



HAL
open science

Preconditioned Gradient Descent for Sketched Mixture Learning

Joseph Gabet, Maxime Ferreira Da Costa

► **To cite this version:**

Joseph Gabet, Maxime Ferreira Da Costa. Preconditioned Gradient Descent for Sketched Mixture Learning. International Symposium on Information Theory, IEEE, Jul 2024, Athens, Greece. hal-04425748v2

HAL Id: hal-04425748

<https://centralesupelec.hal.science/hal-04425748v2>

Submitted on 30 May 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Copyright

Preconditioned Gradient Descent for Sketched Mixture Learning

Joseph Gabet and Maxime Ferreira Da Costa

CentraleSupélec, Université Paris–Saclay, CNRS, Laboratory of Signals and Systems, Gif-sur-Yvette, France

Emails: {joseph.gabet, maxime.ferreira}@centralesupelec.fr

Abstract—Sketching consists of reducing the dimensionality of data samples, generally by summarizing their essential features, such as retaining a small number of empirical moments. The reduced representation is called a sketch. In this paper, a Preconditioned Gradient Descent algorithm (PGD) is proposed to estimate the parameter of mixture models (MM) in arbitrary dimensions by minimizing the non-convex quadratic loss between the sketch and the characteristic function of an MM of varying parameters. Preconditioning is introduced to dynamically adapt the descent direction to the local landscape of the objective function, accelerating convergence, with no computational overhead per iteration compared with vanilla GD. An analysis of the linear convergence rate of PGD is conducted, and numerical simulations showcase the method’s effectiveness, particularly when the weight of the classes is unbalanced or when a substantial number of data samples is available.

I. INTRODUCTION

Clustering is a fundamental task in unsupervised machine learning. It aims to segment data into homogeneous classes, revealing intrinsic structures, patterns, or hidden relationships in data without prior knowledge of class labels or expected outcomes. Clustering is ubiquitous in applied science and engineering, with applications in areas like image segmentation, biological data analysis [1], content recommendations, and anomaly detection. Due to its flexibility in capturing various data shapes, the mixture model (MM) is commonly assumed to model clusters. It is a well-studied statistical prior in statistics and machine learning. When the class distribution is known, specific numerical methods are known to estimate the parameters of an MM from its empirical samples, such as expectation-maximization [2], [3], or hierarchical clustering [4]. However, alternative methods can provide comparative advantages such as a more comprehensive theoretical understanding or a better scaling with the ambient dimension and the dynamic range of the prior distribution [5].

More recently, estimating MM has been approached through *sketching* [6], [7], a scalable method, which retains in the first few empirical moments of the data distribution to reduce the ambient dimensionality. Then, the mixture parameters are estimated by minimizing the non-convex quadratic loss between the sketch and the characteristic function of an MM.

This work is supported in part by ANR PIA funding: ANR-20-IDEEES-0002, AID funding: AID-2023639, and the Orange chair on “Sustainable 6G” held by CentraleSupélec.

A. Contributions and Organization of the Paper

In this study, we consider the sketching framework proposed in [6], and provide novel guarantees on the local geometry of the optimization landscape around the ground truth parameter. Furthermore, we adapt the Preconditioned Gradient Descent (PGD) algorithm [8] into sketching context and prove its linear convergence rate towards the ground truth under sufficiently provided a large number of samples. Additionally, the convergence guarantees are extended to the multidimensional case and to a broad class of convolution kernels.

The rest of the paper is organized as follows. The sketching method and the PGD algorithm are presented in Section II. Section III presents our main results. First, convergence guarantees are provided under infinitely many samples in Theorem 1, then the scaling law of the loss function is discussed under finitely many data samples. Section IV present a proof of Theorem 1. Numerical experiments are conducted in Section V, and a conclusion is drawn in Section VI.

B. Notation and Definitions

Vectors and matrices are denoted by boldface and capital boldface letters, respectively. Elements of vectors $\mathbf{x} \in \mathbb{C}^N$ with odd dimension $N = 2n + 1$ are indexed between $-n$ and n , so that $\mathbf{x} = [x_{-n}, \dots, x_n]^\top$ for convenience. The notation a, b is used to denote the set of all integers from a to b , inclusive. Transpose and Hermitian transpose of a vector or a matrix \mathbf{A} are denoted by \mathbf{A}^\top and \mathbf{A}^H , respectively. We denote by $\mathbb{1}_d$ the all-one vector in \mathbb{R}^d . With a slight abuse of notation, we denote by $|\mathbf{a}|$, $|\mathbf{a}|^p$, the vector with entries equal the modulus, the p th power of the entries of vector \mathbf{a} , respectively. The entrywise multiplication product is denoted $\mathbf{a} \odot \mathbf{a}'$, and the d -dimensional torus is written $\mathbb{T}^d = (\mathbb{R}/\mathbb{Z})^d$.

II. SKETCHED MIXTURE LEARNING

A. Problem Formulation

Given a d -dimension probability distribution with probability density function (PDF) g , we write $\mu(\boldsymbol{\theta}^*)$ a g -mixture of r components parameterized by $\boldsymbol{\theta}^* = [\mathbf{a}^{*T}, \boldsymbol{\tau}^{*T}]^\top$, where $\mathbf{a}^* = [a_1^*, \dots, a_r^*]^\top \in [0, 1]^r$ is the weight vector of the classes and $\boldsymbol{\tau}^* = [\tau_{1,1}^*, \dots, \tau_{1,d}^*, \tau_{2,1}^*, \dots, \tau_{r,d}^*]^\top \in \mathbb{R}^{dr}$ encodes the locations of the centroids. Thus the PDF $\mu(\boldsymbol{\theta})$ is given by

$$\mu(\boldsymbol{\theta}^*)(\boldsymbol{\tau}) = \sum_{j=1}^r a_j^* g(\boldsymbol{\tau} - \boldsymbol{\tau}_j^*), \quad \boldsymbol{\tau} \in \mathbb{R}^d. \quad (1)$$

In this work, the mixture density function g is assumed to be known. Mixture learning is the task of inferring the parameter θ^* from m *i.i.d.* samples $\mathbf{X} = \{x_1, \dots, x_m\}$ of the distribution $\mu(\theta^*)$.

In the paper, rather than directly relying on the empirical probability distribution to recover the mixture parameter θ^* , we propose to minimize the Euclidean loss of n discrete samples of the empirical characteristic function of the observation \mathbf{X} . This sketching technique [6] maps the ambient dimension from $\mathbb{R}^{md} \rightarrow \mathbb{R}^N$, which allows fine control of the dimensionality of the optimization space, and yields substantial computational benefits when $md \gg n$. The empirical multidimensional characteristic function (ECF) $\Phi\{\mathbf{X}\} : \mathbb{R}^d \rightarrow \mathbb{C}$ of \mathbf{X} reads

$$\Phi\{\mathbf{X}\}(\mathbf{u}) = \frac{1}{m} \sum_{i=1}^m e^{-2i\pi\langle \mathbf{u}, \mathbf{x}_i \rangle}, \quad \mathbf{u} \in \mathbb{R}^d, \quad (2)$$

which is an unbiased estimator for the characteristic function of the ground truth MM [9], that we denote $\bar{\Phi}\{\mu(\theta^*)\}$. Furthermore, the characteristic function \hat{g} of the PDF g , the expression for a MM $\mu(\theta)$ is given for all θ by

$$\bar{\Phi}\{\mu(\theta)\}(\mathbf{u}) = \hat{g}(\mathbf{u}) \odot \sum_{j=1}^r a_j e^{-2i\pi\langle \mathbf{u}, \boldsymbol{\tau}_j \rangle}, \quad \mathbf{u} \in \mathbb{R}^d. \quad (3)$$

For the purpose of data processing, one must evaluate the ECF for a finite number of moments according to a sketching scheme. Different sketching schemes have been studied in the literature to discriminate best the classes of the MM, such as uniform sampling and random sampling [6], [10]. For conciseness, we restrict our setup to uniform sampling schemes and let $N = 2n + 1$ be an odd number. We consider the acquisition of N^d samples of the ECF (2) taken over the centrally symmetric uniform sampling set $\Omega = \left[-\frac{n}{N}, \frac{n}{N}\right]^d$.

In the sequel, the number of classes r is assumed to be known. The parameters θ^* are estimated by minimizing the quadratic loss $\mathcal{L}(\cdot)$ between the sketch and the characteristic function of a mixture parameterized by θ . That is

$$\mathcal{L}(\theta) = \frac{1}{2} \left\| [\Phi\{\mathbf{X}\}(\mathbf{u})]_{\mathbf{u} \in \Omega} - [\bar{\Phi}\{\mu(\theta)\}(\mathbf{u})]_{\mathbf{u} \in \Omega} \right\|_2^2. \quad (4)$$

Given the expression (3) of the characteristic function, Equation (4) amounts to finding a r -sparse combination of complex exponentials weighted by the moments of the shape function g that explains best in the quadratic sense the observations $[\Phi\{\mathbf{X}\}(\mathbf{u})]_{\mathbf{u} \in \Omega}$. This problem is also known as *line spectral estimation* or sparse super-resolution [11] in signal processing, and various methods exist to solve it (see *e.g.* [12]–[17], and reference therein). Of particular to the context of this paper, greedy approaches such as Orthogonal Matching Pursuit (OMP) algorithms [18], [19] provide a fast and scalable family of numerical methods for solving sparse inverse problems such as line spectral estimation. However, those methods rely on discretizing the parameter space, yielding a basis mismatch [20] and imperfect reconstruction, even with infinitely many distribution samples ($m = +\infty$).

To circumvent the limitations of greedy algorithms, we propose instead to study the refinement of first-order optimization

Algorithm 1 Preconditioned Gradient Descent

- 1: Initialize θ_0 using OMP; $k \leftarrow 0$.
 - 2: **while** stopping criterion is not met **do**
 - 3: Compute P_k as in (7)
 - 4: $\theta_{k+1} \leftarrow \theta_k - P_k \frac{\partial \mathcal{L}(\theta_k)}{\partial \theta}$.
 - 5: $k \leftarrow k + 1$
 - 6: **Return** θ_k
-

iterates initialized at the output of the OMP algorithm. As the optimization landscape of the loss $\mathcal{L}(\theta)$ given in (4) is non-convex, global convergence guarantees could hardly be derived. Herein, we rely instead on a study of the *local geometry* of the loss around the ground truth θ^* to determine the width of its basin of attraction and the contraction rate to this minimum from an adaptive preconditioning of the descent direction. Further global convergence guarantees demand ensuring an initialization within the basin of attraction.

B. Learning with Preconditioning

We propose refinements and new associated theoretical guarantees to pre-existing work. Given the heterogeneous nature of vector θ —it contains both locations and weights information which scale differently with N —, we use *preconditioning* to adapt the direction of descent of the first-order method to the local landscape of the cost function [8]. This is done by multiplying at step k the gradient by a matrix P_k , which depends on the current estimate θ_k . The optimization procedure is detailed in Algorithm 1. We note that the algorithm doesn't leverage the implicit constraints $\sum_{j=1}^r a_j = 1$ and $a_j \geq 0$, which are unnecessary to establish the local convergence results.

As the computational complexity of Algorithm 1 is driven by that of the matrix-vector multiplication in the fourth line, we restrict our analysis to diagonal preconditioning matrices so that computing the preconditioned descent direction comes with a marginal additional computational cost compared with vanilla gradient descent.

In the sequel, we write $K : \mathbb{R}^d \rightarrow \mathbb{R}$ the discrete autocorrelation of $g(\cdot)$ defined by

$$K(\boldsymbol{\tau}) = \sum_{\mathbf{k} \in \Omega} |\hat{g}(\mathbf{k})|^2 e^{2i\pi\langle \mathbf{k}; \boldsymbol{\tau} \rangle}, \quad (5)$$

which is real-valued since $g(\cdot)$ is real, 1-periodic in every variable and infinitely differentiable.

For convenience, we present the expression of both gradients of the loss (4) with respect to amplitudes and positions, that on direct observations in the frequency domain, $\mathbf{x} = [\Phi\{\mathbf{X}\}(\mathbf{u})]_{\mathbf{u} \in \Omega}$,

$$\frac{d\mathcal{L}(\theta)}{da_j} = \sum_{l=1}^r a_l K(\boldsymbol{\tau}_j - \boldsymbol{\tau}_l) - \sum_{\mathbf{k} \in \Omega} \overline{\hat{g}(\mathbf{k})} x_k e^{2i\pi\langle \boldsymbol{\tau}_j, \mathbf{k} \rangle} \quad (6a)$$

$$\begin{aligned} \frac{d\mathcal{L}(\theta)}{d\tau_{j,d_1}} &= a_j \sum_{l=1}^r a_l \frac{\partial K}{\partial \tau_{.,d_1}}(\boldsymbol{\tau}_j - \boldsymbol{\tau}_l) \\ &\quad - \mathcal{R} \left(\sum_{\mathbf{k} \in \Omega} \overline{\hat{g}(\mathbf{k})} x_k 2i\pi k_{d_1} e^{2i\pi\langle \boldsymbol{\tau}_j, \mathbf{k} \rangle} \right). \end{aligned} \quad (6b)$$

Furthermore, we select the preconditionner,

$$\mathbf{P}_k = \text{diag} \begin{pmatrix} \mathbb{1}_r \\ [-\nabla^2 K(\mathbf{0})]_{1,1}^{-1} |\mathbf{a}_k|^{-2} \\ \vdots \\ [-\nabla^2 K(\mathbf{0})]_{d,d}^{-1} |\mathbf{a}_k|^{-2} \end{pmatrix}. \quad (7)$$

III. ANALYSIS OF ALGORITHM 1

A. Metrics of Analysis

The *dynamic range* of the problem $\kappa > a_{\max}^*/a_{\min}^* > 1$ is defined as the ratio between the largest weight a_{\max} and the smallest weight a_{\min} of the classes of the ground truth MM $\mu(\boldsymbol{\theta}^*)$. Additionally, the *minimal separation* between the classes, denoted $\Delta(\boldsymbol{\tau}^*)$, is defined as the smallest possible ℓ_∞ distance over the torus \mathbb{T}^d between two distinct centroids. That is $\Delta(\boldsymbol{\tau}^*) = \min_{\mathbf{k} \in \mathbb{Z}^d} \min_{j \neq j'} \|\boldsymbol{\tau}_j - \boldsymbol{\tau}_{j'} - \mathbf{k}\|_\infty$. The dynamic range and the minimal separation are known to be quantities of interest to assess the stability of the line spectral estimation problem [21], [22].

Furthermore, we measure the contraction rate in terms of *worst relative distance* of any parameter to the ground truth. This allows a homogeneous control of the parameter estimate across all the classes and transcends certain properties of the problem at hand, such as the number of source points or the dynamic range. Mathematically, we define the diagonal matrix \mathbf{S} as

$$\mathbf{S} = \text{diag} \begin{pmatrix} \mathbf{a}^{*-1} \\ \sqrt{[-\nabla^2 K(\mathbf{0})]_{1,1} \mathbb{1}_r} \\ \vdots \\ \sqrt{[-\nabla^2 K(\mathbf{0})]_{d,d} \mathbb{1}_r} \end{pmatrix},$$

and study the contraction rate of the sequence $\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty$. Specifically, in one-dimensional settings, the error metric reads

$$\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty = \max_j \left\{ \frac{|a_{k,j} - a_j^*|}{|a_j^*|}; \sqrt{-K''(0)} |\tau_{k,j} - \tau_k^*| \right\}.$$

B. Asymptotic Convergence Guarantees

In this section, we study the convergence of Algorithm 1 towards the global minimum of the loss function (4) when the characteristic function of $\mu(\boldsymbol{\theta}^*)$ is observed, that is $m \rightarrow \infty$. To that end, we introduce the notion of admissible density, for which we will establish a guarantee of convergence.

Definition 1 (Admissible density): A density g is said to be *admissible* if there exist functions $h_0, \{h_1^i\}_{1 \leq i \leq d}, \{h_2^{i,j}\}_{1 \leq i,j \leq d}$ of $\mathbb{T}^d \rightarrow \mathbb{R}$, that are positive, summable, and

- decreasing by positive coordinates : $\mathbf{0} \leq \mathbf{x} \leq \mathbf{y} \implies h(\mathbf{x}) \geq h(\mathbf{y})$;
- absolutely dominating the derivatives: $\forall \mathbf{u}, h_0(\mathbf{u}) \geq |K(\mathbf{u})|, h_1^i(\mathbf{u}) \geq \left| \frac{\partial K(\mathbf{u})}{\partial u_i} \right|, h_2^{i,j}(\mathbf{u}) \geq \left| \frac{\partial^2 K(\mathbf{u})}{\partial u_i \partial u_j} \right|,$

and if there exists a constant $C_g > 0$, depending only on g and independent on N such that $\int_{\mathbb{T}^d} h_0 \leq C_g, \int_{\mathbb{T}^d} h_1^i \leq (-[\nabla^2 K(\mathbf{0})]_{i,i})^{\frac{1}{2}} C_g$, and $\int_{\mathbb{T}^d} h_2^{i,j} \leq ([\nabla^2 K(\mathbf{0})]_{i,i} [\nabla^2 K(\mathbf{0})]_{j,j})^{\frac{1}{2}} C_g$.

Note the admissibility conditions in Definition 1 are mild and satisfied by most densities of interest, such as the Gaussian. The next theorem, whose proof is presented in Section IV, guarantees the linear convergence of Algorithm 1 as a function of the dynamic range and the minimal separation, provided a close enough initialization and under an admissibility assumption of the mixture shape g .

Theorem 1 (Asymptotic convergence of PGD): Suppose that g is admissible in the sense of Definition 1. Assume $n \geq 2$, an infinite number of samples ($m = +\infty$), then there exists a constant $C_g > 0$, depending only on g , such that if

$$\gamma := C_g \left(\frac{N \Delta(\boldsymbol{\tau}^*)}{2} \right)^{-d} \frac{a_{\max}^*}{a_{\min}^*} < \frac{1}{2} \quad (8a)$$

$$4 \leq \max_{1 \leq p \leq d} \left\{ \sqrt{[-\nabla^2 K(\mathbf{0})]_{p,p}} \right\} \Delta(\boldsymbol{\tau}^*), \quad (8b)$$

and if the initial point satisfies

$$\|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty \leq 1 - \sqrt{\frac{2}{3}} \quad (8c)$$

then $\{\boldsymbol{\theta}_k\}_{k \in \mathbb{N}}$ converges towards $\boldsymbol{\theta}^*$ and

$$\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \leq \left(\frac{1}{2} + \gamma \right)^k \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty. \quad (9)$$

C. Scaling Law under a Finite Number of Samples

In practice, the characteristic function is estimated from finitely many samples $m < +\infty$, and the ECF is subject to stochastic noise. A benefit of stretching from finite many samples resides in the boundedness of the ECF, which allows the control of its pointwise deviation from the characteristic function as follows.

Theorem 2 (Pointwise concentration of the ECF):

$$\forall \mathbf{u} \in \mathbb{R}^d, \mathbb{P} \left(\left| \Phi\{\mathbf{X}\}(\mathbf{u}) - \mathbb{E} \left(e^{-2i\pi \langle \mathbf{u}, \mathbf{X} \rangle} \right) \right| \geq t \right) \leq 2e^{-\frac{t^2 m}{4}} \quad (10)$$

The proof of Theorem (2) is a direct consequence of Hoeffding's inequality [23] applied to the real and imaginary part of the difference in (10). The deviation between the ECF and the ground truth characteristic function can be uniformly controlled via the union bound, yielding for any $t \geq 0$

$$\mathbb{P} \left(\left\| [\Phi\{\mathbf{X}\}(\mathbf{u})]_{\mathbf{u} \in \Omega} - [\overline{\Phi}\{\mu(\boldsymbol{\theta})\}(\mathbf{u})]_{\mathbf{u} \in \Omega} \right\|_\infty \geq t \right) \leq 2N^d e^{-\frac{t^2 m}{4}}. \quad (11)$$

Equation (10) suggests $m = \Omega(d \log(N))$ samples are for the stochastic noise to vanish in the asymptotic $N \rightarrow \infty$.

IV. PROOF OF THEOREM 1

We start the proof with the following lemma, which relates the contraction of the iterates sequence with the conditioning of the Hessian over the segment between $\boldsymbol{\theta}^*$ and $\boldsymbol{\theta}_k$.

Lemma 1 (Contraction of the iterates):

Let $\mathcal{S}_k = \{\boldsymbol{\theta}^* + u(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*) \mid u \in [0, 1]\}$, we have

$$\|\mathbf{S}(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}^*)\|_\infty \leq \max_{\boldsymbol{\theta} \in \mathcal{S}_k} \left\{ \|\mathbf{I} - \mathbf{S} \mathbf{P}_k \mathbf{H}(\boldsymbol{\theta}) \mathbf{S}^{-1}\|_\infty \right\} \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty. \quad (12)$$

The rest of the proof focuses on bounding the quantity $\|\mathbf{I} - \mathbf{S}\mathbf{P}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty$.

A. Hessian Decomposition

Starting from the expression (6) of the gradient of the loss, the Hessian matrix $\mathbf{H}(\boldsymbol{\theta}) = \nabla^2\mathcal{L}(\boldsymbol{\theta})$ is given by (c.f. [24])

$$\mathbf{H}(\boldsymbol{\theta}) = \mathbf{M}(\mathbf{a})^H\mathbf{D}(\boldsymbol{\tau})\mathbf{M}(\mathbf{a}) + \mathbf{E}(\boldsymbol{\theta}), \quad (13)$$

where the matrices $\mathbf{M}(\mathbf{a})$ and $\mathbf{D}(\boldsymbol{\tau})$ are given by

$$\mathbf{M}(\mathbf{a}) = \text{diag} \begin{pmatrix} \mathbb{1}_r \\ \sqrt{-[\nabla^2 K(0)]_{1,1}\mathbf{a}} \\ \vdots \\ \sqrt{-[\nabla^2 K(0)]_{d,d}\mathbf{a}} \end{pmatrix} \quad (14)$$

$$\mathbf{D}(\boldsymbol{\tau}) = \begin{pmatrix} \mathbf{D}_0(\boldsymbol{\tau}) & \mathbf{D}_1^1(\boldsymbol{\tau}) & \dots & \mathbf{D}_1^d(\boldsymbol{\tau}) \\ \mathbf{D}_1^1(\boldsymbol{\tau})^H & \mathbf{D}_2^{1,1}(\boldsymbol{\tau}) & \dots & \mathbf{D}_2^{1,d}(\boldsymbol{\tau}) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{D}_1^d(\boldsymbol{\tau})^H & \mathbf{D}_2^{1,d}(\boldsymbol{\tau})^H & \dots & \mathbf{D}_2^{d,d}(\boldsymbol{\tau}) \end{pmatrix}, \quad (15)$$

and where the generic term of each block $\mathbf{D}_p(\boldsymbol{\tau})$ is given by

$$[\mathbf{D}_0(\boldsymbol{\tau})]_{i,j} = K(\boldsymbol{\tau}_i - \boldsymbol{\tau}_j) \quad (16a)$$

$$[\mathbf{D}_1^d(\boldsymbol{\tau})]_{i,j} = (-[\nabla^2 K(\mathbf{0})]_{d_1})^{-\frac{1}{2}} \frac{\partial K(\boldsymbol{\tau}_i - \boldsymbol{\tau}_j)}{\partial w_{d_1}} \quad (16b)$$

$$[\mathbf{D}_2^{d_1, d_2}(\boldsymbol{\tau})]_{i,j} = (\nabla^2 K(\mathbf{0})_{d_1} \nabla^2 K(\mathbf{0})_{d_2})^{-\frac{1}{2}} \frac{\partial^2 K(\boldsymbol{\tau}_i - \boldsymbol{\tau}_j)}{\partial w_{d_1} \partial w_{d_2}}. \quad (16c)$$

The term $\mathbf{E}(\boldsymbol{\theta})$ can be interpreted as a perturbation term that vanishes at $\boldsymbol{\theta} = \boldsymbol{\theta}^*$. By the triangle inequality, we have $\|\mathbf{S}\mathbf{P}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \leq \|\mathbf{S}\mathbf{P}_k\mathbf{M}(\mathbf{a})^H\mathbf{D}(\boldsymbol{\tau})\mathbf{M}(\mathbf{a})\mathbf{S}^{-1} - \mathbf{I}\|_\infty + \|\mathbf{S}\mathbf{P}_k\mathbf{E}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty$. We provide the following lemma on $\|\mathbf{S}\mathbf{P}_k\mathbf{E}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty$, and skip its proof conciseness, and focus on the sequel by bounding the remaining term.

Lemma 2: If g is admissible in the sense of Definition 1, and if $\|\mathbf{S}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|_\infty < 1$, then there exists a constant $C'_g > 0$, depending only on g such that

$$\begin{aligned} & \|\mathbf{S}\mathbf{P}_k\mathbf{E}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty \\ & \leq C'_g (N\Delta(\boldsymbol{\tau}))^{-d} \frac{a_{\max}^* \|\mathbf{S}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|_\infty}{a_{\min}^* (1 - \|\mathbf{S}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|_\infty)^2}. \end{aligned} \quad (17)$$

B. Bound on $\|\mathbf{S}\mathbf{P}_k\mathbf{M}(\mathbf{a})^H\mathbf{D}(\boldsymbol{\tau})\mathbf{M}(\mathbf{a})\mathbf{S}^{-1} - \mathbf{I}\|_\infty$

For all j , we have the identities [8]

$$\frac{|a_j^*|}{|a_{k,j}|} \leq \frac{|a_j^*|}{|a_j^*| - |a_{k,j} - a_j^*|} \leq \frac{1}{1 - \|\mathbf{S}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|_\infty} \quad (18a)$$

$$\max\{|a_j^*|, |a_j|\} \leq 1 + \|\mathbf{S}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|_\infty \quad (18b)$$

From a direct calculation, exploiting the diagonal structure of the matrices \mathbf{S} and $\mathbf{M}(\mathbf{a})$, and Equations (18), we have

$$\begin{aligned} & \|\mathbf{S}\mathbf{P}_k\mathbf{M}(\mathbf{a})^H\mathbf{D}(\boldsymbol{\tau})\mathbf{M}(\mathbf{a})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\ & \leq \max_j \left\{ \frac{1}{a_j^*}; \frac{|a_j|}{|a_{k,j}|^2} \right\} \|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty \max_j \{|a_j^*|; |a_j|\} \end{aligned}$$

$$\begin{aligned} & + \max_j \left\{ \left| 1 - \frac{a_j^2}{|a_{k,j}|^2} \right| \right\} \\ & \leq \frac{a_{\max}^*}{a_{\min}^*} \frac{1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} \|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty \\ & \quad + \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} - 1. \end{aligned} \quad (19)$$

We now aim at bounding $\|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty$. To proceed, we exploit the block decomposition (15) and the Hermitian structure of sub-blocks (16) to obtain

$$\begin{aligned} \|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty & \leq \max_i \left\{ \|\mathbf{D}_0(\boldsymbol{\tau}) - \mathbf{I}\|_\infty + \sum_{j=1}^d \|\mathbf{D}_1^j(\boldsymbol{\tau})\|_\infty, \right. \\ & \left. \|\mathbf{D}_1^i(\boldsymbol{\tau})\|_\infty + \|\mathbf{D}_2^{i,i}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty + \sum_{\substack{1 \leq j \leq d \\ i \neq j}} \|\mathbf{D}_2^{i,j}(\boldsymbol{\tau})\|_\infty \right\}. \end{aligned} \quad (20)$$

A key element to control the right-hand side of (20) via the generic terms (16) is to bound the summation of the autocorrelation function K and its derivatives at the points $\{\boldsymbol{\tau}_i - \boldsymbol{\tau}_j\}$, which the following lemma proposes.

Lemma 3 (Summation bounds on the autocorrelation):

If g is admissible in the sense of Definition 1 then

$$\begin{aligned} & \max_{i,j} \left\{ \|\mathbf{D}_0(\boldsymbol{\tau}) - \mathbf{I}\|_\infty, \|\mathbf{D}_1^i(\boldsymbol{\tau})\|_\infty, \right. \\ & \left. \|\mathbf{D}_2^{i,j}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty, \|\mathbf{D}_2^{j,i}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty \right\} \leq (N\Delta(\boldsymbol{\tau}))^{-d} C_g. \end{aligned} \quad (21)$$

Equation (20) and Lemma 3 immediately imply $\|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty \leq (d+1)(N\Delta(\boldsymbol{\tau}))^{-d} C_g$.

C. End of the proof

First of all, Lemma 2, Equation (19) and Lemma 3 yield

$$\begin{aligned} \|\mathbf{S}\mathbf{P}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty & \leq \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} - 1 \\ & \quad + C_g'' (N\Delta(\boldsymbol{\tau}))^{-d} \frac{a_{\max}^*}{a_{\min}^*} (1 + \|\mathbf{S}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|_\infty)^2. \end{aligned} \quad (22)$$

We conclude the proof of Theorem 1 by induction. Assume the conditions (8b) and (8c) hold, which is true of $k=0$. This ensures $\|\boldsymbol{\tau}_0 - \boldsymbol{\tau}^*\| \leq \frac{\Delta(\boldsymbol{\tau}^*)}{4}$, thus $\Delta(\boldsymbol{\tau}) \geq \frac{1}{2}\Delta(\boldsymbol{\tau}^*)$. Thus, the inequality (22) becomes

$$\begin{aligned} & \max_{\boldsymbol{\theta} \in \mathcal{S}_k} \|\mathbf{S}\mathbf{P}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\ & \leq \frac{1}{2} + C_g \left(\frac{N\Delta(\boldsymbol{\tau}^*)}{2} \right)^{-d} \frac{a_{\max}^*}{a_{\min}^*} \end{aligned} \quad (23)$$

$$\leq \frac{1}{2} + \gamma < 1, \quad (24)$$

with $C_g = C_g'' \left(2 - \sqrt{\frac{2}{3}}\right)^2$. We conclude with Lemma 1 that $\|\mathbf{S}(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}^*)\|_\infty \leq \left(\frac{1}{2} + \gamma\right) \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty$ which proves the desired statement. \blacksquare

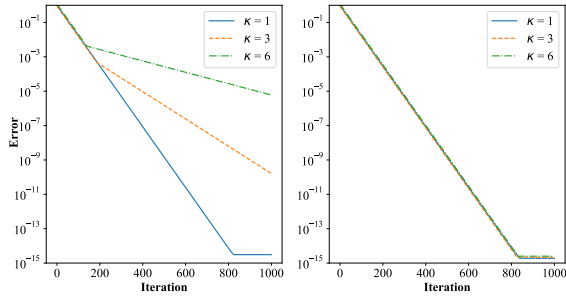


Figure 1. Convergence rates of the iterate sequence of preconditioned GD with fixed step size (left) and adaptive step size (right) towards the ground truth for a 1D Gaussian mixture with variance 1 and centers -2 and 2. The amplitudes are generated to satisfy different values of κ .

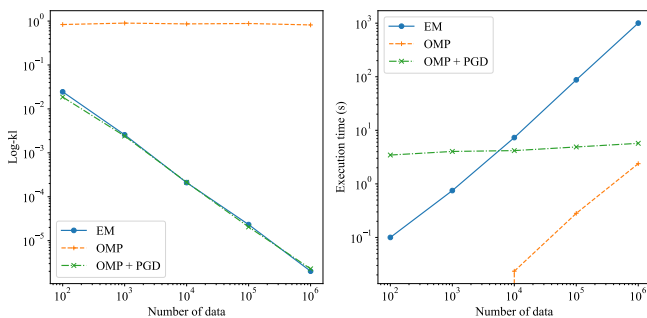


Figure 2. Kullback-Leibler divergence (left) and runtime in seconds (right) per number of samples drawn, displayed in a log-log scale for EM, OMP, and OMP+PGD methods applied to a 2-component 1D Gaussian mixture (with variances of 1, centroids at 0 and 3, and proportions of 0.7 and 0.3, respectively). The results are averaged over 10 random trials.

V. EXPERIMENTAL RESULTS

A. Linear convergence rate

This section conducts a numerical study of the convergence rate of Algorithm 1 compared to a vanilla gradient descent algorithm, where the preconditioning matrix \mathbf{P} doesn't vary with the iterate k . A one-dimensional homoscedastic Gaussian mixture with sufficient separation is considered, while the dynamic range $\kappa = \frac{a_{\max}^*}{a_{\min}^*}$ varies. Figure 1 pictures the error sequence of the iterates $\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_{\infty}$ as k varies. Figure 1 confirms linear convergence of Algorithm 1. Furthermore, we note that adaptive preconditioning remains unaffected by the dynamic scale, in contrast to the fixed step where performance degrades under similar conditions, which corroborates with Theorem 1 for a large enough minimal separation $N\Delta(\tau)$, and reinforces the theoretical underpinnings of our approach.

B. Stochastic noise recovery

We evaluate the performance of our proposed approach, which combines Orthogonal Matching Pursuit (OMP) with Preconditioned Gradient Descent (PGD), denoted as OMP+PGD, against two other clustering methods: OMP alone, and Expectation-Maximization (EM) [3].

We select a Gaussian convolution kernel. The left-hand side graph of Figure 2, shows the Kullback-Leibler divergence

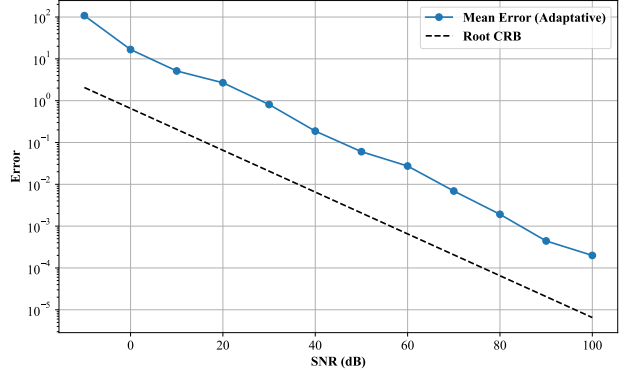


Figure 3. Error $\|\mathbf{S}(\boldsymbol{\theta}_{\infty} - \boldsymbol{\theta}^*)\|_{\infty}$ of Algorithm 1 at convergence as a function of the SNR. g is the Gaussian kernel with $r = 2$ classes. The dynamic range is set to $\kappa = 3$, and the results are averaged over 100 random trials.

(KL) between the reconstructing mixture and the ground truth as a function of the number of samples m for all the considered algorithms. OMP + PGD appears to achieve similar performance as EM and converges as m increases, while OMP is not able to converge to the ground truth due to the mismatch introduced by the discretization. The right-side graph presents the algorithms' execution time. As expected, the execution time increases polynomially with m for all three algorithms. However, sketching in OMP+PGD allows a reduction in the ambient dimension before optimization, yielding better scalability.

C. Signal noise recovery

Finally, we assess the robustness of Algorithm 1 to the stochastic noise caused by the finite number of samples m . Let $\mathbf{w} = [\Phi\{\mathbf{X}\}(\mathbf{u})]_{\mathbf{u} \in \Omega} - [\Phi\{\mu(\boldsymbol{\theta}^*)\}(\mathbf{u})]_{\mathbf{u} \in \Omega}$ be the stochastic noise, which is controlled with high-probability in ℓ_{∞} -norm in Equation (11). We define the signal-to-noise ratio as $\text{SNR} = \|\Phi\{\mu(\boldsymbol{\theta}^*)\}\|_2^2 / \|\mathbf{w}\|_2^2$. Figure 3 shows the error estimate on the parameters of OMP+PGD to estimate a Gaussian mixture with two classes as a function of the SNR. The results are benchmarked against the statistical Cramér-Rao lower bound [25]. Although the empirical error does not achieve the CRB, the statistical error is of the same decay rate. Thus, the proposed method is robust to the stochastic noise induced by a finite number empirical samples.

VI. CONCLUSION

We have presented a preconditioned gradient descent algorithm to identify the components of a mixture model from the empirical moments of the observation through a sketching technique. Local convergence guarantees are provided in arbitrary dimensions and for a very broad class of kernels in Theorem 1, under a good enough initial point, and in the asymptotic number of distribution samples. In particular, we prove the linear contraction rate of the iterate sequence. The analysis amounts to line spectral estimation, bridging classical learning and signal processing.

Future research directions include a study of the convergence in a stochastic framework, where a limited number of samples is assumed.

REFERENCES

- [1] Y. Zhao and G. Karypis, “Data clustering in life sciences”, *Molecular biotechnology*, vol. 31, pp. 55–80, 2005.
- [2] T. K. Moon, “The expectation-maximization algorithm”, *IEEE Signal Processing Magazine*, vol. 13, no. 6, pp. 47–60, 1996.
- [3] F. Dellaert, “The expectation maximization algorithm”, *College of Computing, Georgia Institute of Technology*, 2002.
- [4] J. Goldberger and S. Roweis, “Hierarchical clustering of a mixture model”, *Advances in neural information processing systems*, vol. 17, 2004.
- [5] G. J. McLachlan and K. E. Basford, *Mixture models: Inference and applications to clustering*. M. Dekker New York, 1988, vol. 38.
- [6] R. G. Nicolas Keriven Anthony Bourrier and P. Pérez, “Sketching for large-scale learning of mixture models”, *Information and Inference*, vol. 7, no. 3, pp. 447–508, 2018.
- [7] R. Gribonval, G. Blanchard, N. Keriven, and Y. Traonmilin, “Statistical learning guarantees for compressive clustering and compressive mixture modeling”, *Mathematical Statistics and Learning*, vol. 3, no. 2, pp. 165–257, 2021.
- [8] M. Ferreira Da Costa and Y. Chi, “Local geometry of nonconvex spike deconvolution from low-pass measurements”, *IEEE Journal on Selected Areas in Information Theory*, vol. 4, pp. 1–15, 2023.
- [9] A. Feuerwerker and R. A. Mureika, “The empirical characteristic function and its applications”, *Annals of Statistics*, vol. 5, no. 1, pp. 88–97, 1977.
- [10] R. Gribonval, A. Chatalic, N. Keriven, V. Schellekens, L. Jacques, and P. Schniter, “Sketching data sets for large-scale learning: Keeping only what you need”, *IEEE Signal Processing Magazine*, vol. 38, no. 5, pp. 12–36, 2021.
- [11] G. Tang, B. N. Bhaskar, and B. Recht, “Near minimax line spectral estimation”, *IEEE Transactions on Information Theory*, vol. 61, no. 1, pp. 499–512, 2014.
- [12] R. Prony, “Essai experimental”, *J. de l’Ecole Polytechnique*, vol. 2, p. 929, 1795.
- [13] R. Schmidt, “Multiple emitter location and signal parameter estimation”, *IEEE transactions on antennas and propagation*, vol. 34, no. 3, pp. 276–280, 1986.
- [14] R. Roy and T. Kailath, “Esprit-estimation of signal parameters via rotational invariance techniques”, *IEEE Transactions on acoustics, speech, and signal processing*, vol. 37, no. 7, pp. 984–995, 1989.
- [15] E. J. Candès and C. Fernandez-Granda, “Towards a mathematical theory of super-resolution”, *Communications on pure and applied Mathematics*, vol. 67, no. 6, pp. 906–956, 2014.
- [16] Y. Chi and M. Ferreira Da Costa, “Harnessing sparsity over the continuum: Atomic norm minimization for superresolution”, *IEEE Signal Processing Magazine*, vol. 37, no. 2, pp. 39–57, 2020.
- [17] P.-J. Bénard, Y. Traonmilin, J.-F. Aujol, and E. Soubies, “Estimation of off-the-grid sparse spikes with overparametrized projected gradient descent: Theory and application”, *hal-04220523*, 2023.
- [18] J. A. Tropp and A. C. Gilbert, “Signal recovery from random measurements via orthogonal matching pursuit”, *IEEE Transactions on information theory*, vol. 53, no. 12, pp. 4655–4666, 2007.
- [19] C. Soussen, R. Gribonval, J. Idier, and C. Herzet, “Joint k-step analysis of orthogonal matching pursuit and orthogonal least squares”, *IEEE Transactions on Information Theory*, vol. 59, no. 5, pp. 3158–3174, 2013.
- [20] Y. Chi, L. L. Scharf, A. Pezeshki, and A. R. Calderbank, “Sensitivity to basis mismatch in compressed sensing”, *IEEE Transactions on Signal Processing*, vol. 59, no. 5, pp. 2182–2195, 2011.
- [21] A. Moitra, “Super-resolution, extremal functions and the condition number of Vandermonde matrices”, in *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, 2015, pp. 821–830.
- [22] M. Ferreira Da Costa, “The condition number of weighted non-harmonic Fourier matrices with applications to super-resolution”, *hal-04261330*, Oct. 2023, preprint.
- [23] S. Boucheron, G. Lugosi, and P. Massart, *Concentration Inequalities: A Nonasymptotic Theory of Independence* (Oxford University Press), English. Oxford, UK: Oxford University Press, 2016, p. 496.
- [24] Y. Traonmilin and J.-F. Aujol, “The basins of attraction of the global minimizers of the non-convex sparse spike estimation problem”, *Inverse Problems*, vol. 36, no. 4, p. 045 003, 2020.
- [25] L. L. Scharf and L. McWhorter, “Geometry of the Cramér-Rao bound”, *Signal Processing*, vol. 31, no. 3, pp. 301–311, 1993.