $k$ satisfies the following lower bound:*

$$\frac{4\sigma - 4 - \log n + \log(1 + \phi)}{\log\left(\frac{1}{\rho}\right)} < k.$$

*If the system $Ax \leq b$ is feasible, then,*

$$p \leq H(\sigma, \phi, k, \bar{\rho}) = \sqrt{\frac{1 + \phi}{n}} \, 2^{2\sigma - 2} \, \bar{\rho}^{\frac{k}{2}},$$

*where $p$ is the probability that the current iterate is not a certificate of feasibility. And $\bar{\rho} = \max\{\rho, \rho_2^2\} < 1$, where $\rho$ and $\rho_2$ are defined in Theorem 3.14 and Theorem 3.15 for the choice $\xi \in Q_1$ and $\xi \in Q_2$, respectively. Also note that the function $H(\sigma, \phi, k, \bar{\rho})$ is a decreasing function with respect to $k$.*

*Proof.* See Appendix.

$\square$

**Remark 3.21.** *Note that instead of a normalized system if we consider a non-normalized system $\overline{A}x \leq \overline{b}$, $\|\overline{a_i}\| \neq 1$ for some $i$, then suppose the number of iterations $k$ satisfies the following lower bound:*

$$\frac{4\overline{\sigma} - 4 - \log n + \log(1 + \phi) + 2\log\psi}{\log\left(\frac{1}{\rho}\right)} < k,$$

*where $\overline{\sigma}$ is the binary encoding length for $\overline{A}, \overline{b}$. If the system $\overline{A}x \leq \overline{b}$ is feasible, then,*

$$p \leq \sqrt{\frac{1 + \phi}{n}} \, 2^{2\overline{\sigma} - 2} \, \psi \, \bar{\rho}^{\frac{k}{2}},$$

*where $p =$ probability that the current update $x_k$ is not a certificate of feasibility and $\psi = \max_j \|\overline{a_j}\|$.*

**Corollary 3.21.1.** *(Theorem 1.5 in [11]) Suppose $\overline{A}, \overline{b}$ are rational matrices with binary encoding length, $\overline{\sigma}$, and that we run the SKM method on the system $\overline{A}x \leq \overline{b}$ ($\|\overline{a_i}\| \neq 1$ for some $i$) and $x_0 = 0$. Suppose the number of iterations $k$ satisfies the following lower bound:*

$$\frac{4\overline{\sigma} - 4 - \log n + 2\log\psi}{\log\left(\frac{1}{h(\delta)}\right)} < k,$$

*where $\overline{\sigma}$ is the binary encoding length for $\overline{A}, \overline{b}$. If the system $\overline{A}x \leq \overline{b}$ is feasible, then,*

$$p \leq \sqrt{\frac{1}{n}} \, 2^{2\overline{\sigma} - 2} \, \psi \, [h(\delta)]^{\frac{k}{2}},$$

*where $p =$ the probability that the current update $x_k$ is not a certificate of feasibility and $\psi = \max_j \|\overline{a_j}\|$.*

*Proof.* Take $\xi = 0$ in Theorem 3. Then, we have, $\phi = 0$, $\rho = \phi + \phi_1 = h(\delta) = \rho_2^2$. It can be easily checked that the GSKM method with $\xi = 0$ is just the SKM method. Now, considering Theorem 3.20 with the above parameter choice, we can get the bound of Corollary 3.21.1. $\square$

### 3.3 Convergence Analysis of the PASKM Method

In this subsection, we study convergence properties of the proposed PASKM algorithm, i.e., we study the convergence behavior of the quantities of $\mathbb{E}[\|v_k - \mathcal{P}(v_k)\|^2]$, $\mathbb{E}[\|x_k - \mathcal{P}(x_k)\|^2]$, $\mathbb{E}[\|y_k - \mathcal{P}(y_k)\|^2]$ and $\mathbb{E}[f(x_k)]$ generated by the PASKM method. We proved that for a range of step parameters $\alpha, \gamma, \omega$, the proposed PASKM method enjoys a global linear rate. We also provided convergence analysis of the function values $f(x_k)$, with respect to the Cesaro average. The next Theorem deals with the convergence of the sequences $\{v_k\}$ and $\{y_k\}$ as well as the function values $f(x_k)$ generated by the PASKM algorithm.

**Theorem 3.22.** *Let $\{x_k\}$ be the sequence of random iterates generated by algorithm 3 and let $0 < \delta < 2$ and $0 \le \alpha, \omega \le 1$ such that $\gamma + 3\omega - 2 \le 0$, $\omega h(\delta)(1 - \alpha)(1 + \gamma) < 1$ and the following condition*

$$\omega(1 + \gamma) + h(\delta)(1 - \alpha) + \alpha(1 - \omega) + \alpha\gamma\mu_1(\gamma + 3\omega - 2)$$
$$- \omega h(\delta)(1 - \alpha)(1 + \gamma) < 1, \tag{28}$$

*holds. Define, $\Pi_1 = \omega(1 + \gamma)$, $\Pi_2 = (1 - \omega) + \gamma\mu_1(\gamma + 3\omega - 2)$, $\Pi_3 = \alpha\omega(1 + \gamma)$, $\Pi_4 = (1 - \alpha)h(\delta) + \alpha(1 - \omega) + \alpha\gamma\mu_1(\gamma + 3\omega - 2)$ and $\Gamma_1, \Gamma_2, \Gamma_3, \rho_1, \rho_2$ as in (24). Then the sequence of iterates $\{v_k\}$ and $\{y_k\}$ converges and the following results hold:*

$$\mathbb{E}\begin{bmatrix} d(v_{k+1}, P)^2 \\ d(y_{k+1}, P)^2 \end{bmatrix} \le \begin{bmatrix} \Gamma_2\Gamma_3(\Gamma_1 - 1)\,\rho_1^{k+1} + \Gamma_1\Gamma_3(\Gamma_2 + 1)\,\rho_2^{k+1} \\ \Gamma_3(\Gamma_1 - 1)\,\rho_1^{k+1} + \Gamma_3(\Gamma_2 + 1)\,\rho_2^{k+1} \end{bmatrix} d(y_0, P)^2,$$

*and*

$$\mathbb{E}\left(f(y_{k+1})\right) \le \frac{\mu_2}{2}\left[\Gamma_3(\Gamma_1 - 1)\,\rho_1^{k+1} + \Gamma_3(\Gamma_2 + 1)\,\rho_2^{k+1}\right]\,d(y_0, P)^2.$$

*where $\Gamma_1, \Gamma_3 \ge 0$ and $0 \le \rho_1 \le \rho_2 < 1$.*

*Proof.* See Appendix 3. □

The next Theorem deals with the convergence of the sequences $\{v_k\}$ and $\{x_k\}$ generated by the PASKM algorithm.

**Theorem 3.23.** *Let, $v_{k+1}$ and $x_{k+1}$ are generated by Algorithm 3. If we select the parameters $\omega$, $\gamma$, $\alpha$ as*

$$\omega = 1 - \frac{\zeta\mu_1^2 + 2\gamma\mu_1 - \zeta\mu_1}{1 + \zeta\mu_1^2}, \quad \gamma = \sqrt{\zeta\eta\mu_1}, \quad \alpha = \frac{\eta}{\eta + \gamma},$$

*where, $\zeta$ is chosen as $0 < \zeta < \frac{4\eta\mu_1}{(1-\mu_1)^2}$ if $\mu_1 < 1$, otherwise choose any $\zeta > 0$. Then, for any $0 < \delta < 2$, the sequence of iterates $\{v_k\}$, $\{x_k\}$ converges and the following result holds:*

$$\mathbb{E}\left[d(v_{k+1}, P)^2 + \zeta\mu_1\,d(x_{k+1}, P)^2\right] \le \omega^{k+1}\,\mathbb{E}\left[d(v_0, P)^2 + \zeta\mu_1\,d(x_0, P)^2\right]$$
$$= (1 + \zeta\mu_1)\,\omega^{k+1}\,d(x_0, P)^2.$$

*This theorem implies that the PASKM algorithm converges linearly with a rate of $\omega$, which accumulates to a total of $\mathcal{O}(\frac{1 + \zeta\mu_1^2}{\zeta\mu_1^2 + 2\gamma\mu_1 - \zeta\mu_1}\log 1/\epsilon)$ iterations to bring the given error below $\epsilon > 0$.*

*Proof.* See Appendix 3. □

In the next Theorem, we present the convergence analysis of the function $f(x)$ with respect to the Cesaro average for the PASKM algorithm. We showed that the Cesaro average of the PASKM iterates converges to the optimum at a rate of $\mathcal{O}(1/k)$ where $k$ is the number of iterations.

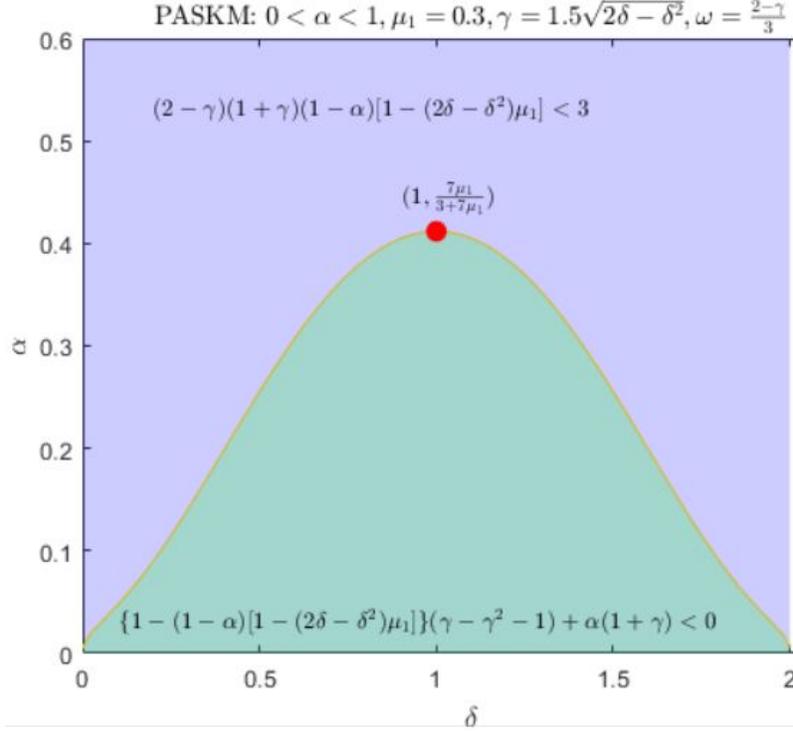**Theorem 3.24.** *Let $\{y_k\}$ be the random sequence generated by Algorithm 3. Take, $0 \le 1 - \alpha, \omega < 1$, $0 < \delta < \frac{2(1-\omega+\alpha\omega)}{1+2\omega-2\alpha\omega}$ and $\alpha\gamma = \alpha\delta + \omega\delta(1 - \alpha)$. Define $\tilde{y}_k = \frac{1}{k}\sum\limits_{l=1}^{k} y_l$ and $f(y)$ as in (9), then*

$$\mathbb{E}\left[f(\bar{y}_k)\right] \le \frac{(1 - \omega + \alpha\omega)^2\,d(y_0, P)^2 + 2\delta(\delta - 2 + 3\omega - 3\alpha\omega + \delta\omega - \delta\alpha\omega)f(y_0)}{2\delta k\,(2 - 2\omega + 2\alpha\omega - 2\delta\omega + 2\delta\alpha\omega - \delta)}.$$

*Proof.* See Appendix 3. □

**Parameter selection for PASKM algorithm** In this section, we discuss allowable parameter selection for the PASKM algorithm based on Theorem 3.22. If the parameters $0 \le \alpha, \omega \le 1$ and $\gamma \ge 0$ satisfies $\gamma + 3\omega - 2$ and the

condition of (28) hold then the PASKM method will converge for any $0 < \delta < 2$ [8]. To simplify the conditions for ease of implementation, let's take $0 \leq \gamma < 2$ and $\omega = \frac{2-\gamma}{3+p}$ for some $0 \leq p \leq \frac{1}{\mu_1}$ [9].



**Figure 4:** Allowable parameter range

In Figure 4, we plot the feasible region considering the above parameter choice and the conditions of Theorem 3.22. Considering the choice of $\gamma$ and $\omega$, the condition $\omega h(\delta)(1 - \alpha)(1 + \gamma) < 1$ simplifies to

$$\alpha > 1 - \frac{3+p}{(2-\gamma)(1+\gamma)h(\delta)} = \frac{2h(\delta) + \gamma h(\delta) - \gamma^2 h(\delta) - 3 - p}{(2-\gamma)(1+\gamma)h(\delta)}$$

$$= \frac{(\gamma - \gamma^2 - 1 - p) + \eta\mu_1(\gamma^2 - 2 - \gamma)}{(2-\gamma)(1+\gamma)h(\delta)} \leq 0,$$

where, we used the fact that the conditions $\gamma - \gamma^2 - 1 \leq 0$ and $\gamma^2 - 2 - \gamma \leq 0$ hold for any $0 \leq \gamma \leq 2$. That implies for any $\alpha \geq 0$, the condition $\omega h(\delta)(1 - \alpha)(1 + \gamma) < 1$ holds. Similarly, we can simplify the condition of (28) as follows:

$$\alpha < \underbrace{\frac{(1 + p - \gamma + \gamma^2)(1 - h(\delta))}{1 - h(\delta) + p + \gamma + (\gamma - p)h(\delta) - \gamma^2 h(\delta) + \mu_1 p\gamma(\gamma - 2)}}_{> 0 \text{ for } 0 < \delta < 2}$$

$$= \alpha(\gamma, \delta, p) \leq 1. \tag{29}$$

Therefore, if we choose $\gamma, \omega$ and $\alpha$ as

$$\gamma = 1.5\sqrt{2\delta - \delta^2}, \ p = 0, \ \omega = \frac{2-\gamma}{3+p}, \ \alpha = 0.99 * \alpha(\gamma, \delta, p) \tag{30}$$

$$\gamma = 2\sqrt{2\delta - \delta^2}, \ p = 0, \ \omega = \frac{2-\gamma}{3+p}, \ \alpha = 0.99 * \alpha(\gamma, \delta, p), \tag{31}$$

---

[8]When, $\delta = 2$, we have $h(\delta) = 1 - \eta\mu_1 = 1$. In that case, we can simplify the condition of (28) as $\omega < \frac{2-\gamma}{3+\frac{1}{\mu_1}}$. In other words, for $\delta = 2$ the PASKM algorithm will converge if we select the parameters as $0 \leq \gamma < 2$, $0 \leq \alpha \leq 1$, $\omega < \frac{2-\gamma}{3+\frac{1}{\mu_1}}$ and $\mu_1 = \frac{\lambda^+_{\min}(A^T A)}{m}$.

[9]Note that for the choice $p > \frac{1}{\mu_1}$ the condition (29) trivially holds as the right hand side of (29) is always greater than 1.

then the convergence result of Theorem 3.22 holds for the PASKM algorithm. We will use these two sets of parameter choices in our numerical experiments. Note that, our choice is empirical in nature. One can probably find a better combination of parameters than (30) and (31). Similarly, if we choose $\gamma, \omega$ and $\alpha$ as

$$\zeta = \frac{3.99\eta\mu_1}{(1-\mu_1)^2}, \ \omega = 1 - \frac{\zeta\mu_1^2 + 2\gamma\mu_1 - \zeta\mu_1}{1 + \zeta\mu_1^2}, \ \ \gamma = \sqrt{\zeta\eta\mu_1}, \ \ \alpha = \frac{\eta}{\eta + \gamma}, \tag{32}$$

then the convergence result of Theorem 3.23 holds for the PASKM algorithm. The choice of (32) is not of practical benefit as the value of $\frac{\lambda_{\min}^+(A^T A)}{m}$ is very small for most test cases. From (32), we have $\gamma \propto \frac{1}{m}$, which is very small for large test instances. Smaller $\gamma$ slows down the convergence of the PASKM algorithm as $\gamma$ can be seen as a projection parameter like $\delta$.

## 4 Numerical Experiments

In this section, we discuss the numerical experiments performed to show the computational efficiency of the proposed algorithms (Algorithm 2 and 3). As mentioned before, we limit our focus on the over-determined systems regime (i.e., $m \gg n$) where iterative methods are competitive in general. However, from our experiments, we see similar computational behavior for the under-determined systems as well.

### 4.1 Experiment Specifications

We implemented the proposed GSKM and PASKM algorithms in *MATLAB R2018b* and performed the experiments in a Dell Precision 7510 workstation with 32GB RAM, Intel Core i7-6820HQ CPU, processor running at 2.70 GHz. To analyze computational performance, we perform the numerical experiments for a wide range of instances including both randomly generated and real-world test problems.

- **Randomly generated problems:** Gaussian and highly correlated systems
- **Real-world test instances:** Standard ML data sets and Sparse Netlib LP instances

We compare SKM with two versions of the proposed GSKM and PASKM algorithms for a better understanding of the algorithmic behavior. In Table 2, we provide the parameter choices for GSKM and PASKM algorithms. Throughout the numerical experiments section, we compared SKM with GSKM-1, GSKM-2 and PASKM-1, PASKM-2.

**Table 2:** Parameter choice of GSKM and PASKM algorithms for the numerical experiments.
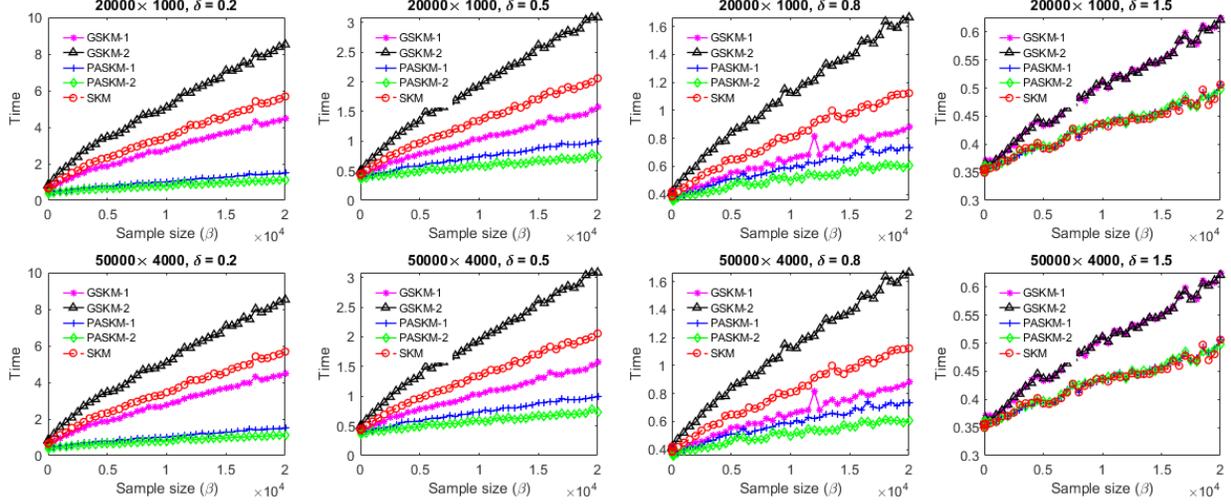
| Parameters | GSKM (Algorithm 2, $\xi \in Q$) | | | PASKM (Algorithm 3, $\alpha, \omega, \gamma$) | |
|---|---|---|---|---|---|
| | SKM | GSKM-1 | GSKM-2 | PASKM-1 | PASKM-2 |
| $1 \le \beta \le m$ $0 < \delta < 2$ | $\xi = 0$ | $\xi = -0.1$ $\xi = -0.2$ | $\xi = 0.5$ | $\alpha, \omega, \gamma$ as in (30) | $\alpha, \omega, \gamma$ as in (31) |

Finally, we investigate the performance behavior of the proposed GSKM and PASKM methods with state-of-the-art methods such as Interior point methods (IPMs) and Active set methods (ASMs) for several Netlib LP instances. The total CPU time is calculated in seconds (s). For a fair comparison, we run the algorithms 10 times and report the averaged performance throughout the experiments. Moreover, all the algorithms start from the same initial point that is far away from the feasible region.
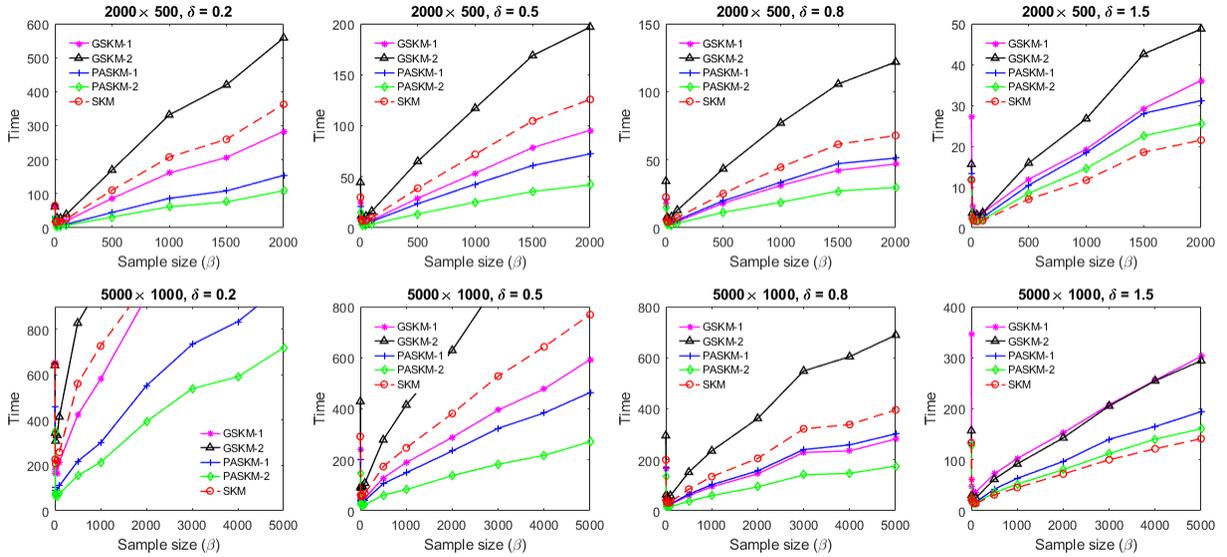
### 4.2 Experiments on Randomly Generated Instances

We considered the linear feasibility $Ax \le b$, where the entries of matrices $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$ are chosen randomly from a certain distribution. To maintain the system consistency (i.e., $b \in \mathcal{R}(\mathbf{A})$), we first generated vectors $x_1, x_2 \in \mathbb{R}^n$ at random from the corresponding distributions, then set $b$ as the convex combination of vectors $Ax_1$ and $Ax_2$ (i.e., $b = \sigma Ax_1 + (1-\sigma)Ax_2, \ 0 \le \sigma \le 1$). Two types of random data sets are considered: highly correlated, and Gaussian. For the correlated systems, data matrices $A$ and $x_1, x_2$ are chosen uniformly at random between $[0.9, 1.0]$ (i.e., $a_{ij}, x_j \in [0.9, 1.0], \ i = 1, 2, ..., m, \ j = 1, 2, ..., n$). For the Gaussian system data matrices, $A$ and $x_1, x_2$ are chosen uniformly at random from standard normal distribution (i.e., $a_{ij}, x_j \in \mathcal{N}(0, 1), \forall i, j$). Moreover, the vector $b \in \mathbb{R}^m$ is generated by following the above-mentioned procedure.

**CPU time vs Sample size** $\beta$    We first compared the total CPU time of the proposed algorithms (GSKM-1, GSKM-2, PASKM-2, PASKM-2) with the original SKM algorithm. The comparison is carried out by varying the sample size $\beta$ from 1 to the total row size $m$. The positive residual error tolerance is chosen as $10^{-05}$ (i.e., $\| (Ax - b)^+ \|_2 \leq 10^{-05}$). The comparison is carried out for $\delta = 0.2, 0.5, 0.8$ and 1.5. In Figure 5, we compared the above-mentioned algorithms for two randomly generated highly correlated linear feasibility problems of size $20000 \times 1000$ and $50000 \times 4000$. From Figure 5, we see that the proposed GSKM-1, PASKM-1, PASKM-2 algorithms outperform the SKM algorithm in terms of average CPU time when $\delta = 0.2, 0.5, 0.8$. For $\delta = 1.5$, the performance of SKM, PASKM-1 and PASKM-2 are fairly similar whereas, the performance of GSKM-1 and GSKM-2 are worse compared to all other algorithms.
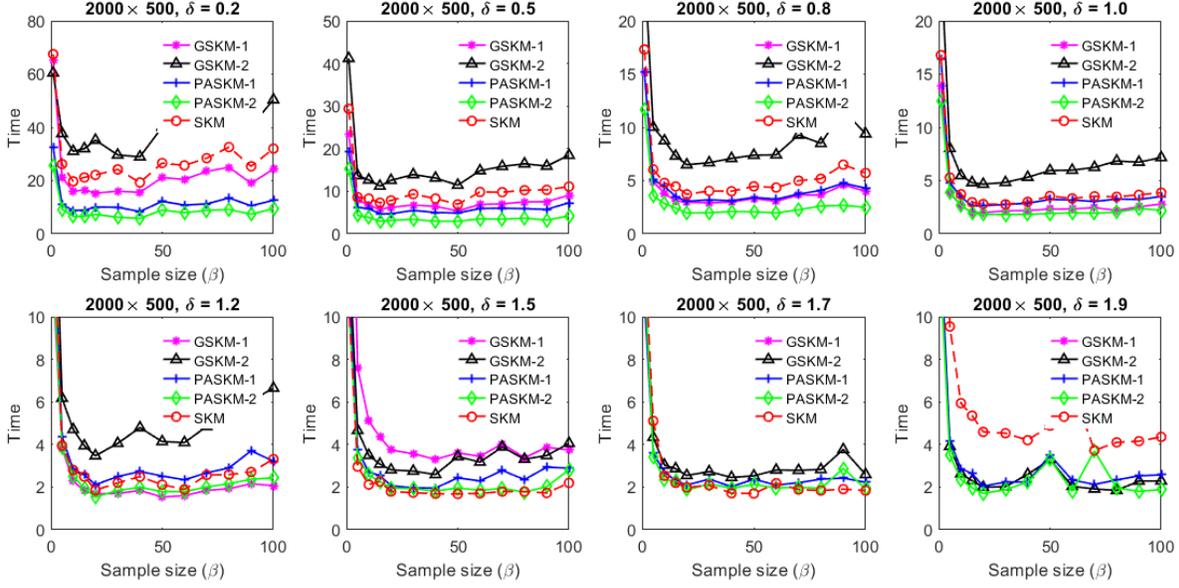


**Figure 5:** Sample size $\beta$ VS average CPU time comparison among SKM, GSKM, PASKM variants for $\delta = 0.2, 0.5, 0.8, 1.5$ on correlated systems. Problem size: $20000 \times 1000$ (Top panel), $50000 \times 4000$ (Bottom panel).



**Figure 6:** Sample size $\beta$ VS average CPU time comparison among SKM, GSKM, PASKM variants for $\delta = 0.2, 0.5, 0.8, 1.5$ on Gaussian systems. Problem size: $2000 \times 500$ (Top panel), $5000 \times 1000$ (Bottom panel).

We present the time versus sample size plot for two randomly generated Gaussian system of size $2000 \times 500$ and $5000 \times 1000$ in Figure 6. All the algorithms show similar performance patterns as shown in the correlated systems (Figure 5) for the choice of $0 < \delta < 1$. However, for the case of $\delta = 1.5$, SKM and PASKM-2 perform marginally better than the other algorithms. Since all of the considered methods perform significantly well whenever $\beta$ is small (i.e., $1 < \beta \leq 100$). For a better understanding, we compare the proposed algorithms for $1 < \beta \leq 100$. In Figure 7, we plot the time vs $\beta$ graph for a $2000 \times 500$ Gaussian problem for smaller $\beta$.

**Figure 7:** Sample size $\beta$ VS average CPU time comparison among SKM, GSKM, PASKM variants for $\delta = 0.2, 0.5, 0.8, 1, 1.2, 1.5, 1.7, 1.9$ and samller sample size (i.e., $1 \leq \beta \leq 100$) on a $2000 \times 500$ Gaussian system.
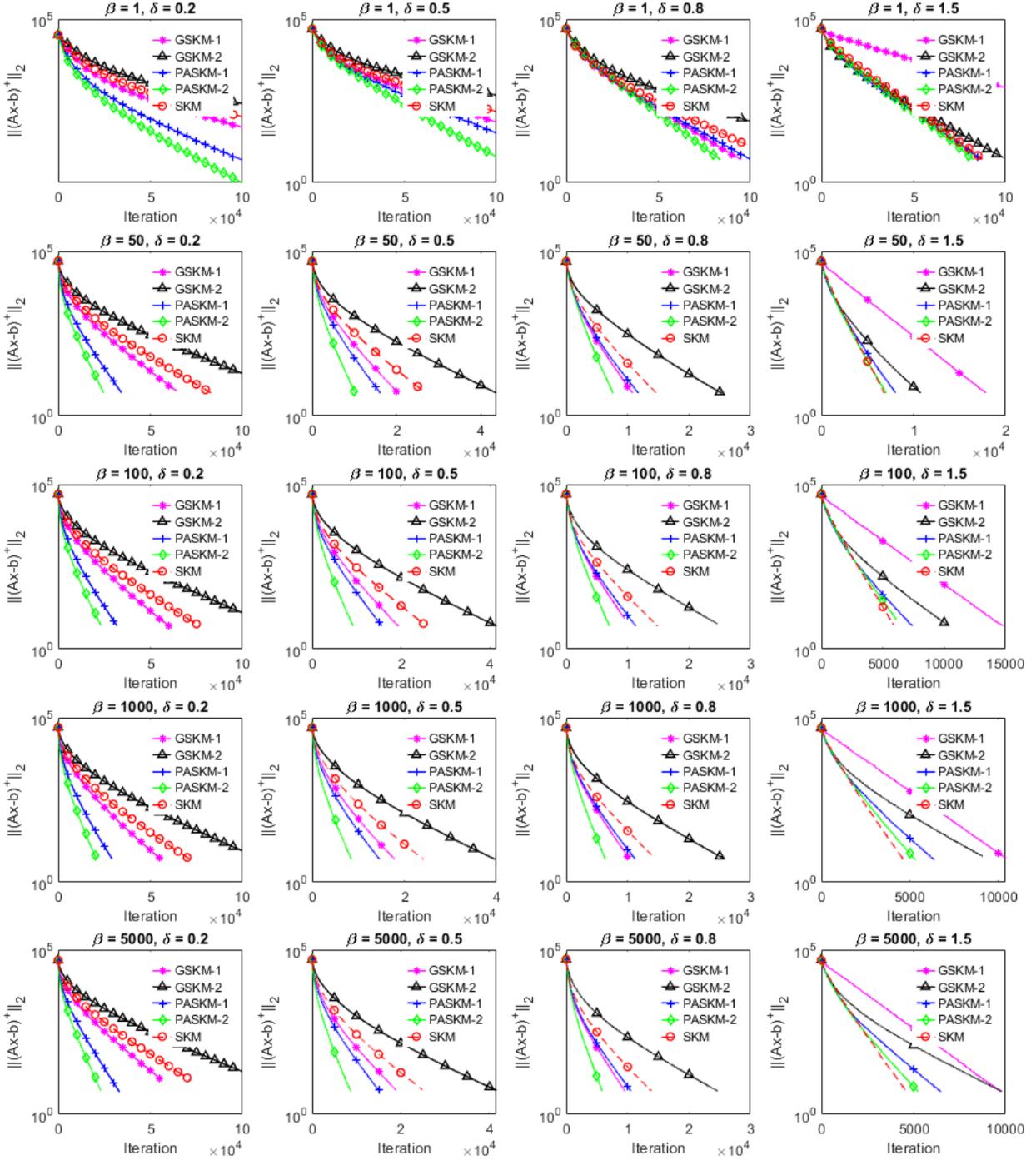
In a nutshell, we can conclude that for the choice of $0 < \delta < 1$, PASKM-1, PASKM-2 and GSKM-1 outperform the original SKM method. And in that region, PASKM-2 is the best performing algorithm. Moreover, for $1.5 \leq \delta \leq 1.7$, all of the proposed algorithms perform similarly as the SKM method. However, for the case of $\delta = 1.9$, the proposed algorithms significantly outperform the SKM method. Furthermore, we believe with correct parameter choice one can find better-performing variants of GSKM and PASKM compared to the SKM algorithm for the case of $1.5 \leq \delta \leq 1.7$. Finally from Figure 7, we can deduce that the best sample size choice for all of the considered methods occurs at $1 < \beta \ll m$. This amplifies the importance of the special sampling distribution selection.

**Positive residual error** $\| (Ax - b)^+ \|_2$ **VS No. of iterations and Time**   Now, we compare the respective convergence trend for the considered algorithms with respect to the number of iterations and CPU time. We choose positive residual error $\| (Ax - b)^+ \|_2$ as the convergence measure and considered $5000 \times 1000$ Gaussian system. We carried out the analysis for several choices of sample sizes, $\beta = 1, 100, 1000, m$ and the choice of $\delta$ values remains the same as before. In Figures 8 and 9, we provide the respective positive residual decay results for different sample sizes and different projection parameters. We plot positive residual error VS iteration and positive residual error VS time in Figures 8 and 9, respectively. From Figures 8 and 9, we see that irrespective of sample size, $\| (Ax_k - b)^+ \|_2$ converges to zero much faster for the proposed PASKM-1 and PASKM-2, GSKM-1 compared to SKM whenever $\delta < 1$. For the case of $\delta = 1.5$, SKM and PASKM-2 has a similar kind of performance whereas the GSKM-1 performs poorly compared to SKM and PASKM method. As expected, the choice $\beta = 1$ produces the slowest rate and the choice $\beta = 100$ produces the best convergence graph.

**Fraction of satisfied constraints (FSC) VS No. of iterations and Time**   To investigate the generated solution quality of the above-mentioned algorithms of Table 2, we measure the number of satisfied constraints at each iteration, for that we define,

$$\text{Fraction of Satisfied Constraints (FSC)} = \frac{\text{Number of satisfied constraints}}{\text{Total number of constraints } (m)}$$
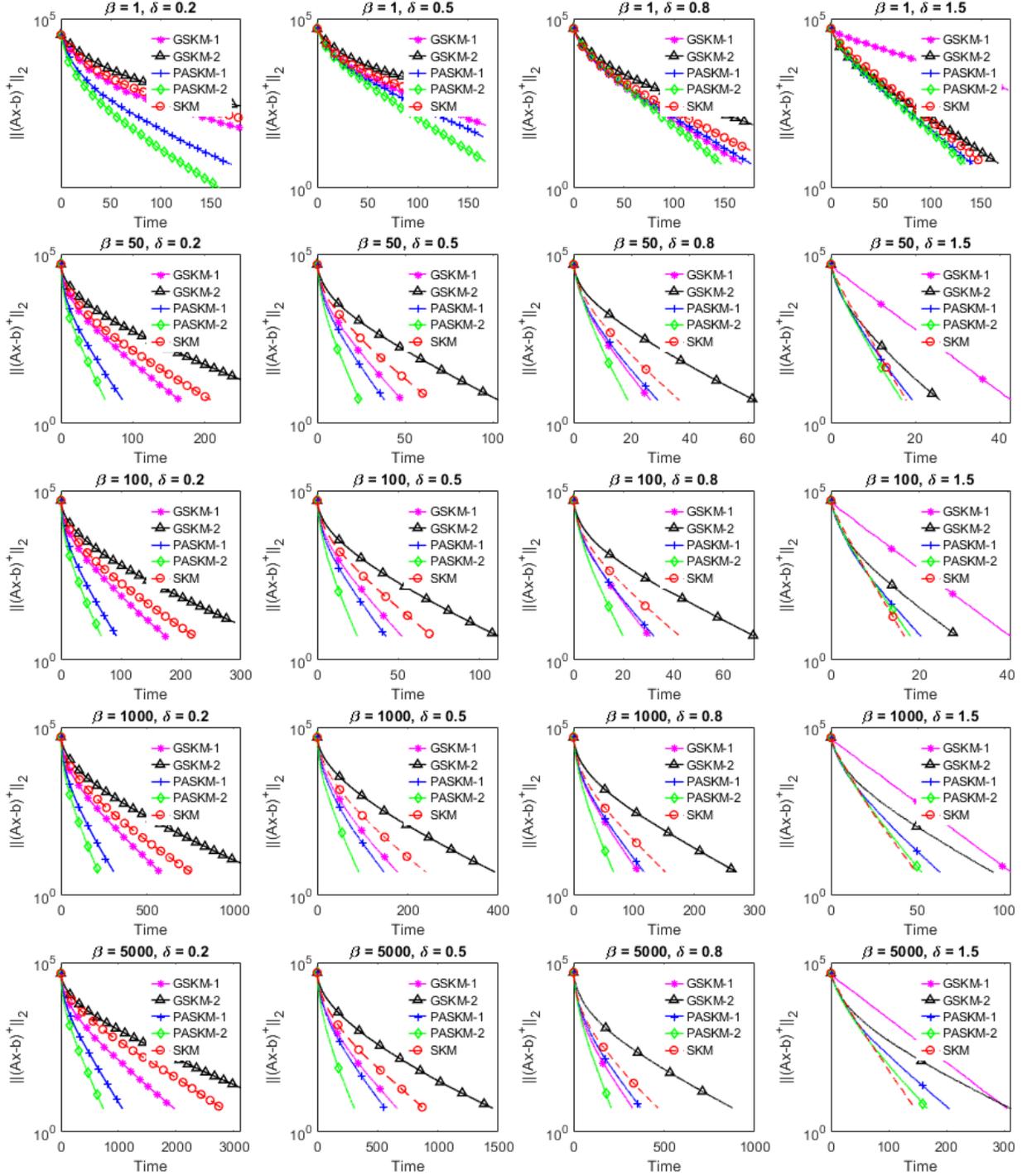
Note that, at any particular iteration we have, $0 \leq \text{FSC} \leq 1$. In Figures 10 and 11, we plot the value of FSC with respect to No. of iterations and CPU time of each algorithm respectively. From Figures 10 and 11, we can see that the choice of $\beta = 1$ is the worst choice for all algorithms as the improvement of FSC is much slower compared to other choices of $\beta$. And for the choice $\beta = 100$, we get the best solution quality for each algorithm. Our proposed GSKM-1, PASKM-1 and PASKM-2 algorithms outperform the other methods significantly for $0 < \delta < 1$ but, for $\delta = 1.5$ only PASKM-2 performs similar to SKM.

**Figure 8:** Positive residual error $\| (Ax - b)^+ \|_2$ VS No. of iteration comparison among SKM, GSKM, PASKM variants for $\delta = 0.2, 0.5, 0.8, 1.5$ and $\beta = 1, 50, 100, 1000, 5000$ on $5000 \times 1000$ Gaussian system.

## 4.3 Experiments on real-world Instances

In this subsection, we consider some nonrandom, real-world test instances. For the sake of unbiased performance analysis, we consider the following two types of real-world data-sets: standard Machine Learning (ML) data-sets for Support Vector Machine (SVM) classifier [11, 66, 67], and sparse linear feasibility problems extracted from benchmark Netlib LP problems [68].

**Figure 9:** Positive residual error $\|(Ax - b)^+\|_2$ VS CPU time comparison among SKM, GSKM, PASKM variants for $\delta = 0.2, 0.5, 0.8, 1.5$ and $\beta = 1, 50, 100, 1000, 5000$ on $5000 \times 1000$ Gaussian system.
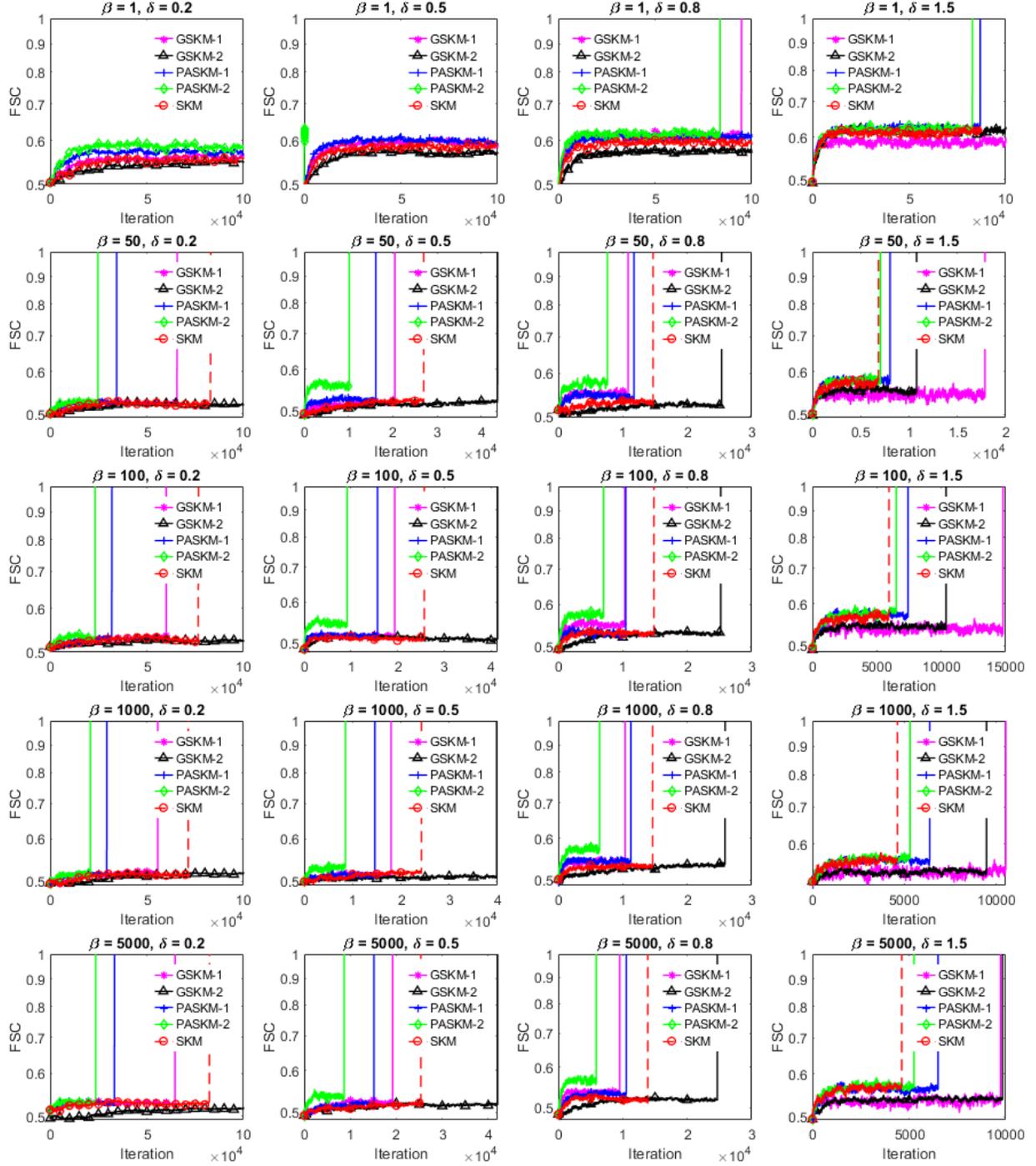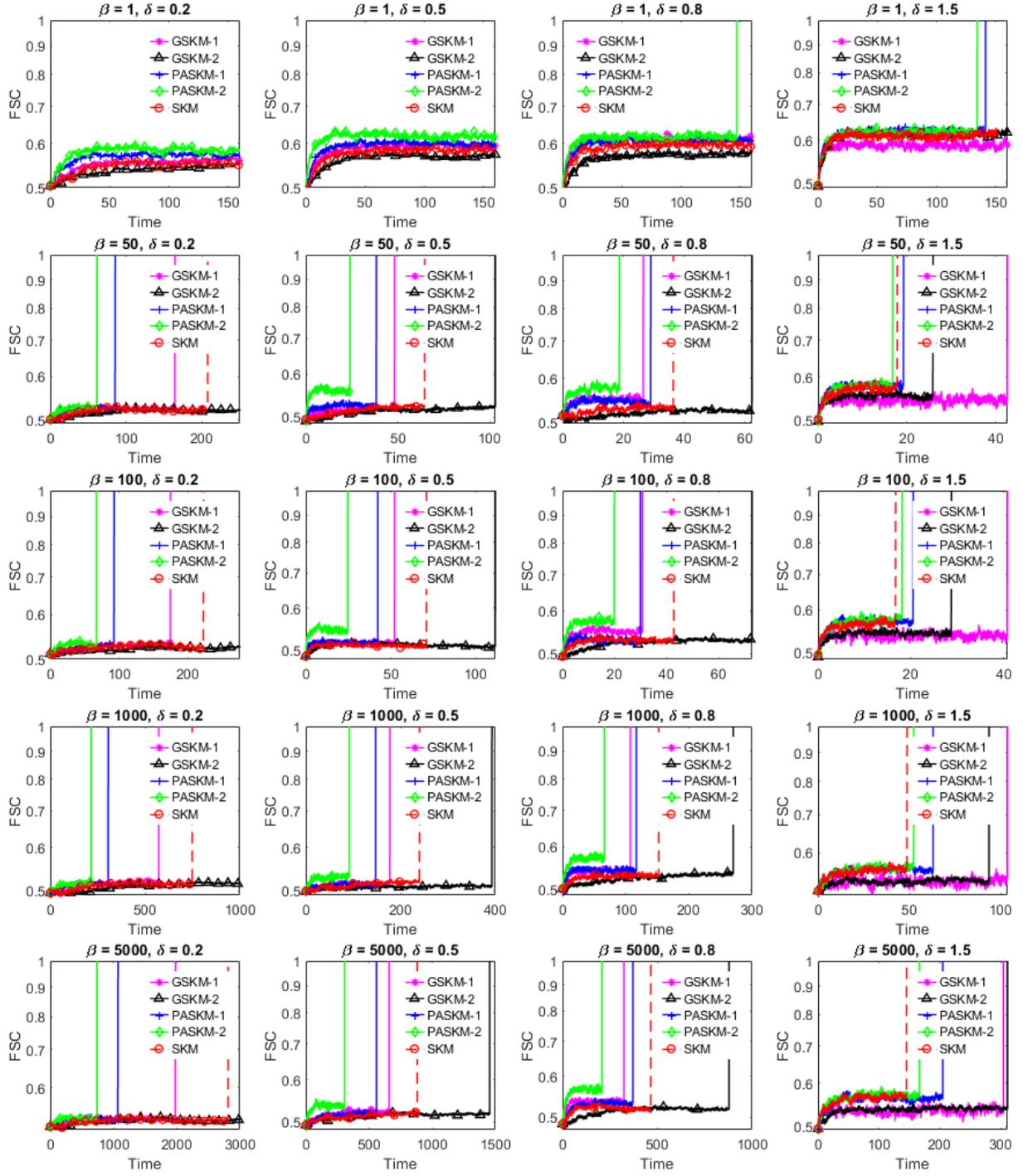
**SVM classifier instances** We first consider two linear feasibility problems arising from binary classification with SVM. We compare the proposed algorithms with SKM to the linear classification problem using the SVM model for the following two data sets: 1) Wisconsin (diagnostic) breast cancer data set and 2) Credit card default data set. The Wisconsin breast cancer data set consists of data points whose features are calculated from images. There are two types of data points: 1) malignant and 2) benign cancer cells. As shown by the researchers [11, 69], the SVM classifier

**Figure 10:** No. of iteration vs fraction of satisfied constraints (FSC) comparison among SKM, GSKM, PASKM variants for $\delta = 0.2, 0.5, 0.8, 1.5$ and $\beta = 1, 50, 100, 1000, 5000$ on $5000 \times 1000$ Gaussian system.

problem can be re-written as an equivalent homogeneous system of linear inequalities ($Ax \leq 0$), which represents the separating hyper-plane between malignant and benign data points. The constraint matrix $A$ has 569 rows (data points) and 30 columns (features). Since the data set is not perfectly separable, we allow tolerance for the positive residual $\|(Ax)^+\|$. For our experiments, we fixed the tolerance as $10^{-3}$ (i.e., we ran the algorithm until $\|(Ax_k)^+\| \leq 10^{-3}$ is satisfied).
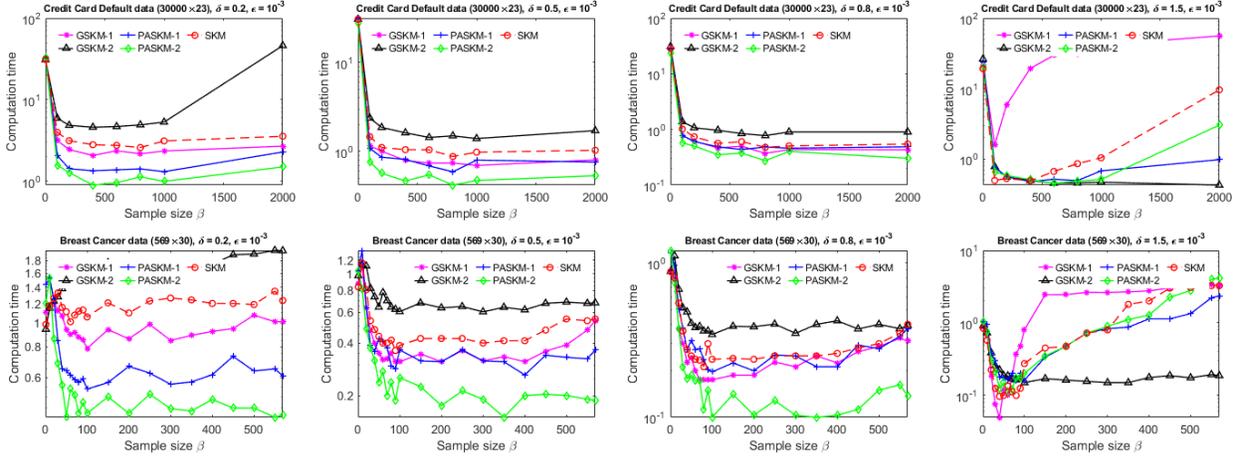
**Figure 11:** CPU time vs fraction of satisfied constraints (FSC) comparison among SKM, GSKM, PASKM variants for $\delta = 0.2, 0.5, 0.8, 1.5$ and $\beta = 1, 50, 100, 1000, 5000$ on $5000 \times 1000$ Gaussian system.

Similarly, we consider the credit card default data set described in [11,66]. This data set consists of features denoting the payment profile of a user and binary variables describing payment conditions in a certain billing cycle: 1 for payment made on time and 0 for late payment. The SVM classification problem for the data set can be transformed into an equivalent homogeneous system of inequalities ($Ax \leq 0$) like before. The solution $x^*$ denotes the coefficients of the separating hyper-plane between on-time and default data points. The transformed data matrix $A$ has 30000 rows (30000

24

user profiles) and 23 columns (22 profile features). As the data set is not separable, like the previous problem we allow a tolerance error. In this case, we ran the algorithms until the condition: $\|(Ax_k)^+\|/\|(Ax_0)^+\| \le 10^{-3}$ is satisfied.

**CPU time vs Sample size** $\beta$    We plot the CPU time VS sample size $\beta$ graphs for SVM problems in Figure 12. To be consistent with our previous experiments, we choose $\delta = 0.2, 0.5, 0.8, 1.5$. From Figure 12, we see that the proposed GSKM-1, PASKM-1 and PASKM-2 algorithms outperform the other algorithms including SKM for $\delta = 0.2, 0.5, 0.8$. However, for $\delta = 1.5$, GSKM-2 performs significantly well compared to the other methods. On the other hand, SKM, PASKM-1 and PASKM-2 follow a similar trend across different sample sizes. PASKM-1 and PASKM-2 marginally outperform SKM for this regime. Another interesting point can be noted that the comparison graphs for the credit card data set are not as smooth as the breast cancer data set graphs [10], which can be attributed to the irregularity of the constraint matrix $A$.



**Figure 12:** Average CPU time VS Sample size $\beta$ comparison among SKM, GSKM, PASKM variants for $\delta = 0.2, 0.5, 0.8, 1.5$ on Support Vector Machine problems; Top panel: Credit card data set, Bottom panel: Wisconsin breast cancer data set.
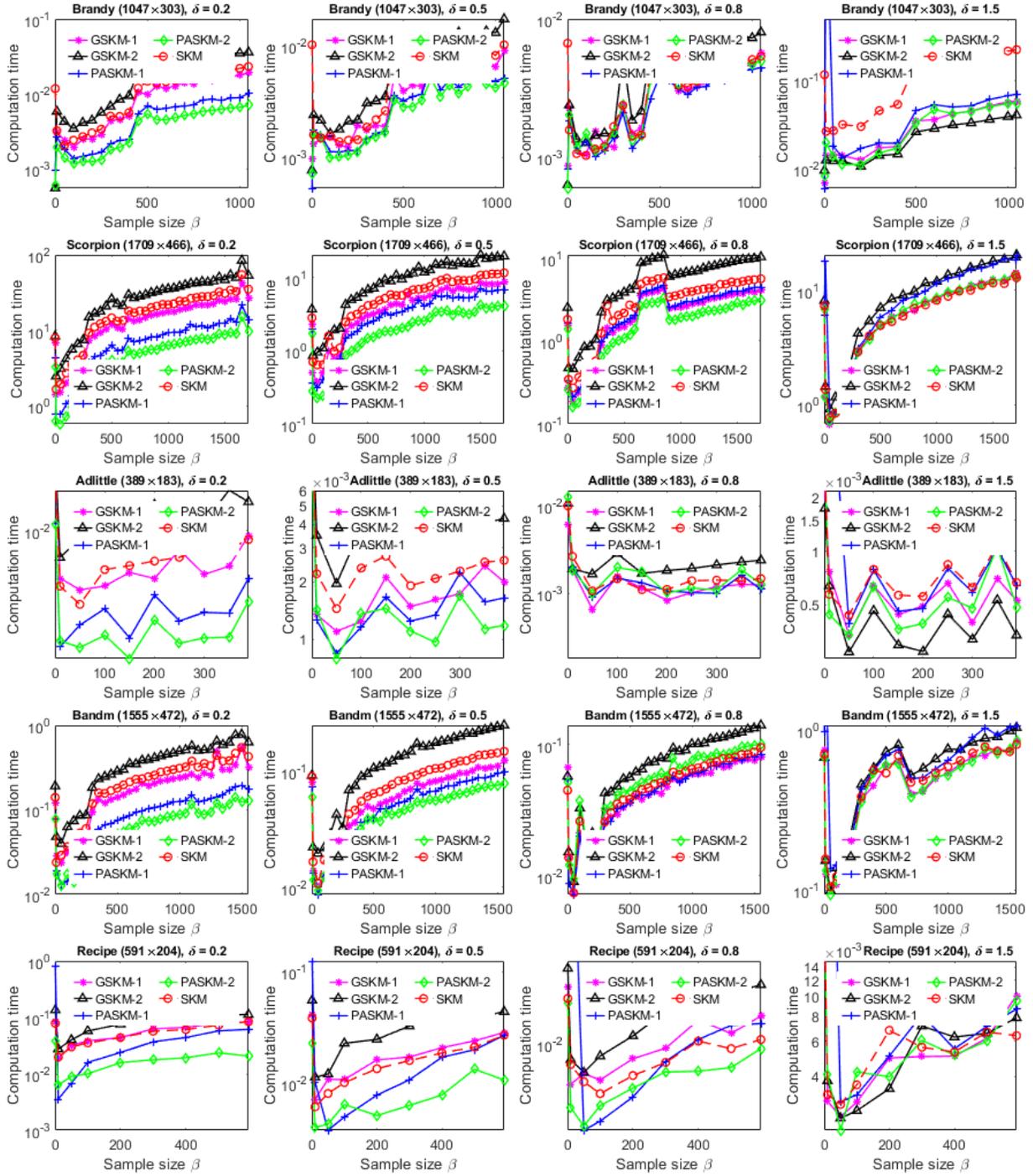
**Netlib LP instances**    We also investigate the comparative performance of the proposed algorithms with SKM on real-world sparse data sets. For this experiment, we consider some Netlib LP [68] test instances. Each of these problems is formulated as a standard linear programming problem ( $\min c^T x$ subject to $Ax = b, \; l \le x \le u$). To conduct the above-mentioned experiments, we transform each of these problems into an equivalent linear feasibility problem.

**CPU time vs Sample size** $\beta$    Now, we plot the CPU time VS sample size $\beta$ graphs for five Netlib LP instances in Figure 13. Later in subsection 4.4, we consider a total of ten Netlib LP instances including the five considered here. In Figure 13, we provide comparison graphs for the following Netlib LP test instances: lp-brandy, lp-addlittle, lp-scorpion, lp-bandm, lp-recipe. Furthermore, we consider different error tolerances for these problems (see Table 3 for details). From Figure 13, we see that the proposed GSKM-1, PASKM-1, PASKM-2 algorithms outperform the SKM algorithm for $\delta = 0.2, 0.5, 0.8$. In the case of $\delta = 1.5$, the performance of SKM, GSKM-1 and PASKM-2 are fairly similar for the problems lp-scorpion, lp-bandm and lp-recipe. For lp-brandy and lp-adlittle, all of the proposed variants of GSKM and PASKM outperform the original SKM.

## 4.4    Comparison with IPM and ASM for Netlib LP instances

In this subsection, we compare the performance of GSKM and PASKM variants with SKM and benchmark commercial solvers for solving Netlib LP test instances. We follow the standard framework used by De Loera *et. al* [11] and Morshed *et. al* [60] in their work for linear feasibility problems. The problem instances are transformed from standard LP problems (i.e., $\min c^T x$ subject to $Ax = b, \; l \le x \le u$ with optimum value $p^*$) to an equivalent linear feasibility formulation (i.e., $\mathbf{A} x \le \mathbf{b}$, where $\mathbf{A} = [A^T \; -A^T \; I \; -I \; c]^T$ and $\mathbf{b} = [b^T \; -b^T \; u^T \; -l^T \; p^*]^T$). For all of the experiments, we compared the proposed algorithms for $0 < \delta < 1$, since from our experiments in subsection 4.2 and 4.3, this is the domain where the proposed GSKM and PASKM variants significantly outperform the SKM method.

---

[10]The credit card data matrix has $30,000$ rows. From our earlier experiments, we observe that the choice of $1 < \beta \le 100$, the proposed algorithms produce the best performance. For that reason, we plot the credit card graph up-to $\beta = 2000$. The irregularity of the credit card graph occurs when $\beta > 2000$.

**Figure 13:** Average CPU time VS Sample size $\beta$ comparison among SKM, GSKM, PASKM variants for $\delta = 0.2, 0.5, 0.8, 1.5$ on Netlib LP instances.

In Table 3, we list the total CPU time in seconds for each of the above-mentioned algorithms in Table 2. In addition to that, we provide the CPU time for Interior point method (IPM) and Active set method (ASM) algorithms for solving the selected Netlib LP problems. For a better and fair comparison, the pseudo-code of the proposed methods and SKM is written in MATLAB and Optimization Toolbox function `fmincon` is used to implement IPM and ASM methods. We first solve the linear feasibility problem ($\mathbf{A}x \le \mathbf{b}$) with SKM, GSKM and PASKM variants and record the CPU time in Table 3. Note that, we can't use `fmincon`'s IPM and ASM algorithms directly to solve the linear feasibility problem

($\min 0$, $s.t$ $\mathbf{A}x \leq \mathbf{b}$) since both methods fail to solve the linear feasibility problems. The reason for that is, in IPM the *Karush Kuhn Tucker* (KKT) system at each iteration becomes singular, and ASM stops in the first step of finding a feasible solution. For a fair comparison, in Table 3, we list the total CPU consumption time as follows: for the SKM

**Table 3:** CPU time comparisons among the state-of-the-art methods (using MATLB's `fmincon` function) solving LP, and SKM, GSKM and PASK solving LF. * implies that the solver was unable to solve the problem with predetermined accuracy within 100,000 function evaluations. CPU time of the best performing algorithm for a problem is represented in bold letters.

| Instance | Dimensions | GSKM $\times 10^{-2}$ | PASKM $\times 10^{-2}$ | SKM $\times 10^{-2}$ | Interior Point | Active set | $\beta$ | $\epsilon$ $\times 10^{-2}$ |
|---|---|---|---|---|---|---|---|---|
| adlittle | $389 \times 138$ | 0.027 | **0.173** | 0.032 | 2.16 | 4.96 | 150 | 0.1 |
| agg | $2207 \times 615$ | 0.22 | **0.196** | 0.23 | 66.54* | 315.91* | 50 | 1 |
| bandm | $1555 \times 472$ | 9.82 | **4.057** | 9.2 | 14.57 | 529.43* | 50 | 1 |
| blend | $337 \times 114$ | 1.48 | **0.581** | 1.28 | 2.28 | 4.62 | 50 | 0.1 |
| brandy | $1047 \times 303$ | 0.53 | **0.491** | 14.06 | 16.97 | 63.11 | 1 | 1 |
| degen2 | $2403 \times 757$ | 26.26 | **10.139** | 20.73 | 7.13 | 21038 | 100 | 1 |
| finnis | $3123 \times 1064$ | 0.53 | 0.532 | **0.527** | 66.16* | 237750* | 10 | 0.1 |
| recipe | $591 \times 204$ | 0.60 | **0.164** | 0.52 | 0.89 | 63.24 | 50 | 0.1 |
| scorpion | $1709 \times 466$ | 156.9 | **42.712** | 125 | 17.68 | 8.02 | 50 | 1 |
| stocfor1 | $565 \times 165$ | 1.05 | **0.553** | 0.95 | 2.13 | 2.52 | 50 | 0.1 |

method we solve the feasibility problem ($\mathbf{A}x \leq \mathbf{b}$) for a certain $\beta$ and $\delta$ [11], for GSKM and PASKM variants we solve the same feasibility problem and report the best performing method from each of the two, and finally for `fmincon` algorithms, we use the original LPs ($\min c^T x$ $s.t$ $Ax \leq b$, $l \leq x \leq u$). Note that, this is not an ideal or obvious comparison, for a better suitable comparison we follow the framework used in [11,60]. We set the stopping criterion for SKM, GSKM and PASKM variants as $\frac{\max(\mathbf{A}x_k - \mathbf{b})}{\max(\mathbf{A}x_0 - \mathbf{b})} \leq \epsilon$ and the halting criterion for the `fmincon`'s algorithms (IPM, ASM) are set as $\frac{\max(Ax_k - b, l - x_k, x_k - u)}{\max(Ax_0 - b, l - x_0, x_0 - u)} \leq \epsilon$ and $\frac{c^T x_k}{c^T x_0} \leq \epsilon$, where $\epsilon$ is the tolerance gap listed in Table 3. To avoid any biased conclusion, for each problem we set the initial update as far as possible from the feasible region.

From the comparison in Table 3, we can see that the proposed algorithms work much faster than IPM and ASM but work marginally better than the existing SKM method. Notice that the improvement of PASKM and GSKM algorithms over the SKM method for most problems are marginal as the proposed algorithms are designed explicitly for dense matrices. One can develop special algorithmic variants of the proposed PASKM and GSKM methods for sparse problems by following some standard aggregation techniques. A possible technique is to combine multiple steps by using the sparsity of the test instances. For instance, after $k^{th}$ iteration when we have $x_k, y_k$ and $v_k$, instead of moving forward with the sequences $x_{k+1}, y_{k+1}$ and $v_{k+1}$, for any $T \gg 1$ we can skip $T$ iterations and update $x_{k+T}, y_{k+T}$ and $v_{k+T}$ using a generalized recurrence relation that can enhance the computational efficiency.

# 5  Conclusion

In this work, we propose a general algorithmic framework (GSKM) for solving linear feasibility problems that unify various SKM type algorithms with the addition of a relaxation parameter $\xi$. From our convergence analysis of the GSKM method, one can recover convergence Theorems of several well-known algorithms such as Randomized Kaczmarz, Motzkin Method and Sampling Kaczmarz Motzkin method. In addition to the general framework, we propose a Nesterov type acceleration scheme in the SKM method called as PASKM. Our proposed PASKM method provides a bridge between Nesterov type acceleration of Machine Learning to sampling Kaczmarz methods for solving linear feasibility problems. To show the effectiveness of the proposed algorithms, we performed a wide range of numerical experiments on various types of random and standard benchmark data sets. For a better understanding of the behavior of the proposed algorithms, we numerically analyze two variants for both GSKM and PASKM algorithms in comparison with the original SKM method. Furthermore, we compare our proposed methods to commercially available methods such as IPM and ASM. In the majority of the test instances, the proposed algorithms significantly outperform the state-of-the-art methods. Furthermore, as shown in our numerical experiments, the correct choice of parameters can lead to much faster and accelerated methods for different types of test instances.

**Future Research**   In the future, the proposed algorithms and the technical analysis can be adopted effectively to various types of extensions such as sparse variants, optimally tuned PASKM, and GSKM, PASKM variants with greedy

---

[11] we note the best possible time from our previous experiments

sampling strategies. First, we plan to extend our work to design efficient sparse variations of the proposed methods that can handle large-scale real-world problems with greater sparsity in the data matrix $A$. Second, we intend to design a test instance dependent scheme for identifying optimal parameter selection (i.e., $\beta, \delta, \xi, \lambda, \tau$) for both GSKM and PASKM. For the GSKM algorithm, adaptive parameter selection (i.e., $\beta_k, \delta_k, \xi_k$) policy can be a great area of future research. One can also derive connecting ideas between the proposed GSKM and induced projection plane generation of Chubanov [30, 31] which can produce faster algorithms. Finally, we aspire to develop adaptive sampling strategies and integrate the greedy Kaczmarz [48] type method into the GSKM framework to further speed up the convergence.

## 6  Acknowledgements

## Appendix 1

**Proof of Lemma 3.5**    Using the expectation expression given in the first section we have,

$$\mathbb{E}_{\mathbb{S}}\left[a_{i*}a_{i*}^T\right] = \frac{1}{\binom{m}{\beta}} \sum_{j=0}^{m-\beta} \binom{\beta-1+j}{\beta-1} (A^T A)_{\underline{\mathbf{i_j}}}$$

$$\preceq \frac{\binom{m-1}{\beta-1}}{\binom{m}{\beta}} \sum_{j=0}^{m-\beta} (A^T A)_{\underline{\mathbf{i_j}}} \preceq \frac{\beta}{m} \sum_{i=1}^{m} a_i a_i^T = \frac{\beta}{m} A^T A.$$

Here, the notation $(A^T A)_{\underline{\mathbf{i_j}}}$ denotes the matrix $a_l a_l^T$ where the index $l$ belongs to the list (5). Furthermore, the index $l$ corresponds to the $(\beta + j)^{th}$ entry on the list (5). This proves the Lemma.

**Proof of Lemma 3.6**    Using the definition of the expectation from (7), we have,

$$\mathbb{E}_{\mathbb{S}}\left[\left|(a_{i*}^T x - b_{i*})^+\right|^2\right] \overset{(7)}{=} \frac{1}{\binom{m}{\beta}} \sum_{j=0}^{m-\beta} \binom{\beta-1+j}{\beta-1} |(Ax - b)_{\underline{\mathbf{i_j}}}^+|^2$$

$$\overset{\text{Lemma 3.2}}{\geq} \frac{1}{\binom{m}{\beta}} \sum_{j=0}^{m-\beta} \frac{\sum_{l=0}^{m-\beta} \binom{\beta-1+l}{\beta-1}}{m-\beta+1} |(Ax - b)_{\underline{\mathbf{i_j}}}^+|^2$$

$$\geq \frac{1}{m-\beta+1} \sum_{j=0}^{m-\beta} \left|(Ax - b)_{\underline{\mathbf{i_j}}}^+\right|^2$$

$$\geq \frac{1}{m-\beta+1} \min\{\frac{m-\beta+1}{m-s}, 1\} \|(Ax - b)^+\|^2$$

$$\overset{\text{Lemma 3.1}}{\geq} \frac{1}{mL^2} d(x, P)^2.$$

Here, $s$ is the number of zero entries in the residual $(Ax - b)^+$, which also corresponds to the number of satisfied constraints for $x$. Since $AP(x) \leq b$, we have the following:

$$\mathbb{E}_{\mathbb{S}}\left[\left|(a_{i*}^T x - b_{i*})^+\right|^2\right] \overset{(7)}{=} \frac{1}{\binom{m}{\beta}} \sum_{j=0}^{m-\beta} \binom{\beta-1+j}{\beta-1} |(Ax - b)_{\underline{\mathbf{i_j}}}^+|^2$$

$$\leq \frac{1}{\binom{m}{\beta}} \sum_{j=0}^{m-\beta} \binom{\beta-1+j}{\beta-1} \left|(Ax - AP(x))_{\underline{\mathbf{i_j}}}\right|^2$$

$$= \frac{1}{\binom{m}{\beta}} (x - P(x))^T \sum_{j=0}^{m-\beta} \binom{\beta-1+j}{\beta-1} (A^T A)_{\underline{\mathbf{i_j}}} (x - P(x))$$

$$= (x - \mathcal{P}(x))^T \, \mathbb{E}_{\mathbb{S}} \left[ a_{i*} a_{i*}^T \right] (x - \mathcal{P}(x))$$

$$\overset{\text{Lemma 3.4 \& 3.5}}{\leq} \min \left\{ 1, \frac{\beta}{m} \lambda_{\max} \right\} \| x - \mathcal{P}(x) \|^2$$

$$= \min \left\{ 1, \frac{\beta}{m} \lambda_{\max} \right\} d(x, P)^2.$$

Combining the above identities and using the expression for $f(x)$ from (9) we get,

$$\frac{\mu_1}{2} \, d(x, P)^2 \leq f(x) \leq \frac{\mu_2}{2} \, d(x, P)^2,$$

which proves the Lemma.

**Proof of Lemma 3.7**   From the definition of $f(x)$, it can be easily checked that $f(x)$ is a convex function. Now, by the convexity property of $f(x)$, for any $x, y \in \mathbb{R}^n$, we have the following:

$$\langle x - y, \nabla f(y) \rangle \leq f(x) - f(y). \tag{33}$$

Therefore, we have

$$\left\langle x - y, \mathbb{E}_{\mathbb{S}} \left[ (a_{i*}^T y - b_{i*})^+ a_{i*} \right] \right\rangle = \langle x - y, \nabla f(y) \rangle \leq f(x) - f(y)$$

$$\overset{\text{Lemma 3.6}}{\leq} \frac{\mu_2}{2} \, d(x, P)^2 - \frac{\mu_1}{2} \, d(y, P)^2.$$

This completes the proof.

**Proof of Lemma 3.9**   Since, $\bar{y} \in P$, from the definition we have,

$$\left\langle \bar{y} - y, \mathbb{E}_{\mathbb{S}} \left[ a_{i*} (a_{i*}^T y - b_{i*})^+ \right] \right\rangle = \mathbb{E}_{\mathbb{S}} \left[ (a_{i*}^T y - b_{i*})^+ \left( a_{i*}^T \bar{y} - a_{i*}^T y \right) \right]$$

$$\leq \mathbb{E}_{\mathbb{S}} \left[ (a_{i*}^T y - b_{i*})^+ \left( b_{i*} - a_{i*}^T y \right) \right]$$

$$= -\mathbb{E}_{\mathbb{S}} \left[ \left| (a_{i*}^T y - b_{i*})^+ \right|^2 \right]$$

$$= -2f(y) \overset{\text{Lemma 3.6}}{\leq} -\mu_1 \, d(y, P)^2.$$

Here, we used the identity $xx^+ = |x^+|^2$. This proves the Lemma.

**Proof of Lemma 3.11**   Since, $\mathcal{P}(x) \in P$, we have

$$\mathbb{E}_{\mathbb{S}} \left[ d(z, P)^2 \right] \overset{\text{Lemma 3.3}}{\leq} \mathbb{E}_{\mathbb{S}} \left[ \| z - \mathcal{P}(x) \|^2 \right] = \mathbb{E}_{\mathbb{S}} \left[ \left\| x - \mathcal{P}(x) - \delta \left( a_{i*}^T x - b_{i*} \right)^+ a_{i*} \right\|^2 \right]$$

$$\overset{(9)}{=} \| x - \mathcal{P}(x) \|^2 + 2\delta^2 f(x) + 2\delta \left\langle \mathcal{P}(x) - x, \nabla f(x) \right\rangle$$

$$\overset{\text{Lemma 3.6}}{\leq} \| x - \mathcal{P}(x) \|^2 - 2(2\delta - \delta^2) f(x)$$

$$\leq \| x - \mathcal{P}(x) \|^2 - (2\delta - \delta^2) \, \mu_1 \| x - \mathcal{P}(x) \|^2 = h(\delta) \, d(x, P)^2.$$

Here, we used the lower bound of the expected value from Lemma 3.6.

**Proof of Theorem 3.12**   Since, $\phi_1, \phi_2 \geq 0$, the largest root $\phi$ of equation $\phi^2 + \phi_1 \phi - \phi_2 = 0$ can written as

$$\phi = \frac{-\phi_1 + \sqrt{\phi_1^2 + 4\phi_2}}{2} \geq \frac{-\phi_1 + \phi_1}{2} = 0.$$

Then using the given recurrence we have,

$$G_{k+1} + \phi G_k \leq (\phi + \phi_1) G_k + \phi_2 G_{k-1}$$

$$= (\phi + \phi_1) \left( G_k + \phi G_{k-1} \right)$$

$$\vdots$$

$$\leq (\phi + \phi_1)^k \left( G_1 + \phi G_0 \right)$$

$$= (\phi + \phi_1)^k (1 + \phi) G_0.$$

This proves the first part. Also note that since $\phi_1 + \phi_2 < 1$, we have,

$$\phi + \phi_1 = \frac{\phi_1 + \sqrt{\phi_1^2 + 4\phi_2}}{2} < \frac{\phi_1 + \sqrt{\phi_1^2 + 4(1 - \phi_1)}}{2} = \frac{\phi_1 + 2 - \phi_1}{2} = 1.$$

For the second part, notice that from the recurrence inequality, we can deduce the following matrix inequality:

$$\begin{bmatrix} G_{k+1} \\ G_k \end{bmatrix} \leq \begin{bmatrix} \phi_1^2 + \phi_2 & \phi_1\phi_2 \\ \phi_1 & \phi_2 \end{bmatrix} \begin{bmatrix} G_{k-1} \\ G_{k-2} \end{bmatrix}. \tag{34}$$

The Jordan decomposition of the matrix in the above expression is given by,

$$\begin{bmatrix} \phi_1^2 + \phi_2 & \phi_1\phi_2 \\ \phi_1 & \phi_2 \end{bmatrix} = \begin{bmatrix} -\phi & \phi + \phi_1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \phi^2 & 0 \\ 0 & \rho^2 \end{bmatrix} \begin{bmatrix} \frac{-1}{\phi_1 + 2\phi} & \frac{1}{2} + \frac{\phi_1}{2(\phi_1 + 2\phi)} \\ \frac{1}{\phi_1 + 2\phi} & \frac{1}{2} - \frac{\phi_1}{2(\phi_1 + 2\phi)} \end{bmatrix}. \tag{35}$$

Next, we discuss two possible cases of values of $k$. Also, we substituted $\phi_2 = \phi(\phi + \phi_1)$ in the Jordan decomposition of equation (35).

**Case 1: $k$ even**

$$\begin{bmatrix} G_{k+1} \\ G_k \end{bmatrix} \overset{(34)}{\leq} \begin{bmatrix} \phi_1^2 + \phi^2 + \phi\phi_1 & \phi^2\phi_1 + \phi\phi_1^2 \\ \phi_1 & \phi^2 + \phi\phi_1 \end{bmatrix} \begin{bmatrix} G_{k-1} \\ G_{k-2} \end{bmatrix}$$

$$\vdots$$

$$\leq \begin{bmatrix} \phi_1^2 + \phi^2 + \phi\phi_1 & \phi^2\phi_1 + \phi\phi_1^2 \\ \phi_1 & \phi^2 + \phi\phi_1 \end{bmatrix}^{\frac{k}{2}} \begin{bmatrix} G_1 \\ G_0 \end{bmatrix}$$

$$\overset{(35)}{=} \begin{bmatrix} -\phi & \phi + \phi_1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \phi^k & 0 \\ 0 & \rho^k \end{bmatrix} \begin{bmatrix} \frac{-1}{\phi_1 + 2\phi} & \frac{1}{2} + \frac{\phi_1}{2(\phi_1 + 2\phi)} \\ \frac{1}{\phi_1 + 2\phi} & \frac{1}{2} - \frac{\phi_1}{2(\phi_1 + 2\phi)} \end{bmatrix} \begin{bmatrix} G_0 \\ G_0 \end{bmatrix}$$

$$= \begin{bmatrix} (1 + \phi)\rho^{k+1} + (1 - \phi - \phi_1)\phi^{k+1} \\ (1 + \phi)\rho^k - (1 - \phi - \phi_1)\phi^k \end{bmatrix} \begin{bmatrix} \frac{G_0}{\phi_1 + 2\phi} \end{bmatrix}$$

$$\overset{(22)}{=} \begin{bmatrix} R_1\rho^{k+1} + R_2\phi^{k+1} \\ R_1\rho^k - R_2\phi^k \end{bmatrix} G_0. \tag{36}$$

Here, we used $G_0 = G_1$.

**Case 2: $k$ odd**

$$\begin{bmatrix} G_{k+1} \\ G_k \end{bmatrix} \overset{(34)}{\leq} \begin{bmatrix} \phi_1^2 + \phi^2 + \phi\phi_1 & \phi^2\phi_1 + \phi\phi_1^2 \\ \phi_1 & \phi^2 + \phi\phi_1 \end{bmatrix} \begin{bmatrix} G_{k-1} \\ G_{k-2} \end{bmatrix}$$

$$\vdots$$

$$\leq \begin{bmatrix} \phi_1^2 + \phi^2 + \phi\phi_1 & \phi^2\phi_1 + \phi\phi_1^2 \\ \phi_1 & \phi^2 + \phi\phi_1 \end{bmatrix}^{\frac{k-1}{2}} \begin{bmatrix} G_2 \\ G_1 \end{bmatrix}$$

$$\overset{(35)}{=} \begin{bmatrix} -\phi & \phi + \phi_1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \phi^{k-1} & 0 \\ 0 & \rho^{k-1} \end{bmatrix} \begin{bmatrix} \frac{-1}{\phi_1 + 2\phi} & \frac{1}{2} + \frac{\phi_1}{2(\phi_1 + 2\phi)} \\ \frac{1}{\phi_1 + 2\phi} & \frac{1}{2} - \frac{\phi_1}{2(\phi_1 + 2\phi)} \end{bmatrix} \begin{bmatrix} (\phi_1 + \phi_2)G_0 \\ G_0 \end{bmatrix}$$

$$= \begin{bmatrix} (\phi + \phi_1 + \phi_2)\rho^k - (\phi - \phi_2)\phi^k \\ (\phi + \phi_1 + \phi_2)\rho^{k-1} + (\phi - \phi_2)\phi^{k-1} \end{bmatrix} \begin{bmatrix} \frac{G_0}{\phi_1 + 2\phi} \end{bmatrix}$$

$$\overset{(22)}{=} \begin{bmatrix} R_3\rho^k - R_4\phi^k \\ R_3\rho^{k-1} + R_4\phi^{k-1} \end{bmatrix} G_0. \tag{37}$$

Here, we used the inequality $G_2 \leq \phi_1 G_1 + \phi_2 G_0$. Now combining the relations from equation (36) and (37), we can prove the second part of Theorem 3.12.

**Proof of Theorem 3.13** From the given recurrence relation, we have

$$
\begin{bmatrix} H_{k+1} \\ F_{k+1} \end{bmatrix} \leq \begin{bmatrix} \Pi_1 & \Pi_2 \\ \Pi_3 & \Pi_4 \end{bmatrix} \begin{bmatrix} H_k \\ F_k \end{bmatrix} \leq \begin{bmatrix} \Pi_1 & \Pi_2 \\ \Pi_3 & \Pi_4 \end{bmatrix}^k \begin{bmatrix} H_1 \\ F_1 \end{bmatrix}.
\tag{38}
$$

Using the definitions of (25), we can write the Jordan decomposition of the above matrix as follows

$$
\begin{bmatrix} \Pi_1 & \Pi_2 \\ \Pi_3 & \Pi_4 \end{bmatrix} = \begin{bmatrix} \Gamma_2 & \Gamma_1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \rho_1 & 0 \\ 0 & \rho_2 \end{bmatrix} \begin{bmatrix} -\Gamma_3 & \Gamma_1\Gamma_3 \\ \Gamma_3 & \Gamma_2\Gamma_3 \end{bmatrix}.
\tag{39}
$$

Now, substituting the matrix decomposition into equation (38) and simplifying we have

$$
\begin{aligned}
\begin{bmatrix} H_{k+1} \\ F_{k+1} \end{bmatrix} &\leq \begin{bmatrix} \Pi_1 & \Pi_2 \\ \Pi_3 & \Pi_4 \end{bmatrix}^k \begin{bmatrix} H_1 \\ F_1 \end{bmatrix} \overset{(39)}{=} \begin{bmatrix} \Gamma_2 & \Gamma_1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \rho_1^k & 0 \\ 0 & \rho_2^k \end{bmatrix} \begin{bmatrix} -\Gamma_3 & \Gamma_1\Gamma_3 \\ \Gamma_3 & \Gamma_2\Gamma_3 \end{bmatrix} \begin{bmatrix} H_1 \\ F_1 \end{bmatrix} \\
&= \begin{bmatrix} \Gamma_2\Gamma_3(\Gamma_1 - 1)\,\rho_1^k + \Gamma_1\Gamma_3(\Gamma_2 + 1)\,\rho_2^k \\ \Gamma_3(\Gamma_1 - 1)\,\rho_1^k + \Gamma_3(\Gamma_2 + 1)\,\rho_2^k \end{bmatrix} \begin{bmatrix} H_1 \\ F_1 \end{bmatrix}.
\end{aligned}
\tag{40}
$$

Since $\Pi_1, \Pi_2, \Pi_3, \Pi_4 \geq 0$, one can easily verify that $\Gamma_1, \Gamma_3 \geq 0$. Now, it remains to show that $0 \leq \rho_1 \leq \rho_2 < 1$. To show that, first note that

$$
\begin{aligned}
(\Pi_1 - \Pi_4)^2 + 4\Pi_2\Pi_3 &\overset{(24)}{<} (\Pi_1 - \Pi_4)^2 + 4 - 4\Pi_1\Pi_4 - 4(\Pi_1 + \Pi_4) \\
&= (2 - \Pi_1 - \Pi_4)^2.
\end{aligned}
\tag{41}
$$

Now, we have,

$$
\begin{aligned}
\rho_1 &= \frac{1}{2}\left[\Pi_1 + \Pi_4 - \sqrt{(\Pi_1 - \Pi_4)^2 + 4\Pi_2\Pi_3}\right] \\
&\geq \frac{1}{2}\left[\Pi_1 + \Pi_4 - \sqrt{(\Pi_1 - \Pi_4)^2 + 4\Pi_1\Pi_4}\right] = \frac{1}{2}[\Pi_1 + \Pi_4 - (\Pi_1 + \Pi_4)] = 0.
\end{aligned}
$$

Moreover, since $2 - \Pi_1 - \Pi_4 \geq 0$, we have

$$
\begin{aligned}
\rho_1 \leq \rho_2 &= \frac{1}{2}\left[\Pi_1 + \Pi_4 + \sqrt{(\Pi_1 - \Pi_4)^2 + 4\Pi_2\Pi_3}\right] \\
&\overset{(41)}{<} \frac{1}{2}\left[\Pi_1 + \Pi_4 + \sqrt{(2 - \Pi_1 - \Pi_4)^2}\right] \\
&\overset{(24)}{=} \frac{1}{2}[\Pi_1 + \Pi_4 + 2 - \Pi_1 - \Pi_4] = 1.
\end{aligned}
$$

As $0 \leq \rho_1 \leq \rho_2 < 1$, considering (40) we can deduce that the sequence $\{H_k\}$ and $\{F_k\}$ converges.

# Appendix 2

**Proof of Theorem 3.14** From the update formula of Algorithm 2, we have $z_k = x_k - (A_{\tau_k} x_k - b_{\tau_k})_{i^*}^+ a_{i^*}$ where,

$$
i^* = \arg\max_{i \in \tau_k}\{a_i^T x_k - b_i, 0\} = \arg\max_{i \in \tau_k}(A_{\tau_k} x_k - b_{\tau_k})_i^+.
\tag{42}
$$

Similarly, the previous update formula can be written as, $z_{k-1} = x_{k-1} - (A_{\tau_{k-1}} x_{k-1} - b_{\tau_{k-1}})_{j^*}^+ a_{j^*}$; where,

$$
j^* = \arg\max_{j \in \tau_{k-1}}\{a_j^T x_{k-1} - b_j, 0\} = \arg\max_{j \in \tau_{k-1}}(A_{\tau_{k-1}} x_{k-1} - b_{\tau_{k-1}})_j^+.
\tag{43}
$$

Note that, the notation is consistent with the definition of (6). Since for any $\xi \in Q_1$, $(1 - \xi)\mathcal{P}(x_k) + \xi\mathcal{P}(x_{k-1}) \in P$ we have,

$$
\begin{aligned}
d(x_{k+1}, P)^2 &= \left\| x_{k+1} - \mathcal{P}(x_{k+1}) \right\|^2 \\
&\overset{\text{Lemma 3.3}}{\leq} \left\| x_{k+1} - (1 - \xi)\mathcal{P}(x_k) - \xi\mathcal{P}(x_{k-1}) \right\|^2 \\
&\overset{(10)}{=} \left\| (1 - \xi)z_k + \xi z_{k-1} - (1 - \xi)\mathcal{P}(x_k) - \xi\mathcal{P}(x_{k-1}) \right\|^2
\end{aligned}
$$

$$\overset{(11)}{=} \left\| (1-\xi)\left\{ x_k - \mathcal{P}(x_k) - \delta\left(a_{i*}^T x_k - b_{i*}\right)^+ a_{i*}\right\} \right.$$
$$\left. + \xi\left\{ x_{k-1} - \mathcal{P}(x_{k-1}) - \delta\left(a_{j*}^T x_{k-1} - b_{j*}\right)^+ a_{j*}\right\} \right\|^2$$
$$\leq (1-\xi)\left\| x_k - \mathcal{P}(x_k) - \delta\left(a_{i*}^T x_k - b_{i*}\right)^+ a_{i*}\right\|^2$$
$$+ \xi\left\| x_{k-1} - \mathcal{P}(x_{k-1}) - \delta\left(a_{j*}^T x_{k-1} - b_{j*}\right)^+ a_{j*}\right\|^2. \tag{44}$$

We used the fact that the function $\|\cdot\|^2$ is convex and $0 \leq \xi \leq 1$. Now, taking expectation in both sides of the equation (44) and using Lemma 3.11, we get the following:

$$\mathbb{E}[d(x_{k+1},P)^2 \mid \mathbb{S}_k, \mathbb{S}_{k-1}] \leq (1-\xi)\,\mathbb{E}_{\mathbb{S}_k}\left[\left\| x_k - \mathcal{P}(x_k) - \delta\left(a_{i*}^T x_k - b_{i*}\right)^+ a_{i*}\right\|^2\right]$$
$$+ \xi\,\mathbb{E}_{\mathbb{S}_{k-1}}\left[\left\| x_{k-1} - \mathcal{P}(x_{k-1}) - \delta\left(a_{j*}^T x_{k-1} - b_{j*}\right)^+ a_{j*}\right\|^2\right]$$
$$\overset{\text{Lemma 3.11}}{\leq} (1-\xi)\,h(\delta)\,d(x_k,P)^2 + \xi\,h(\delta)\,d(x_{k-1},P)^2, \tag{45}$$

where, $h(\delta)$ is defined in Lemma 3.11. Taking expectation again in equation (45) and letting $G_{k+1} = \mathbb{E}\left[d(x_{k+1},P)^2\right]$, we get the following:

$$G_{k+1} \leq \phi_1 G_k + \phi_2 G_{k-1}. \tag{46}$$

Since, $\phi_1, \phi_2 \geq 0$, $\phi_1 + \phi_2 < 1$ and $z_0 = z_1$, using first part of Theorem 3.12, we have the following:

$$\mathbb{E}\left[d(x_{k+1},P)^2\right] \leq (1+\phi)(\phi+\phi_1)^k G_0 = (1+\phi)\rho^k\,\mathbb{E}\left[d(x_0,P)^2\right]. \tag{47}$$

Moreover, considering (47) with Lemma 3.6 we get the bound of $\mathbb{E}[f(x_k)]$ which proves the first part of Theorem 3.14. Furthermore, using the second part of Theorem 3.12 and equation (46), we get the second part of Theorem 3.14. Now, to prove the third part first note that $\frac{1}{k}\sum_{l=1}^k \mathcal{P}(x_l) \in P$, using Lemma 3.3 we have

$$\mathbb{E}[d(\tilde{x}_k,P)^2] = \mathbb{E}[\|\tilde{x}_k - \mathcal{P}(\tilde{x}_k)\|^2] \overset{\text{Lemma 3.3}}{\leq} \mathbb{E}\left[\left\|\frac{1}{k}\sum_{l=1}^k (x_l - \mathcal{P}(x_l))\right\|^2\right]$$
$$\leq \mathbb{E}\left[\frac{1}{k}\sum_{l=1}^k \|x_l - \mathcal{P}(x_l)\|^2\right] = \frac{1}{k}\sum_{l=1}^k \mathbb{E}[d(x_l,P)^2]$$
$$\leq \frac{(1+\phi)d(x_0,P)^2}{k}\sum_{l=1}^k \rho^{l-1} \leq \frac{(1+\phi)d(x_0,P)^2}{k(1-\rho)}. \tag{48}$$

Furthermore, using a more simplifies version of (44) we have the following:

$$G_{l+1} - G_l \leq \xi(G_l - G_{l-1}) + 2\xi\delta(2-\delta)[f(x_l) - f(x_{l-1})] - 2\delta(2-\delta)f(x_l),$$

for any $l \geq 1$. Summing up the above identity for $l = 1, 2, ..., k$, we have the following:

$$2\delta(2-\delta)\sum_{l=1}^k f(x_l) \leq \xi G_0 + G_1 - \xi G_k - G_{k+1} + 2\xi\delta(2-\delta)[f(x_k) - f(x_0)]$$
$$\leq (1+\xi)G_0 + \xi[2\eta f(x_k) - G_k]$$
$$\leq (1+\xi)G_0 + \xi[\eta\mu_2 - 1]G_k \leq (1+\xi)d(x_0,P)^2, \tag{49}$$

where, $\eta = 2\delta - \delta^2$. We used the non-negativity of the sequences $G_k$ and $f(x_k)$. We also used the upper bound from Lemma 3.6. Then, we get

$$\mathbb{E}[f(\tilde{x}_k)] \leq \mathbb{E}\left[\frac{1}{k}\sum_{l=1}^k f(x_l)\right] = \frac{1}{k}\sum_{l=1}^k \mathbb{E}[f(x_l)] \leq \frac{(1+\xi)\,d(x_0,P)^2}{2\delta k(2-\delta)}.$$

This proves the second part of Theorem 3.14.

**Proof of Theorem 3.16** For any natural number $l \geq 1$ define, $\vartheta_l = \frac{\xi}{1+\xi}[x_{l-1} - x_l - \delta(a_{j^*}^T x_{l-1} - b_{j^*})^+ a_{j^*}]$, $\Delta_l = x_l + \vartheta_l$ and $\chi_l = \|x_l + \vartheta_l - \mathcal{P}(\Delta_l)\|^2$, then using the update formulas (10) and (11), we have

$$x_{l+1} + \vartheta_{l+1} \overset{(10) \,\&\, (11)}{=} x_l + \vartheta_l - \frac{\delta}{1+\xi}\left(a_{i^*}^T x_l - b_{i^*}\right)^+ a_{i^*},$$

here, the index $i^*$ and $j^*$ are defined based on (6) respectively for the sequences $x_l$ and $x_{l-1}$. Using the above relation, we can write

$$
\begin{aligned}
\chi_{l+1} = \|x_{l+1} + \vartheta_{l+1} - \mathcal{P}(\Delta_{l+1})\|^2 &\overset{\text{Lemma 3.3}}{\leq} \|x_{l+1} + \vartheta_{l+1} - \mathcal{P}(\Delta_l)\|^2 \\
&= \left\|x_l + \vartheta_l - \frac{\delta}{1+\xi}\left(a_{i^*}^T x_l - b_{i^*}\right)^+ a_{i^*} - \mathcal{P}(\Delta_l)\right\|^2 \\
&= \underbrace{\|x_l + \vartheta_l - \mathcal{P}(\Delta_l)\|^2}_{=\chi_l} + \frac{\delta^2}{(1+\xi)^2}\underbrace{\|(a_{i^*}^T x_l - b_{i^*})^+ a_{i^*}\|^2}_{J_1} \\
&\quad - \frac{2\delta}{1+\xi}\underbrace{\left\langle x_l + \vartheta_l - \mathcal{P}(\Delta_l),\ a_{i^*}(a_{i^*}^T x_l - b_{i^*})^+ \right\rangle}_{J_2} \\
&= \chi_l + \frac{\delta^2}{(1+\xi)^2}J_1 - \frac{2\delta}{1+\xi}J_2.
\end{aligned}
\tag{50}
$$

Taking expectation with respect to $\mathbb{S}_l$ we have,

$$\frac{\delta^2}{(1+\xi)^2}\mathbb{E}_{\mathbb{S}_l}[J_1] \overset{(9)}{=} \frac{2\delta^2}{(1+\xi)^2}f(x_l). \tag{51}$$

Similarly, we can simplify the third term of (50) as

$$
\begin{aligned}
&-\frac{2\delta}{1+\xi}\mathbb{E}_{\mathbb{S}_l}[J_2] \\
&\overset{(9)}{=} -\frac{2\delta}{1+\xi}\left\langle x_l - \mathcal{P}(\Delta_l), \nabla f(x_l)\right\rangle - \frac{2\delta\xi}{(1+\xi)^2}\left\langle x_{l-1} - x_l - \delta\nabla f(x_{l-1}), \nabla f(x_l)\right\rangle \\
&= -\frac{2\delta}{1+\xi}\left\langle x_l - \mathcal{P}(\Delta_l), \nabla f(x_l)\right\rangle - \frac{2\delta\xi}{(1+\xi)^2}\left\langle x_{l-1} - x_l, \nabla f(x_l)\right\rangle \\
&\quad + \frac{\delta^2\xi}{(1+\xi)^2}\left[\|\nabla f(x_l) + \nabla f(x_{l-1})\|^2 - \|\nabla f(x_l)\|^2 - \|\nabla f(x_{l-1})\|^2\right] \\
&\overset{\text{Lemma 3.7 \& 3.9}}{\leq} -\frac{4\delta}{1+\xi}f(x_l) - \frac{2\delta\xi}{(1+\xi)^2}\left[f(x_{l-1}) - f(x_l)\right] - \frac{2\delta^2\xi}{(1+\xi)^2}\left[f(x_{l-1}) + f(x_l)\right] \\
&= -\frac{2\delta\xi(1+\delta)}{(1+\xi)^2}f(x_{l-1}) + \frac{2\delta\xi(1+\delta)}{(1+\xi)^2}f(x_l) - \frac{4\delta(1+\xi+\delta\xi)}{(1+\xi)^2}f(x_l).
\end{aligned}
\tag{52}
$$

Using the expressions of equation (51) and (52) in (50) and simplifying further, we have

$$\mathbb{E}[\chi_{l+1}] - \frac{2\delta\xi(1+\delta)}{(1+\xi)^2}f(x_l) + \varpi f(x_l) \leq \mathbb{E}[\chi_l] - \frac{2\delta\xi(1+\delta)}{(1+\xi)^2}f(x_{l-1}), \tag{53}$$

here,

$$\varpi = \frac{4\delta(1+\xi+\delta\xi)}{(1+\xi)^2} - \frac{2\delta^2}{(1+\xi)^2} = \frac{2\delta(2+2\xi+2\delta\xi-\delta)}{(1+\xi)^2} > 0. \tag{54}$$

Now, taking expectation again in (53) and using the tower property, we get,

$$q_{l+1} + \varpi\,\mathbb{E}[f(x_l)] \leq q_l, \quad l = 1, 2, 3..., \tag{55}$$

where, $q_l = \mathbb{E}[\chi_l] - \frac{2\delta\xi(1+\delta)}{(1+\xi)^2}\mathbb{E}[f(x_{l-1})]$. Summing up (55) for $l = 1, 2, ..., k$ we get

$$\sum_{l=1}^{k}\mathbb{E}[f(x_l)] \leq \frac{q_1 - q_{k+1}}{\varpi} \leq \frac{q_1}{\varpi}. \tag{56}$$

33

Now, using Jensen's inequality, we have

$$\mathbb{E}\left[f(\bar{x_k})\right] = \mathbb{E}\left[f\left(\sum_{l=1}^{k}\frac{x_l}{k}\right)\right] \leq \mathbb{E}\left[\frac{1}{k}\sum_{l=1}^{k}f(x_l)\right] = \frac{1}{k}\sum_{l=1}^{k}\mathbb{E}[f(x_l)] \overset{(56)}{\leq} \frac{q_1}{\varpi k}.$$

Since, $x_0 = x_1$, we have $\vartheta_1 = \frac{-\delta\xi}{1+\xi}(a_{i*}^T x_0 - b_{i*})^+ a_{i*}$. Furthermore,

$$\mathbb{E}[\chi_1] = \mathbb{E}\left[\|x_1 + \vartheta_1 - \mathcal{P}(\Delta_1)\|^2\right] \overset{\text{Lemma 3.3}}{\leq} \mathbb{E}\left[\|x_1 + \vartheta_1 - \mathcal{P}(x_0)\|^2\right]$$

$$= \mathbb{E}\left[\|x_0 - \mathcal{P}(x_0) - \frac{\delta\xi}{1+\xi}(a_{i*}^T x_0 - b_{i*})^+ a_{i*}\|^2\right]$$

$$= \|x_0 - \mathcal{P}(x_0)\|^2 + \frac{\delta^2\xi^2}{(1+\xi)^2}\mathbb{E}[|(a_{i*}^T x_0 - b_{i*})^+|^2]$$

$$- \frac{2\delta\xi}{1+\xi}\langle x_0 - \mathcal{P}(x_0), \mathbb{E}[(a_{i*}^T x_0 - b_{i*})^+ a_{i*}]\rangle$$

$$\overset{\text{Lemma 3.6}}{\leq} \|x_0 - \mathcal{P}(x_0)\|^2 + \frac{2\delta^2\xi^2}{(1+\xi)^2}f(x_0) - \frac{2\delta\xi\mu_2}{1+\xi}\|x_0 - \mathcal{P}(x_0)\|^2. \tag{57}$$

Now, from our construction we get

$$q_1 = \mathbb{E}[\chi_1] - \frac{2\delta\xi(1+\delta)}{(1+\xi)^2}\mathbb{E}[f(x_0)] \leq (1 - \frac{2\delta\xi\mu_2}{1+\xi})\,d(x_0,P)^2 + \frac{2\delta\xi(\delta\xi - 1 - \delta)}{(1+\xi)^2}f(x_0).$$

Substituting the values of $\varpi$ and $q_1$ in the expression of $\mathbb{E}\left[f(\bar{x_k})\right]$, we have the following:

$$\mathbb{E}\left[f(\bar{x}_k)\right] \leq \frac{(1+\xi)(1+\xi - 2\delta\xi\mu_2)\,d(x_0,P)^2 + 2\xi\delta(\delta\xi - \delta - 1)f(x_0)}{2\delta k\,(2 + 2\xi + 2\delta\xi - \delta)}.$$

**Proof of Theorem 3.15** Since, the term $\|x_{k+1} - \mathcal{P}(x_{k+1})\|$ is constant under From the update formula of the GSKM algorithm, we get,

$$\mathbb{E}[\|x_{k+1} - \mathcal{P}(x_{k+1})\| \mid \mathbb{S}_{k+1}, \mathbb{S}_k] = \mathbb{E}[\|x_{k+1} - \mathcal{P}(x_{k+1})\| \mid \mathbb{S}_k]$$

$$\overset{\text{Lemma 3.3}}{\leq} \mathbb{E}_{\mathbb{S}_k}[\|x_{k+1} - \mathcal{P}(x_k)\|]$$

$$= \mathbb{E}_{\mathbb{S}_k}[\|z_k - \mathcal{P}(x_k) - \xi(z_k - z_{k-1})\|]$$

$$\leq \mathbb{E}_{\mathbb{S}_k}[\|z_k - \mathcal{P}(x_k)\|] + |\xi|\,\mathbb{E}_{\mathbb{S}_k}[\|z_k - z_{k-1}\|]$$

$$\leq \left\{\mathbb{E}_{\mathbb{S}_k}[\|z_k - \mathcal{P}(x_k)\|^2]\right\}^{\frac{1}{2}} + |\xi|\|z_k - z_{k-1}\|$$

$$\overset{\text{Lemma 3.11}}{\leq} \sqrt{h(\delta)}\,\|x_k - \mathcal{P}(x_k)\| + |\xi|\|z_k - z_{k-1}\|. \tag{58}$$

We performed the two expectations in order, from the innermost to the outermost. Now, taking expectation in (58) and using the tower property of expectation we have,

$$\mathbb{E}[\|x_{k+1} - \mathcal{P}(x_{k+1})\|] \leq \sqrt{h(\delta)}\,\mathbb{E}[\|x_k - \mathcal{P}(x_k)\|] + |\xi|\,\mathbb{E}[\|z_k - z_{k-1}\|]. \tag{59}$$

Similarly, using the update formula for $z_{k+1}$, we have

$$\mathbb{E}[\|z_{k+1} - z_k\| \mid \mathbb{S}_{k+1}, \mathbb{S}_k] = \mathbb{E}[\mathbb{E}_{\mathbb{S}_{k+1}}[\|x_{k+1} - \delta\left(a_{i*}^T x_{k+1} - b_{i*}\right)^+ a_{i*} - z_k\|] \mid \mathbb{S}_k]$$

$$= \mathbb{E}[\mathbb{E}_{\mathbb{S}_{k+1}}[\| - \xi(z_k - z_{k-1}) - \delta\left(a_{i*}^T x_{k+1} - b_{i*}\right)^+ a_{i*}\|] \mid \mathbb{S}_k]$$

$$\leq |\xi|\,\|z_k - z_{k-1}\| + \delta\,\mathbb{E}[\mathbb{E}_{\mathbb{S}_{k+1}}[|(a_{i*}^T x_{k+1} - b_{i*})^+|] \mid \mathbb{S}_k]$$

$$\leq |\xi|\,\|z_k - z_{k-1}\| + \delta\,\mathbb{E}[\{\mathbb{E}_{\mathbb{S}_{k+1}}[|(a_{i*}^T x_{k+1} - b_{i*})^+|^2]\}^{\frac{1}{2}} \mid \mathbb{S}_k]$$

$$\overset{\text{Lemma 3.6}}{\leq} |\xi|\,\|z_k - z_{k-1}\| + \delta\sqrt{\mu_2}\,\mathbb{E}[\|x_{k+1} - \mathcal{P}(x_{k+1})\| \mid \mathbb{S}_k]. \tag{60}$$

Taking expectation in (60) and using (59) along with the tower property, we have,

$$\mathbb{E}[\|z_{k+1} - z_k\|] \leq |\xi|\,\mathbb{E}[\|z_k - z_{k-1}\|] + \delta\sqrt{\mu_2}\,\mathbb{E}[\|x_{k+1} - \mathcal{P}(x_{k+1})\|]$$

$$\overset{(59)}{\leq} |\xi| \left(1 + \delta\sqrt{\mu_2}\right) \mathbb{E}[\|z_k - z_{k-1}\|] + \delta\sqrt{\mu_2 h(\delta)}\, \mathbb{E}[\|x_k - \mathcal{P}(x_k)\|]. \tag{61}$$

Combining both (59) and (61), we can deduce the following matrix inequality:

$$\mathbb{E}\begin{bmatrix} \|x_{k+1} - \mathcal{P}(x_{k+1})\| \\ \|z_{k+1} - z_k\| \end{bmatrix} \leq \begin{bmatrix} \sqrt{h(\delta)} & |\xi| \\ \delta\sqrt{\mu_2 h(\delta)} & |\xi|\left(1 + \delta\sqrt{\mu_2}\right) \end{bmatrix} \mathbb{E}\begin{bmatrix} \|x_k - \mathcal{P}(x_k)\| \\ \|z_k - z_{k-1}\| \end{bmatrix}. \tag{62}$$

Now, from the definition, it can be easily checked that $\Pi_1, \Pi_2, \Pi_3, \Pi_4 \geq 0$. Since, $\xi \in Q_2$, we have

$$\Pi_2\Pi_3 - \Pi_1\Pi_4 = |\xi|\delta\sqrt{\mu_2 h(\delta)} - |\xi|\sqrt{h(\delta)} - |\xi|\delta\sqrt{\mu_2 h(\delta)} = -|\xi|\sqrt{h(\delta)} \leq 0. \tag{63}$$

Also, we have

$$\Pi_1 + \Pi_4 - \Pi_1\Pi_4 + \Pi_2\Pi_3 = \sqrt{h(\delta)} + |\xi|\left(1 + \delta\sqrt{\mu_2}\right) - |\xi|\sqrt{h(\delta)} < 1. \tag{64}$$

Here, in the last inequality we used the given condition. Considering (64), we can check that $\Pi_1 + \Pi_4 < 1 + |\xi|\sqrt{h(\delta)} = 1 + \min\{1, |\xi|\sqrt{h(\delta)}\} = 1 + \min\{1, \Pi_1\Pi_4 - \Pi_2\Pi_3\}$. Also from (63), we have $\Pi_2\Pi_3 - \Pi_1\Pi_4 \leq 0$, which is precisely the condition provided in (24). Let's define the sequences $F_k = \mathbb{E}[\|z_k - z_{k-1}\|]$ and $H_k = \mathbb{E}[\|x_k - \mathcal{P}(x_k)\|]$. Now, using Theorem 3.13, we have

$$\begin{bmatrix} H_{k+1} \\ F_{k+1} \end{bmatrix} \leq \begin{bmatrix} \Gamma_2\Gamma_3(\Gamma_1 - 1)\,\rho_1^k + \Gamma_1\Gamma_3(\Gamma_2 + 1)\,\rho_2^k \\ \Gamma_3(\Gamma_1 - 1)\,\rho_1^k + \Gamma_3(\Gamma_2 + 1)\,\rho_2^k \end{bmatrix} \begin{bmatrix} H_1 \\ F_1 \end{bmatrix}. \tag{65}$$

where, $\Gamma_1, \Gamma_2, \Gamma_3, \rho_1, \rho_2$ can be derived from (25) using the parameter choice of (27). Note that, from the GSKM algorithm we have, $x_1 = x_0$ and $z_1 = z_0$. Therefore we can easily check that, $F_1 = \mathbb{E}[\|z_1 - z_0\|] = 0$ and $H_1 = \mathbb{E}[\|x_1 - \mathcal{P}(x_1)\|] = \mathbb{E}[\|x_0 - \mathcal{P}(x_0)\|] = \|x_0 - \mathcal{P}(x_0)\| = H_0$. Now, substituting the values of $H_1$ and $F_1$ in (65), we have

$$\begin{bmatrix} H_{k+1} \\ F_{k+1} \end{bmatrix} = \mathbb{E}\begin{bmatrix} d(x_{k+1}, P) \\ \|z_{k+1} - z_k\| \end{bmatrix} \leq \begin{bmatrix} -\Gamma_2\Gamma_3\,\rho_1^k + \Gamma_1\Gamma_3\,\rho_2^k \\ -\Gamma_3\,\rho_1^k + \Gamma_3\,\rho_2^k \end{bmatrix} d(x_0, P). \tag{66}$$

Also from Theorem 3.13 we have, $\Gamma_1, \Gamma_3 \geq 0$ and $0 \leq \rho_1 \leq \rho_2 < 1$. Which proves the Theorem.

**Proof of Theorem 3.20**   Note that, since $Ax \leq b$ is feasible, then from Lemma 3.19, we know that there is a feasible solution $x^*$ with $|x_j^*| \leq \frac{2^\sigma}{2n}$ for $j = 1, ..., n$. Thus, we have,

$$d(x_0, P) = \|x_0 - \mathcal{P}(x_0)\| \leq \|x^*\| \leq \frac{2^{\sigma-1}}{\sqrt{n}}, \tag{67}$$

as $x_0 = 0$. Then if the system $Ax \leq b$ is infeasible, by using Lemma 3.17, we have,

$$\theta(x) \geq 2^{1-\sigma}.$$

This implies when GSKM runs on the system $Ax \leq b$, the system is feasible when $\theta(x) < 2^{1-\sigma}$. Furthermore, since every point of the feasible region $P$ is inside the half-space defined by $\tilde{H}_i = \{x \mid a_i^T x \leq b_i\}$ for all $i = 1, 2, ..., m$, we have the following:

$$\theta(x) = \left[\max_i \{a_i^T x - b_i\}\right]^+ \leq \|a_i^T(x - \mathcal{P}(x))\| \leq d(x, P). \tag{68}$$

Then, for $\xi \in Q_1$ whenever the system $Ax \leq b$ is feasible, we have,

$$\mathbb{E}[\theta(x_k)] \overset{(68)}{\leq} \mathbb{E}[d(x_{k+1}, P)] \leq \sqrt{\mathbb{E}[d(x_{k+1}, P)^2]} \overset{\text{Theorem 3.14}}{\leq} \sqrt{1 + \phi}\rho^{\frac{k}{2}} d(x_0, P). \tag{69}$$

Similarly for $\xi \in Q_2$ whenever the system $Ax \leq b$ is feasible, we have,

$$\mathbb{E}[\theta(x_k)] \overset{(68)}{\leq} \mathbb{E}[d(x_{k+1}, P)] \overset{\text{Theorem 3.15}}{\leq} \sqrt{1 + \phi}\,\rho_2^k\, d(x_0, P). \tag{70}$$

Take, $\bar{\rho} = \max\{\rho, \rho_2^2\}$ [12]. Now combining (69) and (70), for any $\xi \in Q = Q_1 \cup Q_2$, whenever the system $Ax \leq b$ is feasible, we have,

$$\mathbb{E}[\theta(x_k)] \overset{(69)\,\&\,(70)}{\leq} \sqrt{1 + \phi}\,\bar{\rho}^{\frac{k}{2}}\, d(x_0, P) \overset{(67)}{\leq} \sqrt{1 + \phi}\,\bar{\rho}^{\frac{k}{2}}\, \frac{2^{\sigma-1}}{\sqrt{n}}. \tag{71}$$

---

[12] Note that, since $\Gamma_1\Gamma_2 \leq 0$ and $\Gamma_1\Gamma_3 \leq 1 \leq \sqrt{(1 + \phi)}$, from Theorem 3.15 we have $\mathbb{E}[d(x_{k+1}, P)] \leq \sqrt{(1 + \phi)}\,\rho_2^k\, d(x_0, P)$ for any $\xi \in Q_2$.

Here, we used Theorems 3.14 & 3.15 and the identities from equations (67) & (68). Now, for detecting feasibility we need to have, $\mathbb{E}[\theta(x_k)] < 2^{1-\sigma}$. That gives us,

$$\sqrt{1+\phi}\,\bar{\rho}^{\frac{k}{2}}\,\frac{2^{\sigma-1}}{\sqrt{n}} < 2^{1-\sigma}.$$

Simplifying the above identity further we get the following lower bound for $k$:

$$k > \frac{4\sigma - 4 - \log n + \log(1+\phi)}{\log\left(\frac{1}{\bar{\rho}}\right)}.$$

Moreover, if the system $Ax \leq b$ is feasible, then the probability of not having a certificate of feasibility is bounded as follows,

$$p = \mathbb{P}\left(\theta(x_k) \geq 2^{1-\sigma}\right) \leq \frac{\mathbb{E}\left[\theta(x_k)\right]}{2^{1-\sigma}} < \sqrt{\frac{1+\phi}{n}}\,2^{2\sigma-2}\,\bar{\rho}^{\frac{k}{2}}.$$

Here, we used the Markov's inequality $\mathbb{P}(x \geq t) \leq \frac{\mathbb{E}[x]}{t}$. This completes the proof of Theorem 3.20.

## Appendix 3

**Proof of Theorem 3.22**    From the update formula of the PASKM algorithm, we get,

$$\mathbb{E}_{\mathbb{S}_k}[\|v_{k+1} - \mathcal{P}(v_{k+1})\|^2]$$

$$\overset{\text{Lemma 3.3}}{\leq} \mathbb{E}_{\mathbb{S}_k}[\|v_{k+1} - \omega\mathcal{P}(v_k) - (1-\omega)\mathcal{P}(y_k)\|^2]$$

$$= \mathbb{E}_{\mathbb{S}_k}[\|\omega(v_k - \mathcal{P}(v_k)) + (1-\omega)(y_k - \mathcal{P}(y_k)) - \gamma\left(a_{i*}^T y_k - b_{i*}\right)^+ a_{i*}\|^2]$$

$$= \mathbb{E}_{\mathbb{S}_k}[\|\omega(v_k - \mathcal{P}(v_k)) + (1-\omega)(y_k - \mathcal{P}(y_k))\|^2] + \gamma^2\,\mathbb{E}_{\mathbb{S}_k}[|(a_{i*}^T y_k - b_{i*})^+|^2]$$

$$- 2\gamma(1-\omega)\langle y_k - \mathcal{P}(y_k), \mathbb{E}_{\mathbb{S}_k}[(a_{i*}^T y_k - b_{i*})^+ a_{i*}]\rangle$$

$$- 2\gamma\omega\langle v_k - \mathcal{P}(v_k), \mathbb{E}_{\mathbb{S}_k}[(a_{i*}^T y_k - b_{i*})^+ a_{i*}]\rangle$$

$$\leq \omega\|v_k - \mathcal{P}(v_k)\|^2 + (1-\omega)\|y_k - \mathcal{P}(y_k)\|^2 + \gamma^2\,\mathbb{E}_{\mathbb{S}_k}[|(a_{i*}^T y_k - b_{i*})^+|^2]$$

$$- 2\gamma(1-\omega)\,\mathbb{E}_{\mathbb{S}_k}[|(a_{i*}^T y_k - b_{i*})^+|^2] + \omega\gamma\,\mathbb{E}_{\mathbb{S}_k}[|(a_{i*}^T y_k - b_{i*})^+|^2] + \omega\gamma\|v_k - \mathcal{P}(v_k)\|^2$$

$$= \omega(1+\gamma)\|v_k - \mathcal{P}(v_k)\|^2 + (1-\omega)\|y_k - \mathcal{P}(y_k)\|^2 + 2\gamma(\gamma + 3\omega - 2)f(y_k)$$

$$\leq \omega(1+\gamma)\|v_k - \mathcal{P}(v_k)\|^2 + \{1 - \omega + \gamma\mu_1(\gamma + 3\omega - 2)\}\|y_k - \mathcal{P}(y_k)\|^2. \tag{72}$$

Here, we used the condition $\gamma + 3\omega - 2 \leq 0$. Similarly, using the update formula for $y_{k+1}$, we have

$$\mathbb{E}_{\mathbb{S}_k}[\|y_{k+1} - \mathcal{P}(y_{k+1})\|^2]$$

$$\overset{\text{Lemma 3.3}}{\leq} \mathbb{E}_{\mathbb{S}_k}[\|\alpha(v_{k+1} - \mathcal{P}(v_{k+1})) + (1-\alpha)(x_{k+1} - \mathcal{P}(y_k))\|^2]$$

$$\leq \alpha\,\mathbb{E}_{\mathbb{S}_k}[\|v_{k+1} - \mathcal{P}(v_{k+1})\|^2] + (1-\alpha)\,\mathbb{E}_{\mathbb{S}_k}[\|x_{k+1} - \mathcal{P}(y_k)\|^2]$$

$$\overset{\text{Lemma 3.11}}{\leq} \alpha\,\mathbb{E}_{\mathbb{S}_k}[\|v_{k+1} - \mathcal{P}(v_{k+1})\|^2] + (1-\alpha)h(\delta)\|y_k - \mathcal{P}(y_k)\|^2. \tag{73}$$

Following Theorem 3.13, let us define the sequences $H_k = E[\|v_k - \mathcal{P}(v_k)\|^2]$ and $F_k = \mathbb{E}[\|y_k - \mathcal{P}(y_k)\|^2]$. The goal is to prove that $H_k$ and $F_k$ satisfy the condition (24). Now, taking expectation in (72) and using the tower property of expectation we have,

$$H_{k+1} \leq \omega(1+\gamma)H_k + \{1 - \omega + \gamma\mu_1(\gamma + 3\omega - 2)\}F_k. \tag{74}$$

Similarly, taking expectation in (73) and using (74) along with the tower property of expectation we have,

$$F_{k+1} \leq \alpha H_{k+1} + (1-\alpha)h(\delta)F_k$$

$$\leq \alpha\omega(1+\gamma)H_k + \{(1-\alpha)h(\delta) + \alpha(1-\omega) + \alpha\gamma\mu_1(\gamma + 3\omega - 2)\}F_k. \tag{75}$$

Combining both (74) and (75), we can deduce the following matrix inequality:

$$\begin{bmatrix} H_{k+1} \\ F_{k+1} \end{bmatrix} \leq \begin{bmatrix} \Pi_1 & \Pi_2 \\ \Pi_3 & \Pi_4 \end{bmatrix} \begin{bmatrix} H_k \\ F_k \end{bmatrix} \leq \begin{bmatrix} \Pi_1 & \Pi_2 \\ \Pi_3 & \Pi_4 \end{bmatrix}^{k+1} \begin{bmatrix} H_0 \\ F_0 \end{bmatrix}. \tag{76}$$

Here, we use the fact that $\Pi_1, \Pi_2, \Pi_3, \Pi_4 \geq 0$. Now we will use Theorem 3.13 to simplify the expression of (76). Before we can use Theorem 3.13, we need to make sure the sequences $H_k$ and $F_k$ satisfy the condition of (24). From the definition, we have

$$
\begin{aligned}
\Pi_2\Pi_3 - \Pi_1\Pi_4 &= \alpha\omega(1-\omega)(1+\gamma) + \alpha\omega\gamma\mu_1(1+\gamma)(\gamma+3w-2) \\
&\quad - \omega h(\delta)(1-\alpha)(1+\gamma) - \alpha\omega(1-\omega)(1+\gamma) - \alpha\omega\gamma\mu_1(1+\gamma)(\gamma+3w-2) \\
&= -\omega h(\delta)(1-\alpha)(1+\gamma) \leq 0.
\end{aligned}
\tag{77}
$$

Also, we have

$$
\begin{aligned}
\Pi_1 + \Pi_4 - \Pi_1\Pi_4 + \Pi_2\Pi_3 &= \omega(1+\gamma) + h(\delta)(1-\alpha) + \alpha(1-\omega) \\
&\quad + \alpha\gamma\mu_1(\gamma+3w-2) - \omega h(\delta)(1-\alpha)(1+\gamma) < 1.
\end{aligned}
\tag{78}
$$

Here, in the last inequality, we used the given condition. Considering (78), we can check that $\Pi_1 + \Pi_4 < 1 + \omega h(\delta)(1-\alpha)(1+\gamma) = 1 + \min\{1, \omega h(\delta)(1-\alpha)(1+\gamma)\} = 1 + \min\{1, \Pi_1\Pi_4 - \Pi_2\Pi_3\}$. Also from (77), we have $\Pi_2\Pi_3 - \Pi_1\Pi_4 \leq 0$, which is precisely the condition provided in (24). Now, using Theorem 3.13, we have

$$
\begin{bmatrix} H_{k+1} \\ F_{k+1} \end{bmatrix} \leq \begin{bmatrix} \Gamma_2\Gamma_3(\Gamma_1-1)\,\rho_1^{k+1} + \Gamma_1\Gamma_3(\Gamma_2+1)\,\rho_2^{k+1} \\ \Gamma_3(\Gamma_1-1)\,\rho_1^{k+1} + \Gamma_3(\Gamma_2+1)\,\rho_2^{k+1} \end{bmatrix} \begin{bmatrix} H_0 \\ F_0 \end{bmatrix}.
\tag{79}
$$

where, $\Gamma_1, \Gamma_2, \Gamma_3, \rho_1, \rho_2$ can be derived from (25) using the given parameter. Note that, from the PASKM algorithm we have, $x_0 = v_0 = y_0$. Therefore we can easily check that, $H_0 = \|v_0 - \mathcal{P}(v_0)\|^2 = \|y_0 - \mathcal{P}(y_0)\|^2 = F_0$. Now, substituting the values of $H_0$ and $F_0$ in (79), we have

$$
\mathbb{E}\begin{bmatrix} d(v_{k+1}, P)^2 \\ d(y_{k+1}, P)^2 \end{bmatrix} \leq \begin{bmatrix} \Gamma_2\Gamma_3(\Gamma_1-1)\,\rho_1^{k+1} + \Gamma_1\Gamma_3(\Gamma_2+1)\,\rho_2^{k+1} \\ \Gamma_3(\Gamma_1-1)\,\rho_1^{k+1} + \Gamma_3(\Gamma_2+1)\,\rho_2^{k+1} \end{bmatrix} d(y_0, P)^2.
\tag{80}
$$

Also from Theorem 3.13 we have, $\Gamma_1, \Gamma_3 \geq 0$ and $0 \leq \rho_1 \leq \rho_2 < 1$. Which proves the the first part of the Theorem. Now, considering Lemma 3.6, we get

$$
\mathbb{E}[f(x_{k+1})] \leq \frac{\mu_2}{2}\,\mathbb{E}[\|y_{k+1} - \mathcal{P}(y_{k+1})\|^2] = \frac{\mu_2}{2}\,\mathbb{E}[d(y_{k+1}, P)^2].
\tag{81}
$$

Now, substituting the result of (80) in (81), we get the second part of the Theorem.

**Proof of Theorem 3.23**   Let us define, $\mathcal{V} = \omega\mathcal{P}(v_k) + (1-\omega)\mathcal{P}(y_k)$. Since $\mathcal{V} \in P$, using the update formula of $v_{k+1}$ from equation (15), we have,

$$
\begin{aligned}
d(v_{k+1}, P)^2 = \|v_{k+1} - \mathcal{P}(v_{k+1})\|^2 &\overset{\text{Lemma 3.3}}{\leq} \|v_{k+1} - \mathcal{V}\|^2 \\
&\overset{(15)}{=} \|\omega v_k + (1-\omega)y_k - \mathcal{V} - \gamma(a_{i^*}^T y_k - b_{i^*})^+ a_{i^*}\|^2 \\
&= \underbrace{\|\omega v_k + (1-\omega)y_k - \mathcal{V}\|^2}_{I_1} + \gamma^2 \underbrace{\|(a_{i^*}^T y_k - b_{i^*})^+ a_{i^*}\|^2}_{I_2} \\
&\quad - 2\gamma \underbrace{\langle \omega v_k + (1-\omega)y_k - \mathcal{V},\; a_{i^*}(a_{i^*}^T y_k - b_{i^*})^+ \rangle}_{I_3} \\
&= I_1 + \gamma^2 I_2 - 2\gamma I_3.
\end{aligned}
\tag{82}
$$

Since $\|\cdot\|^2$ is a convex function and $0 < \omega < 1$, we can bound the expected first term as follows,

$$
\begin{aligned}
\mathbb{E}_{\mathbb{S}_k}[I_1] = \mathbb{E}_{\mathbb{S}_k}\left[\|\omega v_k + (1-\omega)y_k - \mathcal{V}\|^2\right] &= \mathbb{E}_{\mathbb{S}_k}\left[\|\omega v_k + (1-\omega)y_k - \omega\mathcal{P}(v_k) - (1-\omega)\mathcal{P}(y_k)\|^2\right] \\
&\leq \omega\|v_k - \mathcal{P}(v_k)\|^2 + (1-\omega)\|y_k - \mathcal{P}(y_k)\|^2 \\
&= \omega\,d(v_k, P)^2 + (1-\omega)\,d(y_k, P)^2.
\end{aligned}
\tag{83}
$$

Taking expectation with respect to the sampling distribution in the second term of equation (82) and using Lemma 3.11 with the choice $z = x_{k+1}$, $x = y_k$ and $\eta = 2\delta - \delta^2$, we get,

$$
\gamma^2\,\mathbb{E}_{\mathbb{S}_k}\left[\|(a_{i^*}^T y_k - b_{i^*})^+ a_{i^*}\|^2\right] = \gamma^2\,\mathbb{E}_{\mathbb{S}_k}\left[|(a_{i^*}^T y_k - b_{i^*})^+|^2\right]
$$

$$\overset{\text{Lemma 3.11}}{\leq} \frac{\gamma^2}{\eta} \left[ d(y_k, P)^2 - \mathbb{E}\left[ d(x_{k+1}, P)^2 \right] \right]. \tag{84}$$

Now, taking expectation in the third term of (82) we get,

$$-2\gamma \mathbb{E}_{\mathbb{S}_k}[I_3] = -2\gamma \langle \omega v_k + (1-\omega)y_k - \mathcal{V}, \mathbb{E}_{\mathbb{S}_k}\left[ a_{i^*}(a_{i^*}^T y_k - b_{i^*})^+ \right] \rangle$$

$$\overset{(13) \,\&\, (9)}{=} -2\gamma \langle \frac{\omega}{\alpha}\left[ y_k - (1-\alpha)x_k \right] + (1-\omega)y_k - \mathcal{V}, \nabla f(y_k) \rangle$$

$$= -2\gamma \langle \frac{\omega(1-\alpha)}{\alpha}(y_k - x_k) + y_k - \mathcal{V}, \nabla f(y_k) \rangle. \tag{85}$$

Using Lemma 3.7 and Lemma 3.9 we can simplify equation (85) as follows,

$$-2\gamma \mathbb{E}_{\mathbb{S}_k}[I_3] = -2\gamma \langle \frac{\omega(1-\alpha)}{\alpha}(y_k - x_k) + y_k - \mathcal{V}, \nabla f(y_k) \rangle$$

$$= 2\gamma \frac{\omega(1-\alpha)}{\alpha}\langle x_k - y_k, \nabla f(y_k) \rangle + 2\gamma \langle \mathcal{V} - y_k, \nabla f(y_k) \rangle$$

$$\overset{\text{Lemma 3.7 \& 3.9}}{\leq} \frac{\gamma\omega(1-\alpha)}{\alpha} d(x_k, P)^2 - \frac{\mu_1\gamma\omega(1-\alpha)}{\alpha} d(y_k, P)^2 - 2\mu_1\gamma \, d(y_k, P)^2. \tag{86}$$

Now, substituting the values of equation (83), (84) & (86) in equation (82) we get the following:

$$\mathbb{E}[d(v_{k+1}, P)^2] = I_1 + \gamma^2 \mathbb{E}_{\mathbb{S}_k}[I_2] - 2\gamma \mathbb{E}_{\mathbb{S}_k}[I_3]$$

$$= \omega \, d(v_k, P)^2 + (1-\omega) \, d(y_k, P)^2 + \frac{\gamma^2}{\eta}\left\{ d(y_k, P)^2 - \mathbb{E}\left[ d(x_{k+1}, P)^2 \right] \right\}$$

$$+ \frac{\gamma\omega(1-\alpha)}{\alpha} d(x_k, P)^2 - \gamma\mu_1\left( 2 + \frac{\omega(1-\alpha)}{\alpha} \right) d(y_k, P)^2.$$

With further simplification, the above identity can be written as follows:

$$\mathbb{E}\left[ d(v_{k+1}, P)^2 + \frac{\gamma^2}{\eta} d(x_{k+1}, P)^2 \right] = \omega \left[ d(v_k, P)^2 + \frac{\gamma(1-\alpha)}{\alpha} d(x_k, P)^2 \right]$$

$$+ d(y_k, P)^2 \left\{ 1 - \omega + \frac{\gamma^2}{\eta} - 2\gamma\mu_1 - \frac{\gamma\omega\mu_1(1-\alpha)}{\alpha} \right\}. \tag{87}$$

Now, let's choose the parameters as in equation (32) along with $0 < \zeta < \frac{4\eta\mu_1}{(1-\mu_1)^2}$. We can easily see that $\frac{\gamma^2}{\eta} = \frac{\gamma(1-\alpha)}{\alpha}$ and $\alpha \in (0, 1)$. Also note that,

$$2\mu_1\gamma = 2\mu_1\sqrt{\eta\zeta\mu_1} > 2\mu_1\sqrt{\frac{\zeta(1-\mu_1)^2\zeta}{4}} = \mu_1\zeta(1-\mu_1), \tag{88}$$

which implies $\omega < 1$. Similarly, whenever $\mu_1 < 1$ we have

$$2\gamma - \zeta - \frac{1}{\mu_1} < 2\sqrt{\eta\mu_1\zeta} - \zeta - 1 \leq 2\sqrt{\zeta} - \zeta - 1 = -(\sqrt{\zeta} - 1)^2 \leq 0, \tag{89}$$

which implies $\omega > 0$. Also, using the parameter choice of (32), we have,

$$1 - \omega + \frac{\gamma^2}{\eta} - 2\gamma\mu_1 - \frac{\gamma\omega\mu_1(1-\alpha)}{\alpha} = 1 - \omega + \zeta\mu_1 - 2\gamma\mu_1 - \omega\zeta\mu_1^2$$

$$= 1 - 2\gamma\mu_1 + \zeta\mu_1 - \omega(1 + \zeta\mu_1^2) = 0. \tag{90}$$

Now, using all of the above relations (equation (32), (90)) in equation (87), we get the following:

$$\mathbb{E}\left[ d(v_{k+1}, P)^2 + \frac{\gamma^2}{\eta} d(x_{k+1}, P)^2 \right] \leq \omega \left[ d(v_k, P)^2 + \underbrace{\frac{\gamma(1-\alpha)}{\alpha}}_{=\frac{\gamma^2}{\eta}} d(x_k, P)^2 \right]$$

$$+ \underbrace{\left[ 1 - \omega + \frac{\gamma^2}{\eta} - 2\gamma\mu_1 - \frac{\gamma\omega\mu_1(1-\alpha)}{\alpha} \right]}_{= 0} d(y_k, P)^2$$

$$= \omega \left[ d(v_k, P)^2 + \frac{\gamma^2}{\eta} d(x_k, P)^2 \right]. \tag{91}$$

Finally, taking expectation again with tower rule and substituting $\frac{\gamma^2}{\eta} = \zeta\mu_1$ we have,

$$\begin{aligned}
\mathbb{E}\left[ d(v_{k+1}, P)^2 + \zeta\mu_1\, d(x_{k+1}, P)^2 \right] &\leq \omega^{k+1}\, \mathbb{E}\left[ d(v_0, P)^2 + \zeta\mu_1\, d(x_0, P)^2 \right] \\
&= (1 + \zeta\mu_1)\, \omega^{k+1}\, d(x_0, P)^2.
\end{aligned}$$

This proves the Theorem. Furthermore, for faster convergence, we need to choose parameters such that, $\omega$ becomes as small as possible. In the proof, we assumed $\mu_1 < 1$ holds which is the most probable scenario. Whenever $\mu_1 = 1$, we must have $\mu_1 = \mu_2 = 1$ and Lemma 6 holds with both equality, i.e., $f(x) = d(x, P)^2$. Therefore if we choose $\alpha, \gamma, \omega$ as $\alpha = \frac{\eta}{\eta+\gamma}$, $\gamma = \sqrt{\zeta\eta}$, $\omega = 1 - \frac{2\gamma}{1+\zeta}$, we can check that condition (89) holds and $0 < \omega < 1$ holds for any $\zeta > 0$.

**Proof of Theorem 3.24** For any natural number $l \geq 1$, using the update formula of $v_{l+1}$, we have

$$\begin{aligned}
v_{l+1} &= \omega v_l + (1 - \omega)y_l - \gamma(a_{i^*}^T y_l - b_{i^*})^+ a_{i^*} \\
&\overset{(13)}{=} \left( 1 - \omega + \frac{\omega}{\alpha} \right) y_l - \frac{\omega(1-\alpha)}{\alpha} x_l - \gamma(a_{i^*}^T y_l - b_{i^*})^+ a_{i^*}. \tag{92}
\end{aligned}$$

Let $\varphi = \omega(1 - \alpha)$. It can be easily checked that $0 \leq \varphi < 1$. Now, considering equation (13), we have

$$\begin{aligned}
y_{l+1} &= \alpha v_{l+1} + (1 - \alpha)x_{l+1} \\
&\overset{(14)\ \&\ (92)}{=} (1 + \omega - \alpha\omega)y_l - \omega(1-\alpha)y_{l-1} + \omega\delta(1-\alpha)(a_{j^*}^T y_{l-1} - b_{j^*})^+ a_{j^*} \\
&\qquad - [\alpha\gamma + (1-\alpha)\delta]\, (a_{i^*}^T y_l - b_{i^*})^+ a_{i^*} \\
&= (1 + \omega - \alpha\omega)y_l - \omega(1-\alpha)y_{l-1} + \omega\delta(1-\alpha)(a_{j^*}^T y_{l-1} - b_{j^*})^+ a_{j^*} \\
&\qquad - \delta(1 + \omega - \alpha\omega)\, (a_{i^*}^T y_l - b_{i^*})^+ a_{i^*} \\
&= (1 + \varphi)y_l - \varphi y_{l-1} + \delta\varphi(a_{j^*}^T y_{l-1} - b_{j^*})^+ a_{j^*} - \delta(1+\varphi)\, (a_{i^*}^T y_l - b_{i^*})^+ a_{i^*}. \tag{93}
\end{aligned}$$

here, the index $i^*$ and $j^*$ are defined based on (6) respectively for the sequences $y_l$ and $y_{l-1}$. Furthermore, with the choice of $x_0 = v_0$, the points $y_0$ and $y_1$ generated by the PASKM method (i.e, algorithm 3 with arbitrary parameter choice) can be calculated as

$$\begin{aligned}
y_1 = \alpha v_1 + (1-\alpha)x_1 &= x_0 - (\alpha\gamma + \delta(1-\alpha))(a_{i^*}^T y_0 - b_{i^*})^+ a_{i^*} \\
&= y_0 - \delta(1+\varphi)(a_{i^*}^T y_0 - b_{i^*})^+ a_{i^*}, \tag{94}
\end{aligned}$$

since $y_0 = x_0 = v_0$. Now, let's define, $\bar{\vartheta}_l = \frac{\varphi}{1-\varphi}[y_l - y_{l-1} + \delta(a_{j^*}^T y_{l-1} - b_{j^*})^+ a_{j^*}]$, $\bar{\Delta}_l = y_l + \bar{\vartheta}_l$ and $\bar{\chi}_l = \|y_l + \bar{\vartheta}_l - \mathcal{P}(\bar{\Delta}_l)\|^2$, then using the update formula (93), we have

$$\begin{aligned}
y_{l+1} + \bar{\vartheta}_{l+1} &= y_{l+1} + \frac{\varphi}{1-\varphi}[y_{l+1} - y_l + \delta(a_{i^*}^T y_l - b_{i^*})^+ a_{i^*}] \\
&= \frac{1}{1-\varphi}y_{l+1} - \frac{\varphi}{1-\varphi}y_l + \frac{\delta\varphi}{1-\varphi}(a_{i^*}^T y_l - b_{i^*})^+ a_{i^*} \\
&\overset{(93)}{=} y_l + \bar{\vartheta}_l - \frac{\delta}{1-\varphi}\left( a_{i^*}^T y_l - b_{i^*} \right)^+ a_{i^*}.
\end{aligned}$$

Using the above relation, we can write

$$\begin{aligned}
\bar{\chi}_{l+1} = \|y_{l+1} + \bar{\vartheta}_{l+1} - \mathcal{P}(\bar{\Delta}_{l+1})\|^2 &\overset{\text{Lemma 3.3}}{\leq} \|y_{l+1} + \bar{\vartheta}_{l+1} - \mathcal{P}(\bar{\Delta}_l)\|^2 \\
&= \left\| y_l + \bar{\vartheta}_l - \frac{\delta}{1-\varphi}\left( a_{i^*}^T y_l - b_{i^*} \right)^+ a_{i^*} - \mathcal{P}(\bar{\Delta}_l) \right\|^2 \\
&= \underbrace{\|y_l + \bar{\vartheta}_l - \mathcal{P}(\bar{\Delta}_l)\|^2}_{=\bar{\chi}_l} + \frac{\delta^2}{(1-\varphi)^2} \underbrace{\|(a_{i^*}^T y_l - b_{i^*})^+ a_{i^*}\|^2}_{I_1} \\
&\quad - \frac{2\delta}{1-\varphi} \underbrace{\left\langle y_l + \bar{\vartheta}_l - \mathcal{P}(\bar{\Delta}_l)\,,\ a_{i^*}(a_{i^*}^T y_l - b_{i^*})^+ \right\rangle}_{I_2}
\end{aligned}$$

39

$$= \bar{\chi}_l + \frac{\delta^2}{(1-\varphi)^2} I_1 - \frac{2\delta}{1-\varphi} I_2. \tag{95}$$

Taking expectation with respect to $\mathbb{S}_l$ we have,

$$\frac{\delta^2}{(1-\varphi)^2} \mathbb{E}_{\mathbb{S}_l}[I_1] \overset{(9)}{=} \frac{2\delta^2}{(1-\varphi)^2} f(y_l). \tag{96}$$

Similarly, we can simplify the third term of (95) as

$$-\frac{2\delta}{1-\varphi} \mathbb{E}_{\mathbb{S}_l}[I_2]$$

$$\overset{(9)}{=} -\frac{2\delta}{1-\varphi} \langle y_l - \mathcal{P}(\bar{\Delta}_l), \nabla f(y_l) \rangle + \frac{2\delta\varphi}{(1-\varphi)^2} \langle y_{l-1} - y_l - \delta\nabla f(y_{l-1}), \nabla f(y_l) \rangle$$

$$= -\frac{2\delta}{1-\varphi} \langle y_l - \mathcal{P}(\bar{\Delta}_l), \nabla f(y_l) \rangle + \frac{2\delta\varphi}{(1-\varphi)^2} \langle y_{l-1} - y_l, \nabla f(y_l) \rangle$$

$$\quad - \frac{\delta^2\varphi}{(1-\varphi)^2} \left[ \|\nabla f(y_l) + \nabla f(y_{l-1})\|^2 - \|\nabla f(y_l)\|^2 - \|\nabla f(y_{l-1})\|^2 \right]$$

$$\overset{\text{Lemma 3.7 \& 3.9}}{\leq} -\frac{4\delta}{1-\varphi} f(y_l) + \frac{2\delta\varphi}{(1-\varphi)^2} [f(y_{l-1}) - f(y_l)] + \frac{2\delta^2\varphi}{(1-\varphi)^2} [f(y_{l-1}) + f(y_l)]$$

$$= \frac{2\delta\varphi(1+\delta)}{(1-\varphi)^2} f(y_{l-1}) - \frac{2\delta\varphi(1+\delta)}{(1-\varphi)^2} f(y_l) + \frac{4\delta(\varphi + \delta\varphi - 1)}{(1-\varphi)^2} f(y_l). \tag{97}$$

Using the expressions of equation (96) and (97) in (95) and simplifying further, we have

$$\mathbb{E}[\bar{\chi}_{l+1}] + \frac{2\delta\varphi(1+\delta)}{(1-\varphi)^2} f(y_l) + \varsigma f(y_l) \leq \mathbb{E}[\bar{\chi}_l] + \frac{2\delta\varphi(1+\delta)}{(1-\varphi)^2} f(y_{l-1}), \tag{98}$$

here,

$$\varsigma = \frac{4\delta(1-\varphi-\delta\varphi)}{(1-\varphi)^2} - \frac{2\delta^2}{(1-\varphi)^2} = \frac{2\delta(2 - 2\varphi - 2\delta\varphi - \delta)}{(1-\varphi)^2} > 0. \tag{99}$$

Now, taking expectation again in (98) and using the tower property, we get,

$$\bar{q}_{l+1} + \varsigma \mathbb{E}[f(y_l)] \leq \bar{q}_l, \quad l = 1, 2, 3..., \tag{100}$$

where, $\bar{q}_l = \mathbb{E}[\bar{\chi}_l] + \frac{2\delta\varphi(1+\delta)}{(1-\varphi)^2} \mathbb{E}[f(y_{l-1})]$. Summing up (100) for $l = 1, 2, ..., k$ we get

$$\sum_{l=1}^{k} \mathbb{E}[f(y_l)] \leq \frac{\bar{q}_1 - \bar{q}_{k+1}}{\varsigma} \leq \frac{\bar{q}_1}{\varsigma}. \tag{101}$$

Now, using Jensen's inequality, we have

$$\mathbb{E}[f(\bar{y}_k)] = \mathbb{E}\left[ f\left( \sum_{l=1}^{k} \frac{y_k}{k} \right) \right] \leq \mathbb{E}\left[ \frac{1}{k} \sum_{l=1}^{k} f(y_l) \right] = \frac{1}{k} \sum_{l=1}^{k} \mathbb{E}[f(y_l)] \overset{(101)}{\leq} \frac{\bar{q}_1}{\varsigma k}.$$

From (94), $y_1 = y_0 - \delta(1+\varphi)(a_{i^*}^T y_0 - b_{i^*})^+ a_{i^*}$ and $\bar{\vartheta}_1 = \frac{-\varphi^2 \delta}{1-\varphi}(a_{i^*}^T y_0 - b_{i^*})^+ a_{i^*}$. Then,

$$\mathbb{E}[\bar{\chi}_1] = \mathbb{E}\left[ \|y_1 + \bar{\vartheta}_1 - \mathcal{P}(\bar{\Delta}_1)\|^2 \right] \overset{\text{Lemma 3.3}}{\leq} \mathbb{E}\left[ \|y_1 + \bar{\vartheta}_1 - \mathcal{P}(y_0)\|^2 \right]$$

$$= \mathbb{E}\left[ \|y_0 - \mathcal{P}(y_0) - \frac{\delta}{1-\varphi}(a_{i^*}^T y_0 - b_{i^*})^+ a_{i^*}\|^2 \right]$$

$$= \|y_0 - \mathcal{P}(y_0)\|^2 + \frac{\delta^2}{(1-\varphi)^2} \mathbb{E}[|(a_{i^*}^T y_0 - b_{i^*})^+|^2]$$

$$\quad - \frac{2\delta}{1-\varphi} \langle y_0 - \mathcal{P}(y_0), \mathbb{E}[(a_{i^*}^T y_0 - b_{i^*})^+ a_{i^*}] \rangle$$

$$\overset{\text{Lemma 3.9}}{\leq} \|y_0 - \mathcal{P}(y_0)\|^2 + \frac{2\delta^2}{(1-\varphi)^2} f(y_0) - \frac{4\delta}{1-\varphi} f(y_0). \tag{102}$$

40

Now, from our construction we get

$$\bar{q}_1 = \mathbb{E}[\bar{\chi}_1] + \frac{2\delta\varphi(1+\delta)}{(1-\varphi)^2}\mathbb{E}[f(y_0)] \le d(y_0, P)^2 + \frac{2\delta(\delta - 2 + 3\varphi + \delta\varphi)}{(1-\varphi)^2}f(y_0).$$

Substituting the values of $\varsigma$ and $q_1$ in the expression of $\mathbb{E}\left[f(\bar{y}_k)\right]$, we have the following:

$$\mathbb{E}\left[f(\bar{y}_k)\right] \le \frac{(1 - \omega + \alpha\omega)^2\, d(y_0, P)^2 + 2\delta(\delta - 2 + 3\omega - 3\alpha\omega + \delta\omega - \delta\alpha\omega)f(y_0)}{2\delta k\,(2 - 2\omega + 2\alpha\omega - 2\delta\omega + 2\delta\alpha\omega - \delta)}.$$

## References

[1] Thomas Strohmer and Roman Vershynin. A randomized kaczmarz algorithm with exponential convergence. *Journal of Fourier Analysis and Applications*, 15(2):262, Apr 2008.

[2] Dennis Leventhal and Adrian S. Lewis. Randomized methods for linear constraints: Convergence rates and conditioning. *Mathematics of Operations Research*, 35(3):641–654, 2010.

[3] Deanna Needell. Randomized kaczmarz solver for noisy linear systems. *BIT Numerical Mathematics*, 50(2):395–403, Jun 2010.

[4] Petros Drineas, Michael W. Mahoney, Shan Muthukrishnan, and Tamás Sarlós. Faster least squares approximation. *Numerische Mathematik*, 117(2):219–249, Feb 2011.

[5] Anastasios Zouzias and Nikolaos M. Freris. Randomized extended kaczmarz for solving least squares. *SIAM Journal on Matrix Analysis and Applications*, 34(2):773–793, 2013.

[6] Yin Tat Lee and Aaron Sidford. Efficient accelerated coordinate descent methods and faster algorithms for solving linear systems. In *Proceedings of the 2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, FOCS '13, pages 147–156, Washington, DC, USA, 2013. IEEE Computer Society.

[7] Anna Ma, Deanna Needell, and Aaditya Ramdas. Convergence properties of the randomized extended gauss seidel and kaczmarz methods. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1590–1604, Jan 2015.

[8] Robert M. Gower and Peter Richtárik. Randomized iterative methods for linear systems. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690, 2015.

[9] Zheng Qu, Peter Richtarik, Martin Takac, and Olivier Fercoq. SDNA: Stochastic Dual Newton Ascent for Empirical Risk Minimization. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pages 1823–1832, New York, USA, 20–22 Jun 2016. PMLR.

[10] Deanna Needell, Nathan Srebro, and Rachel Ward. Stochastic gradient descent, weighted sampling, and the randomized kaczmarz algorithm. *Mathematical Programming*, 155(1):549–573, Jan 2016.

[11] Jesús De Loera, Jamie Haddock, and Deanna Needell. A sampling kaczmarz–motzkin algorithm for linear feasibility. *SIAM Journal on Scientific Computing*, 39(5):S66–S87, 2017.

[12] Meisam Razaviyayn, Mingyi Hong, Navid Reyhanian, and Zhi-Quan Luo. A linearly convergent doubly stochastic gauss–seidel algorithm for solving linear equations and a certain class of over-parameterized optimization problems. *Mathematical Programming*, 176(1):465–496, Jul 2019.

[13] Stefan Kaczmarz. Angenaherte auflsung von systemen linearer gleichungen. *Bulletin International de l'Acadmie Polonaise des Sciences et des Letters*, 35:355–357, 1937.

[14] Richard Gordon, Robert Bender, and Gabor T. Herman. Algebraic reconstruction techniques (art) for three-dimensional electron microscopy and x-ray photography. *Journal of Theoretical Biology*, 29(3):471 – 481, 1970.

[15] Yair Censor. Parallel application of block-iterative methods in medical imaging and radiation therapy. *Mathematical Programming*, 42(1):307–325, Apr 1988.

[16] Gabor T. Herman. *Fundamentals of Computerized Tomography: Image Reconstruction from Projections*. Springer Publishing Company, Incorporated, 2nd edition, 2009.

[17] D. A. Lorenz, S. Wenger, F. Schöpfer, and M. Magnor. A sparse kaczmarz solver and a linearized bregman method for online compressed sensing. In *2014 IEEE International Conference on Image Processing (ICIP)*, pages 1347–1351, Oct 2014.

[18] Joseph M. Elble, Nikolaos V. Sahinidis, and Panagiotis Vouzis. Gpu computing with kaczmarz's and other iterative algorithms for linear systems. *Parallel Computing*, 36(5):215 – 231, 2010. Parallel Matrix Algorithms and Applications.

[19] Fabio Pasqualetti, Ruggero Carli, and Francesco Bullo. Distributed estimation via iterative projections with application to power network monitoring. *Automatica*, 48(5):747 – 758, 2012.

[20] Yair Censor. Row-action methods for huge and sparse systems and their applications. *SIAM Review*, 23(4):444–466, 1981.

[21] Shmuel Agamon. The relaxation method for linear inequalities. *Canadian J. Math*, pages 382–392, 1954.

[22] Theodore S. Motzkin and Issac J. Schoenberg. The relaxation method for linear inequalities. *Canadian J. Math*, pages 393–404, 1954.

[23] Frank Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, pages 65–386, 1958.

[24] Aaditya Ramdas and Javier Peña. Margins, kernels and non-linear smoothed perceptrons. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32, pages 244–252, Bejing, China, 22–24 Jun 2014. PMLR.

[25] Aaditya Ramdas and Javier Peña. Towards a deeper geometric, analytic and algorithmic understanding of margins. *Optimization Methods and Software*, 31(2):377–391, 2016.

[26] Julie Nutini, Behrooz Sepehry, Issam Laradji, Mark Schmidt, Hoyt Koepke, and Alim Virani. Convergence rates for greedy kaczmarz algorithms, and faster randomized kaczmarz rules using the orthogonality graph. In *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*, UAI'16, pages 547–556, Arlington, Virginia, United States, 2016. AUAI Press.

[27] Stefania Petra and Constantin Popa. Single projection kaczmarz extended algorithms. *Numerical Algorithms*, 73(3):791–806, Nov 2016.

[28] Jan Telgen. On relaxation methods for systems of linear inequalities. *European Journal of Operational Research*, 9(2):184 – 189, 1982.

[29] J. F. Maurras, K. Truemper, and M. Akgül. Polynomial algorithms for a class of linear programs. *Mathematical Programming*, 21(1):121–136, Dec 1981.

[30] Sergei Chubanov. A strongly polynomial algorithm for linear systems having a binary solution. *Mathematical Programming*, 134(2):533–570, Sep 2012.

[31] Sergei Chubanov. A polynomial projection algorithm for linear feasibility problems. *Mathematical Programming*, 153(2):687–713, Nov 2015.

[32] Ji Liu and Stephen J. Wright. An accelerated randomized kaczmarz algorithm. *Math. Comput.*, 85(297):153–178, 2016.

[33] Deanna Needell, Ran Zhao, and Anastasios Zouzias. Randomized block kaczmarz method with projection for solving least squares. *Linear Algebra and its Applications*, 484:322 – 343, 2015.

[34] Anna. Ma, Deanna. Needell, and Aaditya. Ramdas. Convergence properties of the randomized extended gauss–seidel and kaczmarz methods. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1590–1604, 2015.

[35] Ahmed. Hefny, Deanna. Needell, and Aaditya. Ramdas. Rows versus columns: Randomized kaczmarz or gauss–seidel for ridge regression. *SIAM Journal on Scientific Computing*, 39(5):S528–S542, 2017.

[36] Robert M. Gower and Peter Richtárik. Linearly convergent randomized iterative methods for computing the pseudoinverse, 2016.

[37] Robert M. Gower and Peter. Richtárik. Randomized quasi-newton updates are linearly convergent matrix inversion algorithms. *SIAM Journal on Matrix Analysis and Applications*, 38(4):1380–1409, 2017.

[38] Peter Richtárik and Martin Takáč. Stochastic reformulations of linear systems: Algorithms and convergence theory, 2017.

[39] Robert Gower, Filip Hanzely, Peter Richtarik, and Sebastian U Stich. Accelerated stochastic matrix inversion: General theory and speeding up bfgs rules for faster second-order optimization. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 1619–1629. Curran Associates, Inc., 2018.

[40] Deanna Needell and Joel A. Tropp. Paved with good intentions: Analysis of a randomized block kaczmarz method. *Linear Algebra and its Applications*, 441:199 – 221, 2014. Special Issue on Sparse Approximate Solution of Linear Systems.

[41] Jonathan Briskman and Deanna Needell. Block kaczmarz method with inequalities. *J. Math. Imaging Vis.*, 52(3):385–396, July 2015.

[42] Deanna Needell and Elizaveta Rebrova. On block gaussian sketching for the kaczmarz method, 2019.

[43] Amitabh Basu, Jesús A. De Loera, and Mark Junod. On chubanov's method for linear programming. *INFORMS Journal on Computing*, 26(2):336–350, 2014.

[44] László A. Végh and Giacomo Zambelli. A polynomial projection-type algorithm for linear programming. *Operations Research Letters*, 42(1):91 – 96, 2014.

[45] Yonina C. Eldar and Deanna Needell. Acceleration of randomized kaczmarz method via the johnson–lindenstrauss lemma. *Numerical Algorithms*, 58(2):163–177, Oct 2011.

[46] A. Agaskar, C. Wang, and Y. M. Lu. Randomized kaczmarz algorithms: Exact mse analysis and optimal sampling probabilities. In *2014 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 389–393, Dec 2014.

[47] Zhong-Zhi Bai and Wen-Ting Wu. On relaxed greedy randomized kaczmarz methods for solving large sparse linear systems. *Applied Mathematics Letters*, 83:21 – 26, 2018.

[48] Zhong-Zhi. Bai and Wen-Ting. Wu. On greedy randomized kaczmarz method for solving large sparse linear systems. *SIAM Journal on Scientific Computing*, 40(1):A592–A606, 2018.

[49] Nikola B Kovachki and Andrew M Stuart. Analysis of momentum methods. *arXiv preprint arXiv:1906.04285*, 2019.

[50] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.

[51] Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.

[52] Yuri Nesterov. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. *Soviet Mathematics Doklady*, Vol. 27:p(372–376), 1983.

[53] Yuri Nesterov. Smooth minimization of non-smooth functions. *Mathematical Programming*, 103(1):127–152, May 2005.

[54] Yuri Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1):125–161, Aug 2013.

[55] Yuri Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Springer Publishing Company, Incorporated, 1 edition, 2014.

[56] Yuri Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.

[57] Nicolas Loizou and Peter Richtárik. Momentum and stochastic momentum for stochastic gradient, newton, proximal point and subspace descent methods, 2017.

[58] Md Sarowar Morshed and Md. Noor-E-Alam. Generalized affine scaling algorithms for linear programming problems. *Computers & Operations Research*, 114:104807, 2020.

[59] Michael Rabbat Nicolas Loizou and Peter Richtárik. Provably accelerated randomized gossip algorithms. *Arxiv*, 2018.

[60] Md Sarowar Morshed, Md Saiful Islam, and Md. Noor-E-Alam. Accelerated sampling kaczmarz motzkin algorithm for the linear feasibility problem. *Journal of Global Optimization*, Oct 2019.

[61] Jamie Haddock and Anna Ma. Greed works: An improved analysis of sampling kaczmarz-motkzin, 2019.

[62] Md Sarowar Morshed and Md. Noor-E-Alam. Heavy ball momentum induced sampling kaczmarz motzkin methods for linear feasibility problems. *arXiv preprint arXiv:200908251*, 2020.

[63] Alan J Hoffman. On approximate solutions of systems of linear inequalities. In *Selected Papers Of Alan J Hoffman: With Commentary*, pages 174–176. World Scientific, 2003.

[64] Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In Paolo Frasconi, Niels Landwehr, Giuseppe Manco, and Jilles Vreeken, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 795–811, Cham, 2016. Springer International Publishing.

[65] L.G. Khachiyan. Polynomial algorithms in linear programming. *USSR Computational Mathematics and Mathematical Physics*, 20(1):53 – 72, 1980.

[66] I-Cheng Yeh and Che-hui Lien. The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. *Expert Syst. Appl.*, 36(2):2473–2480, Mar 2009.

[67] Moshe Lichman. UCI machine learning repository, 2013.

[68] Netlib. The netlib linear programming library.

[69] Giuseppe Calafiore and Laurent El Ghaoui. *Optimization Models*. Control systems and optimization series. Cambridge University Press, October 2014.