

Tuning-Free Structured Sparse Recovery of Multiple Measurement Vectors using Implicit Regularization

Lakshmi Jayalal

Indian Institute of Technology, Madras

Sheetal Kalyani

Indian Institute of Technology, Madras

Abstract

Recovering jointly sparse signals in the multiple measurement vectors (MMV) setting is a fundamental problem in machine learning, but traditional methods like MMV orthogonal matching pursuit (M-OMP) and MMV-FOCal Underdetermined System Solver (M-FOCUSS) often require careful parameter tuning or prior knowledge of the sparsity of the signal and/or noise variance. We introduce a novel tuning-free framework that leverages Implicit Regularization (IR) from overparameterization to overcome this limitation. Our approach reparameterizes the estimation matrix into factors that decouple the shared row-support from individual vector entries. We show that the optimization dynamics inherently promote the desired row-sparse structure by applying gradient descent to a standard least-squares objective on these factors. We prove that with a sufficiently small and balanced initialization, the optimization dynamics exhibit a “momentum-like” effect, causing the norms of rows in the true support to grow significantly faster than others. This formally guarantees that the solution trajectory converges towards an idealized row-sparse solution. Additionally, empirical results demonstrate that our approach achieves performance comparable to established methods without requiring any prior information or tuning.

1 Introduction

Recovering structured signals from incomplete linear measurements is a central theme in modern signal processing, compressed sensing, and machine learning. A particularly important setting is the multiple measurement vectors (MMV) problem corresponding to sig-

nals sharing a common underlying structure, particularly a shared sparse support. This joint sparsity structure, prevalent in applications from neuroimaging (Chen et al., 2022) to multiple-input multiple-output (MIMO) channel estimation (Yang and Li, 2024) allows for an improved recovery over processing single vectors individually. Classical algorithmic approaches for MMV including MMV subspace pursuit (M-SP) (Liu et al., 2019), MMV orthogonal matching pursuit (M-OMP) (Zhang et al., 2022; Chen et al., 2020) and MMV-FOCal Underdetermined System Solver (M-FOCUSS) introduced by Cotter et al. (2005), often depend on the optimal tuning of regularization parameters or require prior knowledge of the sparsity level. These limitations motivate the search for tuning-free methods for achieving row sparsity.

An alternative paradigm based on Implicit Regularization (IR) through overparameterization has recently gained traction (Arora et al., 2018, 2019; Li et al., 2018). Specifically, gradient descent is used on an overparameterized or factorized representation of the target signal rather than adding explicit penalty terms (like ℓ_p norms) to the loss function to encourage sparsity. The dynamics of the optimization process, often combined with specific initialization strategies or early stopping, implicitly guide the solution towards desirable structures like sparsity or low rank. This phenomenon has been analyzed in context of single measurement vector (SMV) sparse recovery by Vaskevicius et al. (2019); Li et al. (2021); Chou et al. (2023).

Motivated by the success of implicit regularization in the single measurement vector (SMV) domain, this paper explores its application to the row-sparse MMV recovery task. We consider the standard least-squares loss function $\mathcal{L}(\mathbf{X}) = \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2$ and optimize it using gradient descent. We propose a novel overparameterization that leverages the Hadamard product (\odot) and Hadamard power (\odot°), which represent element-wise multiplication and product, respectively. This reparameterization is $\mathbf{X} = (\mathbf{g}^{\odot^2} \mathbf{1}_L) \odot \mathbf{V}$, where $\mathbf{1}_L = [1, 1, \dots, 1]_{1 \times L}$ is a row vector of ones that replicates the squared elements of $\mathbf{g} \in \mathbb{R}^N$ across the

columns of $\mathbf{V} \in \mathbb{R}^{N \times L}$. The core idea is that the gradient dynamics acting on the vector \mathbf{g} , which captures the shared row support, will implicitly drive most of its elements towards zero. Due to the multiplicative structure of our parameterization, the corresponding rows of \mathbf{V} are simultaneously driven to zero, thereby inducing row sparsity in the reconstructed matrix, \mathbf{X} . This parametrization is designed specifically to handle multiple vectors sharing a common structure, distinguishing it from other IR methods that focus on SMV by factorizing a single vector. Furthermore, it contrasts with matrix and tensor factorization methods that induce a low-rank bias rather than the desired row sparsity.

Our main contributions are: (1) we introduce a novel Hadamard parameterization for the MMV problem and derive its gradient descent update rules; (2) we provide a theoretical analysis of the associated gradient flow, establishing formal guarantees that the optimization trajectory converges towards an ideal row-sparse solution; and (3) we further demonstrate the practical effectiveness of our approach through simulations, showing it achieves comparable performance, without any tuning, against established benchmarks like M-OMP, M-SP, and M-FOCUSS.

This research thus aims to bridge the gap between the evolving theory of implicit regularization and its application in the domain of MMV recovery, offering a novel, regularization-free algorithmic framework.

2 Related Works

This work builds upon research in several related areas, primarily MMV sparse recovery algorithms, the increasingly studied field of IR in machine learning and the techniques used to analyze IR dynamics.

Sparse MMV Recovery Algorithms: The challenge of recovering jointly sparse signals from MMV is addressed by diverse algorithms, many adapted from SMV methods. Approaches range from greedy pursuit algorithms like M-SP, and M-OMP to Bayesian methods like M-FOCUSS. These traditional approaches often rely on explicit sparsity-promoting penalties (like l_p norms) or greedy selection rules, sometimes requiring knowledge of the sparsity level or structure. Our work offers an alternative by achieving row sparsity implicitly through optimization dynamics.

Implicit Regularization via Overparameterization: A significant body of recent work focuses on the phenomenon of Implicit Regularization (IR), where the dynamics of gradient descent on an overparameterized model find structured solutions without explicit regularizers (Arora et al., 2018, 2019). This

effect has been shown to implicitly favor low-rank solutions in matrix sensing (Soltanolkotabi et al., 2023) and tensor sensing (Razin et al., 2021; Hariz et al., 2024). Crucially, implicit regularization has also been investigated for sparse recovery in the SMV setting (Li et al., 2021; Zhao et al., 2022; Vaskevicius et al., 2019; Wu et al., 2020). This concept critically extends to achieving structured sparsity, such as group sparsity, through neural reparameterizations, where gradient descent implicitly promotes the desired structure as demonstrated by Li et al. (2023). The scale of initialization is a pivotal factor in these processes, often determining the nature of the learned solution. However, the factorizations used in prior work are suboptimal for capturing the joint sparsity structure inherent in MMV applications. Our work addresses this gap by introducing a parameterization specifically designed for this context. In this framework, the dynamics of gradient descent on our proposed factors implicitly promote the row-sparse structure critical for joint recovery.

Analysis Techniques for Implicit Regularization: Analyzing the dynamics and convergence of gradient descent in overparameterized models requires specialized techniques. While methods based on Singular Value Decomposition (SVD) (e.g., Arora et al. (2019); Soltanolkotabi et al. (2023); Li et al. (2018)) and direct element-wise analysis (Li et al., 2021; Zhao et al., 2022; Vaskevicius et al., 2019; Li et al., 2023) are effective for some sparse recovery problems, they become intractable for the complex, non-linear parametrization such as the one we propose. We, therefore, adopt a dynamical systems perspective using gradient flow inspired by the works Razin et al. (2021) and Hariz et al. (2024). This approach characterizes the evolution of the norms of different signal components over time, which reveals a crucial “momentum-like” effect, where components with larger norms grow significantly faster than smaller ones. This effect is the key mechanism driving the implicit row-sparse regularization in our framework.

3 Preliminaries

We define μ -coherence to characterize the sensing matrix and β -smoothness to describe the properties of the loss function. We also introduce the unbalancedness and row-unbalancedness constants, which are crucial for analyzing the dynamics of our parameterization.

Definition 3.1. *μ -coherence is a measure of the correlation between the columns of a sensing matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ in compressed sensing. It is defined as follows:*

$$\mu(\mathbf{A}) = \max_{i \neq j} |\langle \mathbf{A}_{:,i}, \mathbf{A}_{:,j} \rangle|$$

where $\mathbf{A}_{:,i}$ is the i -th column of sensing matrix \mathbf{A} and

$\langle \cdot, \cdot \rangle$ denotes the inner product.

Throughout this work, we assume that the sensing matrix \mathbf{A} is μ -coherent with ℓ_2 -normalized columns.

Definition 3.2. β -smoothness. We say that a continuously differentiable function $f(\mathbf{g}, \mathbf{V})$ is β -smooth if its gradient $\nabla f(\mathbf{g}, \mathbf{V})$ satisfies

$$\begin{aligned} \|\nabla f(\mathbf{g}, \mathbf{V}) - \nabla f(\tilde{\mathbf{g}}, \tilde{\mathbf{V}})\| &= \left\| \begin{pmatrix} \nabla_{\mathbf{g}} f(\mathbf{g}, \mathbf{V}) - \nabla_{\mathbf{g}} f(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \\ \nabla_{\mathbf{V}} f(\mathbf{g}, \mathbf{V}) - \nabla_{\mathbf{V}} f(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \end{pmatrix} \right\| \\ &\leq \beta \|\mathbf{g}, \mathbf{V} - (\tilde{\mathbf{g}}, \tilde{\mathbf{V}})\|, \end{aligned}$$

where $\|\mathbf{g}, \mathbf{V} - (\tilde{\mathbf{g}}, \tilde{\mathbf{V}})\|^2 = L\|\mathbf{g} - \tilde{\mathbf{g}}\|_2^2 + \|\mathbf{V} - \tilde{\mathbf{V}}\|_F^2$,

$$\begin{aligned} &\left\| \begin{pmatrix} \nabla_{\mathbf{g}} f(\mathbf{g}, \mathbf{V}) - \nabla_{\mathbf{g}} f(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \\ \nabla_{\mathbf{V}} f(\mathbf{g}, \mathbf{V}) - \nabla_{\mathbf{V}} f(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \end{pmatrix} \right\|^2 \\ &= \left\| \nabla_{\mathbf{g}} f(\mathbf{g}, \mathbf{V}) - \nabla_{\mathbf{g}} f(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \right\|_2^2 \\ &\quad + \left\| \nabla_{\mathbf{V}} f(\mathbf{g}, \mathbf{V}) - \nabla_{\mathbf{V}} f(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \right\|_F^2. \end{aligned}$$

Definition 3.3. The constant $\epsilon \geq 0$ be defined as the unbalancedness constant (introduced by Razin et al. (2021)), where $\epsilon = \left| \frac{1}{2}\|\mathbf{g}(t)\|_2^2 - \|\mathbf{V}(t)\|_F^2 \right|$ for any time $t \geq 0$.

Definition 3.4. The constant $\epsilon_r \geq 0$ be defined as the row-unbalancedness constant, where $\epsilon_r = \max_{i \in [N]} \left| \frac{1}{2}g_i^2(t) - \sum_j V_{ij}^2(t) \right|$ for any time $t \geq 0$ and all rows i .

Notations: Boldface capital letters are reserved for matrices. Boldface lowercase letters denote vectors, and lowercase letters denote scalars. \mathbf{X}^\top is the transpose of matrix \mathbf{X} . X_{ij} is the (i, j) -th element of \mathbf{X} . $\mathbf{X}_{i\cdot}$ denotes the i -th row of any matrix \mathbf{X} . $\|\cdot\|_F, \|\cdot\|_2$ represent the Frobenius norm of a matrix, and the Euclidean norm of a vector. $\text{diag}(\mathbf{g})$ denotes a diagonal matrix with vector \mathbf{g} as its diagonal elements. We use $\phi(\mathbf{X}) = \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2$ to denote the squared Frobenius norm of the residual $\mathbf{Y} - \mathbf{A}\mathbf{X}$. The term $\mathbf{1}_{N \times L}$ is a matrix of all ones of size $N \times L$.

With these preliminaries established, we now introduce our IR framework for the MMV problem.

4 IR-MMV

Consider the task of recovering an MMV $\mathbf{X} \in \mathbb{R}^{N \times L}$, characterized by a common sparsity pattern across its rows. The measurement model is given by:

$$\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{W},$$

where $\mathbf{Y} \in \mathbb{R}^{M \times L}$ is the measurement matrix, $\mathbf{A} \in \mathbb{R}^{M \times N}$ is the sensing matrix (with $M \leq N$), and $\mathbf{W} \in \mathbb{R}^{M \times L}$ represents additive measurement noise. The objective is to recover \mathbf{X} , which is characterized by having at most K non-zero rows.

Algorithm 1 Implicit Regularization for MMV (IR-MMV)

- 1: **Input:** Sensing matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$, measurement matrix $\mathbf{Y} \in \mathbb{R}^{M \times L}$
 - 2: **Initialize:** $\alpha_g = 10^{-4}, \alpha_v = \frac{1}{\sqrt{2L}}\alpha_g$, learning rates $\eta_g = \eta_v = 10^{-2}$, number of iterations $T = 5 \times 10^6$, $\mathbf{g}(0) \in \mathbb{R}^N$ with $\mathbf{g}(0) = \alpha_g \mathbf{1}_L$, $\mathbf{V}(0) \in \mathbb{R}^{N \times L}$ with $\mathbf{V}(0) = \alpha_v \mathbf{1}_{N \times L}$.
 - 3: **Reconstruct initial signal:** $\mathbf{X}(0) = (\mathbf{g}(0)^{\odot 2} \mathbf{1}_L) \odot \mathbf{V}(0)$.
 - 4: **for** $t = 0, 1, \dots, T-1$ **do**
 - 5: $\mathbf{\Lambda} \leftarrow \mathbf{A}^\top (\mathbf{Y} - \mathbf{A}\mathbf{X})$
 - 6: Update parameters:

$$\begin{aligned} \mathbf{g}(t+1) &= \mathbf{g}(t) + 4\eta_g (\mathbf{g}(t) \odot ((\mathbf{\Lambda} \odot \mathbf{V}(t)) \mathbf{1}_L^\top)) \\ \mathbf{V}(t+1) &= \mathbf{V}(t) + 2\eta_v ((\mathbf{g}^{\odot 2}(t+1) \mathbf{1}_L) \odot \mathbf{\Lambda}) \end{aligned}$$
 - 7: Reconstruct signal: $\mathbf{X}(t+1) = (\mathbf{g}^{\odot 2}(t+1) \mathbf{1}_L) \odot \mathbf{V}(t+1)$
 - 8: **end for**
 - 9: **Output:** Estimated signal matrix \mathbf{X} .
-

Instead of directly optimizing for \mathbf{X} , we introduce an overparameterized representation of \mathbf{X} based on a Hadamard product factorization. This reparameterization is designed such that the inherent dynamics of gradient descent implicitly promote row sparsity. Specifically, we propose to factorize \mathbf{X} as:

$$\mathbf{X} = (\mathbf{g}^{\odot 2} \mathbf{1}_L) \odot \mathbf{V},$$

where $\mathbf{g} = [g_1, g_2, \dots, g_N]^\top \in \mathbb{R}^N$ acts as a row-wise scaling vector, and $\mathbf{V} \in \mathbb{R}^{N \times L}$ is a component matrix. The term $\mathbf{1}_L = [1, 1, \dots, 1]_{1 \times L}$ ensures that the squared elements of \mathbf{g} are replicated across columns. The square term on \mathbf{g} (i.e., $\mathbf{g}^{\odot 2}$) is introduced to keep the row activations positive. Recall that the proposed reparameterization when applied to gradient descent optimization will implicitly drive elements of \mathbf{g} corresponding to inactive rows towards zero. This phenomenon will encourage row sparsity in the reconstructed \mathbf{X} , leveraging the dynamics of the optimization process rather than explicit regularization terms. We define the mean squared loss as:

$$\mathcal{L}(\mathbf{g}, \mathbf{V}) = \|\mathbf{Y} - \mathbf{A}((\mathbf{g}^{\odot 2} \mathbf{1}_L) \odot \mathbf{V})\|_F^2.$$

The optimization is performed using gradient descent on the factors \mathbf{g} and \mathbf{V} . At each iteration, the updated factors are used to reconstruct \mathbf{X} . The gradient descent update rules for \mathbf{g} and \mathbf{V} are derived from the loss function $\mathcal{L}(\mathbf{g}, \mathbf{V})$. The gradient descent updates for \mathbf{g} and \mathbf{V} are:

$$\mathbf{g}(t+1) = \mathbf{g}(t) - \eta_g \nabla_{\mathbf{g}} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) \quad (1)$$

$$\mathbf{V}(t+1) = \mathbf{V}(t) - \eta_v \nabla_{\mathbf{V}} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) \quad (2)$$

The reconstructed signal matrix \mathbf{X} at time t is then given by:

$$\mathbf{X}(t+1) = (\mathbf{g}^{\odot 2}(t+1) \mathbf{1}_L) \odot \mathbf{V}(t+1). \quad (3)$$

The approach is outlined in Algorithm 1. In line with the analysis of matrix and tensor factorization by Gu-nasekar et al. (2017); Arora et al. (2018); Li et al. (2018); Arora et al. (2019); Razin et al. (2021); Hariz et al. (2024), we also model small learning rate for gradient descent through infinitesimal limit, i.e. through gradient flow:

$$\begin{aligned} \frac{d}{dt} g_l(t) &= -\nabla_{g_l} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)), \\ \frac{d}{dt} V_{lm}(t) &= -\nabla_{V_{lm}} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)), \end{aligned}$$

where $g_l(t)$ represents the element of factor \mathbf{g} corresponding to the l -th element at time t and $V_{lm}(t)$ corresponds to the (l, m) -th element of factor \mathbf{V} at time t . In the following section, we will formally analyze the dynamics of gradient flow to show that they guide the solution towards a row-sparse structure.

5 Theoretical Analysis

Let $\mathbf{\Lambda} = \mathbf{A}^\top (\mathbf{Y} - \mathbf{A}\mathbf{X})$ be the residual pre-multiplied by \mathbf{A}^\top .

Lemma 5.1 (Balancedness). *For the system evolving under continuous gradient flow for the set of update equations (1), (2) and (3), derived from the loss function $\mathcal{L}(\mathbf{g}, \mathbf{V}) = \|\mathbf{Y} - \mathbf{A}(\mathbf{g}^{\odot 2} \mathbf{1}_L \mathbf{V})\|_F^2$, the following properties hold for any time $t \geq 0$*

1. *Global balancedness: The quantity $\frac{1}{2} \|\mathbf{g}(t)\|_2^2 - \|\mathbf{V}(t)\|_F^2$ is conserved throughout the optimization process:*

$$\frac{1}{2} \|\mathbf{g}(t)\|_2^2 - \|\mathbf{V}(t)\|_F^2 = \frac{1}{2} \|\mathbf{g}(0)\|_2^2 - \|\mathbf{V}(0)\|_F^2$$

2. *Row-wise balancedness: For each individual row $i \in [N]$, the quantity $\frac{1}{2} g_i^2(t) - \sum_{j \in [L]} V_{ij}^2(t)$ is also conserved:*

$$\frac{1}{2} g_i^2(t) - \sum_{j \in [L]} V_{ij}^2(t) = \frac{1}{2} g_i^2(0) - \sum_{j \in [L]} V_{ij}^2(0)$$

Proof sketch (for detailed proof see B.1). The result follows by showing that under gradient flow, $\frac{1}{2} \sum_{l=1}^N \frac{d}{dt} g_l^2(t) = \sum_{l=1}^N \frac{d}{dt} \|\mathbf{V}_{l:}(t)\|_2^2$ for all $t \geq 0$. \square

Lemma 5.1 proves that if the parameters are balanced in the beginning at time $t = 0$ (i.e., $\frac{1}{2} \|\mathbf{g}(0)\|_2^2 = \|\mathbf{V}(0)\|_F^2$), then they will remain perfectly balanced over time, meaning $\frac{1}{2} \|\mathbf{g}(t)\|_2^2 = \|\mathbf{V}(t)\|_F^2$ for all subsequent time steps. It also shows that for each row i , the relative proportions between the scaling power of g_i and the “size” of the i -th row of \mathbf{V} (captured by its squared Frobenius norm) are preserved during learning. Here, the proposed IR mechanism for inducing sparsity acts consistently at a row-level. If the optimization drives $g_i(t)$ towards zero for a particular row i , then, due to the row-wise balancedness, the corresponding row $\mathbf{V}_{i:}$ must also approach zero, causing the entire i -th row of \mathbf{X} to become sparse.

Remark 5.1. *In the general case, ϵ and ϵ_r are independent quantities. However, under the specific initialization of Algorithm 1 (where $g(0)$ and $V(0)$ are initialized with constant values α_g and α_v), it strictly holds that $\epsilon_r \leq \epsilon$ for all $t \geq 0$.*

Lemma 5.2 (Dynamics of row norm). *Consider a system evolving under continuous gradient flow with updates as in (1), (2) and (3) with initialization $g(0) = \alpha_g \mathbf{1}_L$ and $V(0) = \alpha_v \mathbf{1}_{N \times L}$ where α_g and α_v are small positive scalars. Let $\mathbf{\Lambda}(t) := \mathbf{A}^\top (\mathbf{Y} - \mathbf{A}(\mathbf{g}(t)^{\odot 2} \mathbf{1}_L \odot \mathbf{V}(t)))$ and $\lambda_i(t)$ denote the i -th row of $\mathbf{\Lambda}(t)$, and*

$$\hat{\mathbf{x}}_i(t) := \begin{cases} \frac{\mathbf{X}_{i:}(t)}{\|\mathbf{X}_{i:}(t)\|_2} & \text{if } \|\mathbf{X}_{i:}(t)\|_2 \neq 0 \\ \mathbf{0} & \text{if } \|\mathbf{X}_{i:}(t)\|_2 = 0 \end{cases}.$$

Then, for any time $t \geq 0$ and for all $i \in [N]$ where $g_i^2(t) > 0$ and $\sum_{m \in [N]} V_{im}^2(t) > 0$,

if $\langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \geq 0$ then,

$$\begin{aligned} \frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 &\leq 24 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3} \right)^2 \\ \frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 &\geq 6 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\frac{\|\mathbf{X}_{i:}(t)\|_2^2}{\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3}} \right), \end{aligned}$$

and if $\langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle < 0$ then,

$$\begin{aligned} \frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 &\geq 24 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3} \right)^2 \\ \frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 &\leq 6 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\frac{\|\mathbf{X}_{i:}(t)\|_2^2}{\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3}} \right). \end{aligned}$$

Proof sketch (for detailed proof see Lemma B.3). The proof begins by expressing the time derivative of the row norm, $\frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2$, in terms of the underlying

parameters $g_i(t)$ and $\sum_j V_{ij}(t)$. The upper and lower bounds are then established by relating the parameter norms (g_i^2 and $\sum_j V_{ij}^2$) back to the matrix row norm ($\|\mathbf{X}_{i:}\|_2$) using the row-unbalancedness constant, ϵ . \square

Lemma 5.3 (Dynamics of row norms under perfect balancedness). *Assume that the system evolves under continuous gradient flow, as described in Lemma 5.1, with perfect row balancedness, i.e., $\frac{1}{2}g_i^2(t) = \sum_{j \in [L]} V_{ij}^2(t)$ for all $i \in [N]$ and for any time $t \geq 0$. Then, the rate of change of the Euclidean norm of the i -th row of $\mathbf{X}(t)$ is characterized by:*

$$\frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 = 2^{2/3} 6 \langle \lambda(t), \hat{\mathbf{x}}_i(t) \rangle \|\mathbf{X}_{i:}(t)\|_2^{4/3}.$$

Proof sketch (for detailed proof see Lemma B.4). The proof begins by expressing the derivative $\frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2$ in terms of the underlying parameters ($\mathbf{g}(t), \mathbf{V}(t)$). Then the result follows by applying row-balancedness assumption and substituting back the row norm. \square

The Lemmas 5.2 and 5.3 shows that the evolution rate of each row's norm is roughly proportional to its current norm raised to the power $4/3$. This implies that when the unbalancedness constant is small, this proportionality creates a ‘‘momentum-like’’ effect: smaller row norms evolve more slowly, while larger row norms change more rapidly. This dynamic, in turn, fosters an incremental learning effect: a phenomenon where different parts of the signal are learned sequentially rather than all at once. This process is realized by certain rows growing significantly in magnitude while others remain small, ultimately leading to an implicit regularization that favors low-rank solutions.

Lemma 5.4 (Distance bound in terms of parameters). *Let the estimated trajectory, $\mathbf{X}(t)$ be parametrized by $(\mathbf{g}(t), \mathbf{V}(t))$ and a reference trajectory $\tilde{\mathbf{X}}(t)$, be parameterized by $(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$ both of which evolve under the continuous gradient flow, as described in Lemma C.2. Let the parameters $(\mathbf{g}(t), \mathbf{V}(t))$ and $(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$ be at a distance D from the origin. At any time $t \in [0, T]$, the following inequality holds:*

$$\begin{aligned} \|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\|_F^2 &\leq 8D^4 \sum_{n \in [N]} \sum_{l \in [L]} \left((V_{nl}(t) - \tilde{V}_{nl}(t))^2 \right. \\ &\quad \left. + (g_n(t) - \tilde{g}_n(t))^2 \right) \end{aligned}$$

Proof sketch (for detailed proof see Lemma C.3). The proof begins by substituting the parameterizations for $\mathbf{X}(t)$ and $\tilde{\mathbf{X}}(t)$ into the squared Frobenius

distance, $\|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\|_F^2$. The result then follows from an algebraic rearranging of terms and the application of standard norm inequalities, including the triangle inequality. \square

This lemma provides the necessary link between the matrix trajectories and their underlying parameters. Specifically, the lemma bounds the distance between the reference trajectory, $\tilde{\mathbf{X}}(t)$, and the estimated trajectory, $\mathbf{X}(t)$ by the distance between their respective underlying parameters $(\mathbf{g}(t), \mathbf{V}(t)), (\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$.

Theorem 5.1. *Consider the system evolving under continuous gradient flow, where parameters are initialized such that perfect row-balancedness holds for all rows, i.e., $\frac{1}{2}g_i^2(0) = \sum_{j \in [L]} V_{ij}^2(0)$ for all $i \in [N]$. Let the sensing matrix \mathbf{A} have ℓ_2 -normalized columns and be μ -coherent. Let $\mathbf{X}(t) \in \mathbb{R}^{N \times L}$ be the estimated trajectory obtained by the gradient flow with the proposed parameterization $\mathbf{X}(t) = (\mathbf{g}(t)^{\odot 2} \mathbf{1}_L) \odot \mathbf{V}(t)$. Let β be the smoothness constant of the loss function, ρ be some constant $0 < \rho \leq \alpha_V^3$, $\tilde{D} = N^{\frac{1}{3}} \sqrt{2L+1} \left(\frac{D}{2} + 0.5 \right)^{\frac{1}{3}}$ and $D \geq 2\rho\sqrt{K}$. If the initial value α_V, α_g (where $g(0) = \alpha_g \mathbf{1}_L$ and $V(0) = \alpha_V \mathbf{1}_{N \times L}$, and $\alpha_V = \frac{1}{\sqrt{2L}} \alpha_g$) is sufficiently small, specifically satisfying:*

$$\alpha_V \leq \frac{\epsilon_{app}}{2(\tilde{D} + 2)^2} \frac{\exp(-\beta T)}{\sqrt{2LN(2L+3)}}, \quad (4)$$

where $\epsilon_{app} \in (0, 1)$, then for time $T > 0$, there exists a rank- K trajectory $\tilde{\mathbf{X}}(t)$ such that the estimate $\mathbf{X}(t)$ is close to $\tilde{\mathbf{X}}(t)$ until $t \geq T$ or $\|\mathbf{X}(t)\|_F \geq D$ i.e.,

$$\|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\|_F < \epsilon_{app}$$

for all $t \in [0, \min(T, T_D)]$, where $T_D := \inf\{t \geq 0 \mid \|\mathbf{X}(t)\|_F \geq D\}$.

Proof sketch (for detailed proof see Appendix C.1). We prove this by first constructing a reference rank- K trajectory, $\tilde{\mathbf{X}}(t)$, and then showing that the estimated trajectory, $\mathbf{X}(t)$, remains close to the reference rank- K trajectory for a sufficiently small initialization as defined in (4). The parameters of the estimated trajectory are initialized uniformly and balanced, i.e., $\mathbf{g}(0) = \alpha_g \mathbf{1}_L$ and $\mathbf{V}(0) = \alpha_V \mathbf{1}_{N \times L}$. To construct the reference rank- K trajectory, $\tilde{\mathbf{X}}(t)$, we begin by defining its own parameters $(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$, such that $\tilde{\mathbf{X}}(t) = (\tilde{\mathbf{g}}(t)^{\odot 2} \mathbf{1}_L) \odot \tilde{\mathbf{V}}(t)$. The dynamics of these parameters are governed by the gradient flow on the loss function $\mathcal{L}(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$. Crucially, for any support set \mathcal{S}_K , the corresponding parameters are initialized as, for $i \in [N]$,

$$\tilde{g}_i(0) = \begin{cases} \sqrt{2}\rho^{1/3} & \forall i \in \mathcal{S}_K, \\ 0 & \forall i \notin \mathcal{S}_K, \end{cases} \quad (5)$$

and

$$\tilde{V}_{ij}(0) = \begin{cases} \rho^{1/3} \frac{V_{ij}(0)}{\|\mathbf{V}_{i:}(0)\|} & \forall i \in \mathcal{S}_K, \forall j \in [L], \\ 0 & \forall i \notin \mathcal{S}_K \end{cases} \quad (6)$$

where, $\rho < \alpha_V^3$ is a small positive scalar. Because these non-support rows start at zero, Lemma 5.1 and Lemma 5.3, ensure they remain zero for all time $t \geq 0$. This construction guarantees the existence of a $\tilde{\mathbf{X}}(t)$, which is at most rank- K matrix throughout its evolution. Note that this construction is purely for this theoretical proof; the proposed algorithm does not assume any knowledge of the support set or sparsity. To bound the distance between $\mathbf{X}(t)$ and $\tilde{\mathbf{X}}(t)$, we begin by bounding the distance between the corresponding parameters, namely, $(\mathbf{g}(t), \mathbf{V}(t))$ and $(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$. To bound the evolution of this distance, we rely on Razin et al. (2021)[Lemma 6]. It states that for a locally smooth objective, the distance between two gradient flow trajectories will grow at most exponentially provided the initial distance between the trajectories remains sufficiently small. Our initialization scheme (4), along with (5) and (6), ensures that at $t = 0$, the distance between the parameters remains sufficiently small. Since our loss function is locally smooth, the parameter trajectories remain close.

Finally, the closeness of the parameter trajectories is translated to the distance between the trajectories $\mathbf{X}(t)$ and $\tilde{\mathbf{X}}(t)$ by applying the bound from Lemma 5.4. This leads to the theorem’s main guarantee: for a sufficiently small initialization α_V, α_g , the estimated solution $\mathbf{X}(t)$ is guaranteed to remain within a user-defined distance ϵ_{app} of the reference row-sparse trajectory $\tilde{\mathbf{X}}(t)$ for a specified time duration T , or until the trajectory leaves a defined ball of radius D around the origin. \square

This theorem provides the theoretical guarantee for the proposed framework. It formalizes the implicit regularization effect by establishing that with a sufficiently small and balanced initialization, the algorithm’s estimated trajectory, $\mathbf{X}(t)$ is guaranteed to stay within a predefined error (ϵ_{app}) of a reference rank- K trajectory, $\tilde{\mathbf{X}}(t)$. By starting near the origin, the dynamics ensure that only a few rows that happen to have a slightly larger initial correlation with the true signal are selected to grow rapidly. Meanwhile balanced initialization ensures the proportional evolution of the factors required for the momentum-like dynamics to guide the algorithm’s trajectory toward the row-sparse solution.

The following corollary extends the theoretical guarantees of Theorem 5.1. It establishes that if all reference trajectories converges uniformly to a global minimum

of the objective function, then the estimated trajectory produced by the gradient flow of our proposed over-parameterization will likewise converge to that same solution.

Corollary 5.1.1. *Assume the conditions of Theorem 5.1 and in addition, assume that all reference trajectories $\tilde{\mathbf{X}}(t)$ converge to a solution $\mathbf{X}^* \in \mathbb{R}^{N \times L}$ and this convergence is uniform in the sense that the trajectories are confined to a bounded domain, and for any $\epsilon_{app} > 0$ there exists a time $T_c \leq T$ after which they are all within a distance ϵ_{app} from \mathbf{X}^* . Then for any $\epsilon_{app} > 0$, if the initialization scales α_V, α_g are sufficiently small, for any time $t \in [T_c, T]$ it holds that $\|\mathbf{X}(t) - \mathbf{X}^*\|_F \leq 2\epsilon_{app}$.*

Proof sketch (for detailed proof see Appendix C.2).

We assume that all reference trajectories $\tilde{\mathbf{X}}(t)$ converge to the oracle solution \mathbf{X}^* within a finite time $T > 0$. This means that for any desired approximation error $\epsilon_{app} > 0$ we have $\|\tilde{\mathbf{X}}(t) - \mathbf{X}^*\| \leq \epsilon_{app}$.

From theorem 5.1, if the initialization values α_V, α_g are sufficiently small, the estimated trajectory is guaranteed to be within a distance of ϵ_{app} from the reference trajectory at least until time T . This is a conditional guarantee that holds as long as both trajectories remain within a bounded domain. We show that the estimated trajectory $\mathbf{X}(t)$ is also confined to a bounded domain. We know from our assumptions that the reference trajectory $\tilde{\mathbf{X}}(t)$ is bounded. The result follows from applying the triangle inequality. \square

6 Simulation

In this section, we present simulations to reinforce our theoretical results. We also compare the performance of our proposed IR-MMV approach against three established MMV recovery algorithms: M-OMP, M-SP and M-FOCUSS. Unless otherwise specified, the default values for simulation parameters are: $M = 50$, $N = 25$, $L = 100$, row sparsity level $K = 3$ ($K < L$). The ground truth MMV, \mathbf{X} , is constructed such that its active rows are set to constant values of $\mathbf{1}_L$. The sensing matrix \mathbf{A} is generated by drawing its entries from an i.i.d. sub-Gaussian distribution (to satisfy μ -coherence property) with ℓ_2 -normalized columns. In line with our theoretical analysis, the algorithm’s hyperparameters are configured as follows: \mathbf{g} is initialized with $\alpha_g = 10^{-4}$ to keep the initialization close to origin, while the component matrix \mathbf{V} is initialized to $\alpha_V = \frac{1}{\sqrt{2L}}\alpha_g$ to ensure the perfect row-wise balancedness. The learning rates for the gradient descent updates for g (i.e., η_g) and V (i.e., η_v) are both set to 10^{-2} , a small, constant value chosen for stable

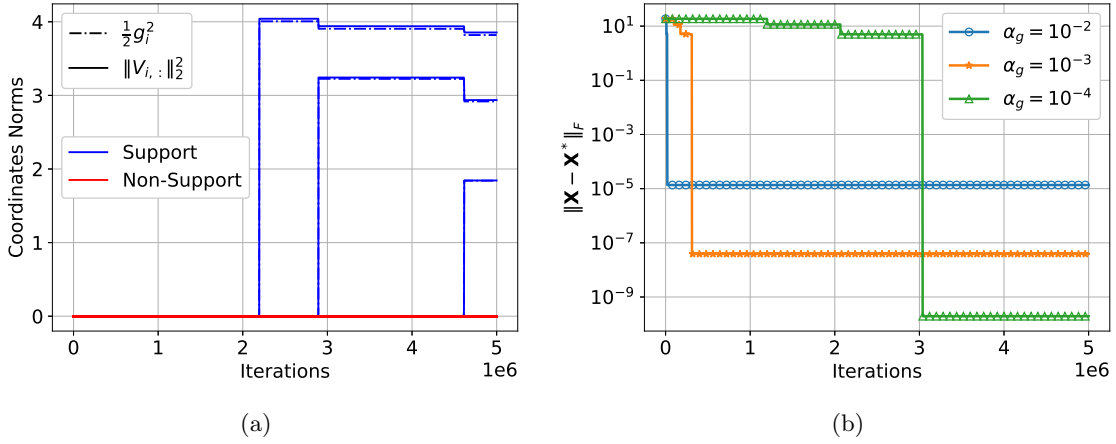


Figure 1: (a) Evolution of the norms of the components for rows in the support and non-support sets. (b) Loss versus iterations for different initialization values of α .

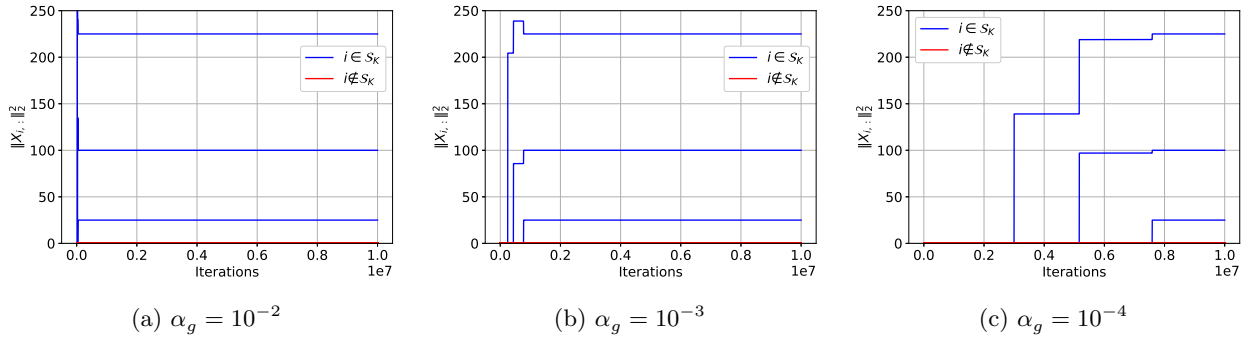


Figure 2: The evolution of coordinate norms during training for different initialization parameters α_g .

gradient descent. Finally, the maximum number of iterations is set to $T = 10^7$ to provide the algorithm with sufficient time to converge from its deliberately small initial state. For our experiments, the parameters of M-FOCUSS are set to $p = 0.8$ and λ set to noise variance as suggested by Cotter et al. (2005).

6.1 Balancedness and Incremental Learning

For this experiment, the ground truth MMV, \mathbf{X} , is constructed such that its active rows are set to constant values of $\mathbf{1}_L, 2 \cdot \mathbf{1}_L$ and $3 \cdot \mathbf{1}_L$. Recall that the parameters are initialized to ensure perfect row-wise balancedness. As illustrated in Figure 1a, the evolution of $\frac{1}{2}g_i^2(t)$ and $\|V_{i,:}(t)\|_2^2$ is tracked throughout the optimization process. The empirical result support the theoretical findings of Lemma 5.1, by demonstrating that this balance is preserved over time. This dynamic leads to a momentum-like effect and an incremental learning phenomenon, where rows of the MMV are learned one at a time, a behavior consistent with prior findings by Gissin et al. (2019); Razin et al. (2021);

Hariz et al. (2024) in the context of implicit regularization via overparametrization. The magnitude of the row norms directly influences the learning speed, with smaller norms evolving more slowly and larger norms changing more rapidly. This dynamic implicitly promotes structured sparsity. The time-invariant balancedness property (Lemma 5.1) guarantees that if an element $g_i(t)$ approaches zero, the corresponding row $V_{i,:}$ will also converge to zero, rendering the entire i -th row of \mathbf{X} sparse.

6.2 Effect of Initialization

To analyze the effect of initialization, the experiment is repeated for a set of different initializations, specifically $\alpha_g \in \{10^{-2}, 10^{-3}, 10^{-4}\}$ and the corresponding α_V . All other parameters, are kept constant across these runs. For this experiment, the ground truth MMV, \mathbf{X} , is generated as in the previous section. The simulation results highlight the importance of initialization values. As seen in Figure 1b, the convergence behavior of the loss function varies significantly with

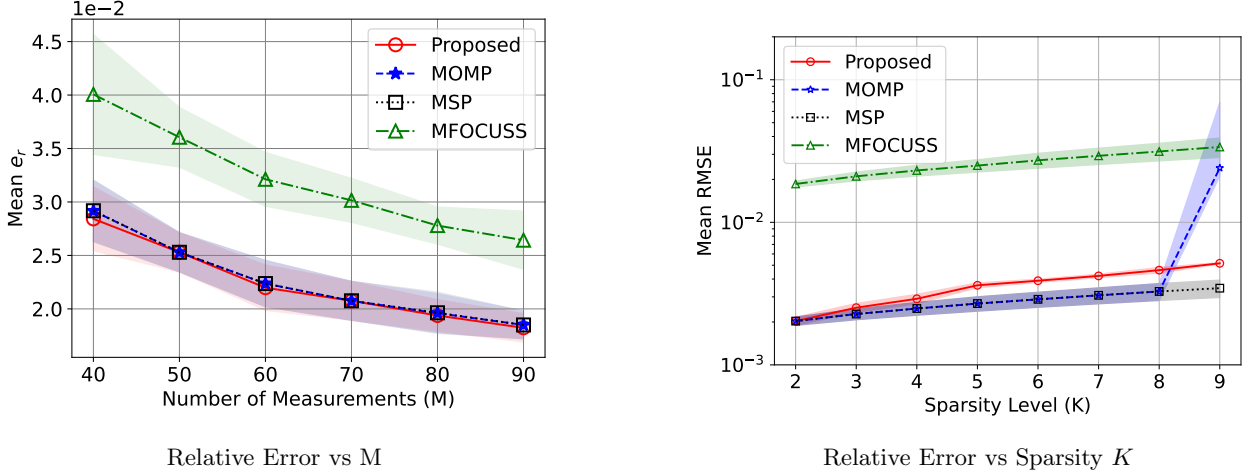


Figure 3: Performance of IR-MMV averaged over 20 trials.

the choice of α_g . It demonstrates that smaller initialization values lead to convergence to a solution with a lower final loss. This empirical finding is consistent with Theorem 5.1, which formally establishes that a smaller initialization value (α_V) is necessary to guarantee a smaller approximation error (ϵ_{app}). Furthermore, Figure 2 illustrates how growth of coordinate norms is influenced by the initialization. A smaller value results in a slower, more gradual convergence, while larger initial values may lead to faster but less stable convergence.

6.3 Performance with Varying Noise

The empirical performance of IR-MMV is compared against different MMV reconstruction algorithms (MOMP, M-SP and M-FOCUSS) for various values of M and K . For these experiments, the measurement matrix \mathbf{Y} is generated by adding noise to the signal ($\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{W}$), such that the Signal-to-Noise Ratio is 40dB. Performance analysis is evaluated using relative estimation error $e_r = \frac{\|\mathbf{X} - \hat{\mathbf{X}}\|_F^2}{\|\mathbf{X}\|_F^2}$. The results, averaged over 20 trials, are shown in Figure 3; the shaded regions represent the standard deviation across the trials. The proposed method demonstrates performance comparable to the other approaches. Critically, it achieves this without any parameter tuning or prior knowledge of data-specific parameters. In contrast, the benchmark algorithms require such information to perform optimally; both M-OMP and M-SP requires prior knowledge of true sparsity level K . M-FOCUSS requires proper tuning of a parameter λ and/or knowledge of noise variance. We defer the results on MNIST data to Appendix D.

7 Conclusion

We introduce a novel overparameterized factorization based on the Hadamard product to induce IR for structured sparse recovery of MMV. We show that the optimization dynamics naturally induce an implicit bias toward row-sparse solutions. Our characterization of the row norm dynamics identifies a momentum-like effect, where rows with larger norms grow significantly faster than those with smaller norms, leading to an incremental learning process that provably favors low-rank solutions. The theoretical analysis, validated through simulations, establishes two key properties of the optimization dynamics. We prove that the relative scaling between the parameters is preserved both globally and for each row throughout the optimization process. We also provide a formal convergence guarantee that with a sufficiently small and balanced initialization, the algorithm’s trajectory is guaranteed to remain close to an idealized row-sparse solution. We show that the performance of the proposed method is comparable to the traditional algorithms like M-OMP, M-SP and M-FOCUSS. Critically, the performance of the proposed framework is achieved without requiring prior knowledge of the sparsity level of the signal or its noise variance. This research successfully bridges the gap between the theoretical understanding of implicit regularization and its practical application, offering a versatile and tuning-free tool for the complex domain of MMV recovery.

References

Arora, S., Cohen, N., and Hazan, E. (2018). On the optimization of deep networks: Implicit acceleration by overparameterization. In *International conference on machine learning*, pages 244–253. PMLR.

- Arora, S., Cohen, N., Hu, W., and Luo, Y. (2019). Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32.
- Chen, D., Tian, L., and Xu, C. (2020). MMV-based OMP for DOA estimation with 1-bit measurement. In *Journal of Physics: Conference Series*, volume 1550, page 032150. IOP Publishing.
- Chen, Z., Xiang, J., Bagnaninchi, P.-O., and Yang, Y. (2022). MMV-Net: A multiple measurement vector network for multifrequency electrical impedance tomography. *IEEE Transactions on Neural Networks and Learning Systems*, 34(11):8938–8949.
- Chou, H.-H., Maly, J., and Rauhut, H. (2023). More is less: inducing sparsity via overparameterization. *Information and Inference: A Journal of the IMA*, 12(3):1437–1460.
- Cotter, S. F., Rao, B. D., Engan, K., and Kreutz-Delgado, K. (2005). Sparse solutions to linear inverse problems with multiple measurement vectors. *IEEE Transactions on signal processing*, 53(7):2477–2488.
- Gissin, D., Shalev-Shwartz, S., and Daniely, A. (2019). The implicit bias of depth: How incremental learning drives generalization. *arXiv preprint arXiv:1909.12051*.
- Gunasekar, S., Woodworth, B. E., Bhojanapalli, S., Neyshabur, B., and Srebro, N. (2017). Implicit regularization in matrix factorization. *Advances in neural information processing systems*, 30.
- Hariz, K., Kadri, H., Ayache, S., Moakher, M., and Artières, T. (2024). Implicit regularization in deep tucker factorization: Low-rankness via structured sparsity. In *International Conference on Artificial Intelligence and Statistics*, pages 2359–2367. PMLR.
- Li, J., Nguyen, T., Hegde, C., and Wong, K. W. (2021). Implicit sparse regularization: The impact of depth and early stopping. *Advances in Neural Information Processing Systems*, 34:28298–28309.
- Li, J., Nguyen, T. V., Hegde, C., and Wong, R. K. (2023). Implicit regularization for group sparsity. *arXiv preprint arXiv:2301.12540*.
- Li, Y., Ma, T., and Zhang, H. (2018). Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47. PMLR.
- Liu, S., Zheng, L., Liu, L., and Lin, Q. (2019). MMV subspace pursuit (M-SP) algorithm for joint sparse multiple measurement vectors recovery. In *2019 IEEE 13th International Conference on ASIC (ASICON)*, pages 1–4. IEEE.
- Razin, N., Maman, A., and Cohen, N. (2021). Implicit Regularization in Tensor Factorization. *Proceedings of Machine Learning Research*, 139:8913–8924.
- Soltanolkotabi, M., Stöger, D., and Xie, C. (2023). Implicit balancing and regularization: Generalization and convergence guarantees for overparameterized asymmetric matrix sensing. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 5140–5142. PMLR.
- Teschl, G. (2012). *Ordinary differential equations and dynamical systems*, volume 140. American Mathematical Soc.
- Vaskevicius, T., Kanade, V., and Rebeschini, P. (2019). Implicit regularization for optimal sparse recovery. *Advances in Neural Information Processing Systems*, 32.
- Wu, X., Dobriban, E., Ren, T., Wu, S., Li, Z., Gunasekar, S., Ward, R., and Liu, Q. (2020). Implicit regularization and convergence for weight normalization. *Advances in Neural Information Processing Systems*, 33:2835–2847.
- Yang, D. and Li, H. (2024). MMV-Net: A multiple measurement vector network for MIMO channel estimation. In *2024 9th International Conference on Computer and Communication Systems (ICCCS)*, pages 1027–1032. IEEE.
- Zhang, X., Xie, L., and Wang, J. (2022). Some results on OMP algorithm for MMV problem. *Mathematical Methods in the Applied Sciences*, 45(9).
- Zhao, P., Yang, Y., and He, Q.-C. (2022). High-dimensional linear regression via implicit regularization. *Biometrika*, 109(4):1033–1046.

Checklist

1. For all models and algorithms presented, check if you include:
 - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes] See Definition 3.1 and Theorem 5.1 statement.
 - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Yes] See section 5
 - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes]
2. For any theoretical claim, check if you include:
 - (a) Statements of the full set of assumptions of all theoretical results. [Yes]

- (b) Complete proofs of all theoretical results. [Yes] See appendix B and C
 - (c) Clear explanations of any assumptions. [Yes]
3. For all figures and tables that present empirical results, check if you include:
- (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes] See section E
 - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes] See section 6
 - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes]
 - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Yes] See section D
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
- (a) Citations of the creator If your work uses existing assets. [Yes]
 - (b) The license information of the assets, if applicable. [Yes]
 - (c) New assets either in the supplemental material or as a URL, if applicable. [Not Applicable]
 - (d) Information about consent from data providers/curators. [Not Applicable]
 - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Not Applicable]
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
- (a) The full text of instructions given to participants and screenshots. [Not Applicable]
 - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Not Applicable]
 - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Not Applicable]

A Gradient Updates

This section details the derivation of the gradient descent update rules for the parameters \mathbf{g} and \mathbf{V} . We consider the loss function for the proposed parameterization, $\mathbf{X} = \mathbf{g}^{\odot 2} \mathbf{1}_L \odot \mathbf{V}$, which is reproduced here for readability:

$$\mathcal{L}(\mathbf{g}, \mathbf{V}) = \|\mathbf{Y} - \mathbf{A}(\mathbf{g}^{\odot 2} \mathbf{1}_L \odot \mathbf{V})\|_F^2.$$

Let $\boldsymbol{\Lambda} := \mathbf{A}^\top (\mathbf{Y} - \mathbf{A}\mathbf{X})$. We begin by deriving the partial derivative of the loss function $\mathcal{L}(\mathbf{g}, \mathbf{V})$ with respect to a single element g_l (the l -th element of \mathbf{g}):

$$\begin{aligned} \nabla_{g_l} \mathcal{L}(\mathbf{g}, \mathbf{V}) &= \nabla_{g_l} \left(\sum_{i \in [M]} \sum_{j \in [L]} \left(Y_{ij} - \sum_{k \in [N]} A_{ik} g_k^2 V_{kj} \right)^2 \right) \\ &= -4 \sum_{i \in [M]} \sum_{j \in [L]} \left(Y_{ij} - \sum_{k \in [N]} A_{ik} g_k^2 V_{kj} \right) A_{il} g_l V_{lj} \\ &= -4 g_l \sum_{i \in [M]} \sum_{j \in [L]} A_{il} \left(Y_{ij} - \sum_{k \in [N]} A_{ik} g_k^2 V_{kj} \right) V_{lj} \\ &= -4 g_l \sum_{j \in [L]} \mathbf{A}^\top (\mathbf{Y} - \mathbf{A}(\mathbf{g}^{\odot 2} \mathbf{1}_L \odot \mathbf{V}))_{lj} V_{lj} \\ &= -4 g_l \sum_{j \in [L]} \Lambda_{lj} V_{lj}. \end{aligned} \tag{7}$$

In vector form, the gradient with respect to \mathbf{g} is given by:

$$\nabla_{\mathbf{g}} \mathcal{L}(\mathbf{g}, \mathbf{V}) = -4 \mathbf{g} \odot ((\boldsymbol{\Lambda} \odot \mathbf{V}) \mathbf{1}_L^\top).$$

Similarly, we compute the gradient of the loss function with respect to the element V_{lm} (the (l, m) -th element of \mathbf{V}):

$$\begin{aligned} \nabla_{V_{lm}} \mathcal{L}(\mathbf{g}, \mathbf{V}) &= \nabla_{V_{lm}} \sum_{i \in [M]} \sum_j \left(Y_{ij} - \sum_{k \in [N]} A_{ik} g_k^2 V_{kj} \right)^2 \\ &= -2 \sum_{i \in [M]} \left(Y_{im} - \sum_{k \in [N]} A_{ik} g_k^2 V_{km} \right) A_{il} g_l^2 \\ &= -2 g_l^2 \sum_{i \in [M]} A_{il} \left(Y_{im} - \sum_{k \in [N]} A_{ik} g_k^2 V_{km} \right) \\ &= -2 g_l^2 (\mathbf{A}^\top (\mathbf{Y} - \mathbf{A}(\mathbf{g}^{\odot 2} \mathbf{1}_L \odot \mathbf{V}))_{lm}) \\ &= -2 g_l^2 \Lambda_{lm}. \end{aligned} \tag{8}$$

In matrix form this can be written as:

$$\nabla_{\mathbf{V}} \mathcal{L}(\mathbf{g}, \mathbf{V}) = -2 \mathbf{g}^{\odot 2} \mathbf{1}_L \odot \boldsymbol{\Lambda}$$

B Characterization of the gradient flow

This section provides the lemmas and proofs that characterize the behavior of the system under the continuous gradient flow. In the following lemma, we establish a key ‘‘balancedness’’ property that is conserved by the gradient flow dynamics. Throughout this work, we denote the i -th row of any matrix \mathbf{Z} as $\mathbf{Z}_{i:}$ and let $\mathbf{1}_{N \times L}$ denote a matrix of ones of size $N \times L$.

Lemma B.1 (Lemma 5.1 restated). *For the system evolving under continuous gradient flow for the set of updates (1), (2) and (3), derived from the loss function $\mathcal{L}(\mathbf{g}, \mathbf{V}) = \|\mathbf{Y} - \mathbf{A}(\mathbf{g}^{\odot 2} \mathbf{1}_L \mathbf{V})\|_F^2$, the following properties hold for any time $t \geq 0$:*

1. *Global balancedness: The quantity $\frac{1}{2} \|\mathbf{g}(t)\|_2^2 - \|\mathbf{V}(t)\|_F^2$ is conserved throughout the optimization process:*

$$\frac{1}{2} \|\mathbf{g}(t)\|_2^2 - \|\mathbf{V}(t)\|_F^2 = \frac{1}{2} \|\mathbf{g}(0)\|_2^2 - \|\mathbf{V}(0)\|_F^2$$

2. *Row-wise balancedness: For each individual row $i \in [N]$, the quantity $\frac{1}{2} g_i^2(t) - \sum_{j \in [L]} V_{ij}^2(t)$ is also conserved:*

$$\frac{1}{2} g_i^2(t) - \sum_{j \in [L]} V_{ij}^2(t) = \frac{1}{2} g_i^2(0) - \sum_{j \in [L]} V_{ij}^2(0)$$

Proof. Let $\mathbf{\Lambda}(t) := \mathbf{A}^\top (\mathbf{Y} - \mathbf{A}(\mathbf{g}^{\odot 2}(t) \mathbf{1}_L \odot \mathbf{V}(t)))$. Gradient flow is related to gradient as follows:

$$\frac{d}{dt} g_l(t) = -\nabla_{g_l} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)).$$

From (7),

$$\frac{d}{dt} g_l(t) = 4g_l(t) \sum_{j \in [L]} \Lambda_{lj}(t) V_{lj}(t). \quad (9)$$

Similarly, from (8),

$$\begin{aligned} \frac{d}{dt} V_{lm}(t) &= -\nabla_{V_{lm}} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)), \\ &= 2g_l^2(t) \Lambda_{lm}(t). \end{aligned} \quad (10)$$

$$\begin{aligned} \frac{d}{dt} \|\mathbf{V}_{l:}(t)\|_2^2 &= \frac{d}{dt} \left(\sum_{m \in [L]} V_{lm}^2(t) \right) = \sum_{m \in [L]} \frac{d}{dt} V_{lm}^2(t) \\ &= 2 \sum_{m \in [L]} V_{lm}(t) \frac{d}{dt} V_{lm}(t). \end{aligned}$$

Substituting (10)

$$\frac{d}{dt} \|\mathbf{V}_{l:}(t)\|_2^2 = 4g_l^2(t) \sum_{m \in [L]} V_{lm}(t) \Lambda_{lm}(t). \quad (11)$$

From (9) and (11),

$$\begin{aligned} g_l(t) \frac{d}{dt} g_l(t) &= \frac{d}{dt} \|\mathbf{V}_{l:}(t)\|_2^2 \\ \frac{1}{2} \frac{d}{dt} g_l^2(t) &= \frac{d}{dt} \|\mathbf{V}_{l:}(t)\|_2^2 \\ \frac{1}{2} \sum_{l \in [N]} \frac{d}{dt} g_l^2(t) &= \sum_{l \in [N]} \frac{d}{dt} \|\mathbf{V}_{l:}(t)\|_2^2 \\ \frac{1}{2} \frac{d}{dt} \|\mathbf{g}(t)\|_2^2 &= \frac{d}{dt} \|\mathbf{V}(t)\|_F^2. \end{aligned} \quad (12)$$

$$\frac{1}{2} \frac{d}{dt} \|\mathbf{g}(t)\|_2^2 = \frac{d}{dt} \|\mathbf{V}(t)\|_F^2. \quad (13)$$

Here, $\mathbf{V}_{l:}(t)$ is the l -th row of matrix $\mathbf{V}(t)$. Integrating both sides of Equation (13) with respect to t from 0 to t directly yields the relationship:

$$\frac{1}{2} (\|\mathbf{g}(t)\|_2^2 - \|\mathbf{g}(0)\|_2^2) = \|\mathbf{V}(t)\|_F^2 - \|\mathbf{V}(0)\|_F^2.$$

Rearranging we get that

$$\frac{1}{2}\|\mathbf{g}(t)\|_2^2 - \|\mathbf{V}(t)\|_F^2 = \frac{1}{2}\|\mathbf{g}(0)\|_2^2 - \|\mathbf{V}(0)\|_F^2.$$

The second result follows from integrating both sides of (12) with respect to t from 0 to t . \square

Remark B.1 (Remark 5.1 restated). *In the general case, ϵ and ϵ_r are independent quantities. However, under the specific initialization of Algorithm 1 (where $\mathbf{g}(0)$ and $\mathbf{V}(0)$ are initialized with constant values α_g and α_V), it strictly holds that $\epsilon_r \leq \epsilon$ for all $t \geq 0$.*

Proof. Recall that the algorithm initializes the parameters as constants across all indices: $g_i(0) = \alpha_g$ for all $i \in [N]$ and $V_{ij}(0) = \alpha_V$ for all $i \in [N], j \in [L]$.

Let $e_i(t)$ denote the unbalancedness of the i -th row at time t :

$$e_i(t) = \frac{1}{2}g_i^2(t) - \sum_{j \in [L]} V_{ij}^2(t).$$

At initialization ($t = 0$), we have:

$$e_i(0) = \frac{1}{2}\alpha_g^2 - L\alpha_V^2 = c, \quad \forall i \in [N],$$

where c is a constant scalar. Lemma B.1 demonstrates that this row-wise unbalancedness is conserved throughout the gradient flow. Therefore, $e_i(t) = e_i(0) = c$ for all $t \geq 0$.

Substituting this into the definition of the global unbalancedness constant ϵ (Definition 3.3), we obtain:

$$\epsilon(t) = \left| \sum_{i=1}^N e_i(t) \right| = \left| \sum_{i=1}^N c \right| = N|c|.$$

Similarly, for the row-unbalancedness constant ϵ_r (Definition 3.4), we have:

$$\epsilon_r(t) = \max_{i \in [N]} |e_i(t)| = \max_{i \in [N]} |c| = |c|.$$

Comparing the two quantities, since $N \geq 1$:

$$\epsilon(t) = N\epsilon_r(t) \implies \epsilon_r(t) \leq \epsilon(t).$$

Finally, since these quantities are conserved throughout time, the time indexing (t) can be dropped, yielding:

$$\epsilon_r(t) \leq \epsilon(t).$$

\square

Lemma B.2. *For the system evolving under continuous gradient flow, as described in Lemma 5.1, with $\mathbf{\Lambda}(t) := \mathbf{A}^\top(\mathbf{Y} - \mathbf{A}(\mathbf{g}^{\odot 2}(t)\mathbf{1}_L \odot \mathbf{V}(t)))$ and any row $l \in [N]$, the following relationship hold for any time $t \geq 0$:*

$$\frac{1}{2} \frac{d}{dt} g_l^2(t) = \sum_{j \in [L]} \frac{d}{dt} V_{lj}^2(t) = 4 \sum_{j \in [L]} (\mathbf{\Lambda}(t) \odot \mathbf{X}(t))_{lj} = 4 (\mathbf{\Lambda}(t) \mathbf{X}^\top(t))_{ll}$$

Proof.

From (11) and (12),

$$\frac{1}{2} \frac{d}{dt} g_l^2(t) = \frac{d}{dt} \|\mathbf{V}_{l:}(t)\|_2^2 = 4g_l^2(t) \sum_{j \in [L]} V_{lj}(t) \Lambda_{lj}(t) \tag{14}$$

$$= 4 \sum_{j \in [L]} \Lambda_{lj}(t) X_{lj}(t). \tag{15}$$

The last line follows from rearranging and substituting for $X_{lj}(t) = g_l^2(t)V_{lj}(t)$ in (14).

In matrix form, equation (15) can be written as $4 \sum_{j \in [L]} (\mathbf{\Lambda}(t) \odot \mathbf{X}(t))_{lj}$ or $4 (\mathbf{\Lambda}(t) \mathbf{X}^\top(t))_{ll}$. \square

Lemma B.3 (Lemma 5.2 restated). *For the system evolving under continuous gradient flow with updates as in (1), (2) and (3) with initialization $g(0) = \alpha_g \mathbf{1}_L$ and $V(0) = \alpha_V \mathbf{1}_{N \times L}$ where α_g and α_V are small positive scalars. Let $\lambda_i(t)$ denote the i -th row of $\Lambda(t)$, and*

$$\hat{\mathbf{x}}_i(t) := \begin{cases} \frac{\mathbf{X}_{i:}(t)}{\|\mathbf{X}_{i:}(t)\|_2} & \text{if } \|\mathbf{X}_{i:}(t)\|_2 \neq 0 \\ \mathbf{0} & \text{if } \|\mathbf{X}_{i:}(t)\|_2 = 0 \end{cases}.$$

Then, for any time $t \geq 0$ and for all $i \in [N]$ where $g_i^2(t) > 0$ and $\sum_{m \in [N]} V_{im}^2(t) > 0$, if $\langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \geq 0$,

$$\begin{aligned} \frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 &\leq 24 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3} \right)^2 \\ \frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 &\geq 6 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\frac{\|\mathbf{X}_{i:}(t)\|_2^2}{\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3}} \right), \end{aligned}$$

if $\langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle < 0$:

$$\begin{aligned} \frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 &\geq 24 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3} \right)^2 \\ \frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 &\leq 6 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\frac{\|\mathbf{X}_{i:}(t)\|_2^2}{\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3}} \right). \end{aligned}$$

Proof. Note that

$$\langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle = \left(\Lambda(t) \frac{\mathbf{X}^\top(t)}{\|\mathbf{X}_{i:}(t)\|_2} \right)_{ii}. \quad (16)$$

With this identity established, we now derive the dynamics of the row norm $\|\mathbf{X}_{i:}(t)\|_2$:

$$\frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 = \frac{1}{2\|\mathbf{X}_{i:}(t)\|_2} \frac{d}{dt} \left(\sum_j X_{ij}^2(t) \right). \quad (17)$$

Recall that $X_{ij}(t) = g_i^2(t) V_{ij}(t)$. Using Lemma B.2 and equation (9), we have

$$\begin{aligned} \frac{d}{dt} \sum_j X_{ij}^2(t) &= 4g_i^4(t) (\Lambda(t) \mathbf{X}^\top(t))_{ii} + 16 \sum_{j \in [L]} V_{ij}^2(t) g_i^2(t) \sum_k \Lambda_{ik}(t) X_{ik}(t) \\ &= 4 \left(g_i^4(t) (\Lambda(t) \mathbf{X}^\top(t))_{ii} + 4 \sum_{j \in [L]} V_{ij}^2(t) g_i^2(t) (\Lambda(t) \mathbf{X}^\top(t))_{ii} \right) \\ &= 4 (\Lambda(t) \mathbf{X}^\top(t))_{ii} \left(g_i^4(t) + 4 \sum_{j \in [L]} V_{ij}^2(t) g_i^2(t) \right) \end{aligned} \quad (18)$$

Substituting equation (18) in (17),

$$\frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 = \frac{2}{\|\mathbf{X}_{i:}(t)\|_2} (\Lambda(t) \mathbf{X}^\top(t))_{ii} \left(g_i^4(t) + 4 \sum_{j \in [L]} V_{ij}^2(t) g_i^2(t) \right) \quad (19)$$

$$\begin{aligned} &\stackrel{a}{\leq} 24 \sum_{j \in [L]} \Lambda_{ij}(t) \frac{X_{ij}(t)}{\|\mathbf{X}_{i:}(t)\|_2} \left(\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3} \right)^2 \\ &\leq 24 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3} \right)^2. \end{aligned} \quad (20)$$

Note that in this case, $\langle \lambda_i(t), \hat{\mathbf{x}}(t)_i \rangle$ is non-negative. From the definition of unbalancedness, we have $\left| \frac{1}{2}g_i^2(t) - \sum_{m \in [L]} V_{im}^2(t) \right| \leq \epsilon_r$. The inequality in step (a) follows by substituting (16), (21) and (22) in (19).

$$\begin{aligned} \sum_{m \in [L]} V_{im}^2(t) &\leq \epsilon_r + \min \left\{ \sum_{m \in [L]} V_{im}^2(t), \frac{1}{2}g_i^2(t), \frac{1}{2}g_i^2(t) \right\} \\ &\leq \epsilon_r + \left(\frac{1}{4}g_i^4(t) \sum_{m \in [L]} V_{im}^2(t) \right)^{\frac{1}{3}}. \end{aligned}$$

Recall that $X_{im}(t) = g_i^2(t)V_{im}(t)$:

$$\sum_{m \in [L]} V_{im}^2(t) \leq \epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3}. \quad (21)$$

Similarly, we have:

$$\frac{1}{2}g_i^2(t) \leq \epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3}. \quad (22)$$

For the lower bound, we start from (19) and rewrite the expression as follows:

$$\begin{aligned} \frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 &= \frac{2}{\|\mathbf{X}_{i:}(t)\|_2} (\mathbf{\Lambda}(t) \mathbf{X}^\top(t))_{ii} \left(g_i^4(t) + 4 \sum_{j \in [L]} V_{ij}^2(t) g_i^2(t) \right) \\ &= \frac{2}{\|\mathbf{X}_{i:}(t)\|_2} (\mathbf{\Lambda}(t) \mathbf{X}^\top(t))_{ii} \left(\frac{g_i^4(t) \sum_{j \in [L]} V_{ij}^2(t)}{\sum_{j \in [L]} V_{ij}^2(t)} + 4 \frac{\sum_{j \in [L]} V_{ij}^2(t) g_i^4(t)}{g_i^2(t)} \right) \\ &= \frac{2}{\|\mathbf{X}_{i:}(t)\|_2} (\mathbf{\Lambda}(t) \mathbf{X}^\top(t))_{ii} \sum_{j \in [L]} X_{ij}^2(t) \left(\frac{1}{\sum_{j \in [L]} V_{ij}^2(t)} + 2 \frac{1}{\frac{1}{2}g_i^2(t)} \right). \end{aligned}$$

Substituting the inequalities (21) and (22),

$$\frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 \geq 6 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \left(\frac{\|\mathbf{X}_{i:}(t)\|_2^2}{\epsilon + \|\mathbf{X}_{i:}(t)\|_2^{2/3}} \right). \quad (23)$$

When $\langle \lambda_i, \hat{\mathbf{x}}_i \rangle < 0$, the direction of inequality in (20) and (23) reverses resulting in the second part of the lemma. \square

Now, in the following lemma we look at a special case when $\epsilon = 0$

Lemma B.4 (Lemma 5.3 restated). *Assume that the system evolves under continuous gradient flow, as described in Lemma B.3, with perfect row balancedness (i.e., $\frac{1}{2}g_i^2(t) = \sum_j V_{ij}^2(t) \forall i \in [N]$ and time $t \geq 0$). The rate of change of the Euclidean norm of the i -th row of $\mathbf{X}(t)$ is characterized by:*

$$\frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 = 2^{2/3} \cdot 6 \langle \lambda_i(t), \hat{\mathbf{x}}_i(t) \rangle \|\mathbf{X}_{i:}(t)\|_2^{4/3}$$

Proof.

Rewriting (19),

$$\frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 = \frac{2}{\|\mathbf{X}_{i:}(t)\|_2} (\mathbf{\Lambda}(t) \mathbf{X}^\top(t))_{ii} \left(g_i^4(t) + 4 \sum_{j \in [L]} V_{ij}^2(t) g_i^2(t) \right).$$

Substituting perfect-row balancedness condition $\frac{1}{2}g_i^2(t) = \sum_{j \in [L]} V_{ij}(t)$:

$$\frac{d}{dt} \|\mathbf{X}_{i:}(t)\|_2 = \frac{6}{\|\mathbf{X}_{i:}(t)\|_2} (\mathbf{\Lambda}(t) \mathbf{X}^\top(t))_{ii} g_i^4(t) \stackrel{a}{=} 2^{2/3} 6 \langle \boldsymbol{\lambda}_i(t), \hat{\mathbf{x}}_i(t) \rangle \|\mathbf{X}_{i:}(t)\|_2^{4/3},$$

where equality (a) follows by substituting $g_i^4(t) = 2^{2/3} \|\mathbf{X}_{i:}(t)\|_2^{4/3}$ since $\|\mathbf{X}_{i:}(t)\|_2^2 = \frac{1}{2}g_i^6(t)$, and recalling $\langle \boldsymbol{\lambda}_i(t), \hat{\mathbf{x}}_i(t) \rangle = \frac{(\mathbf{\Lambda}(t) \mathbf{X}^\top(t))_{ii}}{\|\mathbf{X}_{i:}(t)\|_2}$.

□

Lemma B.5. *For the system evolving under continuous gradient flow, the loss function $\mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) = \|\mathbf{Y} - \mathbf{A}(\mathbf{g}(t)^{\odot 2} \mathbf{1}_L \odot \mathbf{V}(t))\|_F^2$ is locally Lipschitz continuous with respect to each component $g_n(t)$ (for all $n \in [N]$) and $V_{nl}(t)$ (for all $n \in [N], l \in [L]$) over a domain of interest $\mathcal{D} := \{(g_n(t), V_{nl}(t)) \in \mathbb{R}^N \times \mathbb{R}^{N \times L} : |g_n(t)| \leq B_g, |V_{nl}(t)| \leq B_V \forall n \in [N], l \in [L]\}$ for some constants $B_g > 0$ and $B_V > 0$. Furthermore, $\mathcal{L}(g(t), V(t))$ is uniformly continuous with respect to time t .*

Proof. We begin by proving the local Lipschitz continuity of the loss function. This requires bounding the term $\left| \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) - \mathcal{L}(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)) \right|$ by the distance between the parameters $(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$ and $(\mathbf{g}(t), \mathbf{V}(t))$. Let the distance between the parameters at time t be defined as

$$\left\| (\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)) - (\mathbf{g}(t), \mathbf{V}(t)) \right\|^2 = \sum_{j \in [L]} \sum_{k \in [N]} \left(\left(\tilde{V}_{kj}(t) - V_{kj}(t) \right)^2 + (\tilde{g}_k(t) - g_k(t))^2 \right). \quad (24)$$

For brevity, we omit the explicit time dependency (t) for variables \mathbf{g} , \mathbf{V} , \mathbf{X} , and $\mathbf{\Lambda}$ unless otherwise specified.

We now expand term $\left| \mathcal{L}(\mathbf{g}, \mathbf{V}) - \mathcal{L}(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \right|$:

$$\begin{aligned} \left| \mathcal{L}(\mathbf{g}, \mathbf{V}) - \mathcal{L}(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \right| &= \left| \sum_{i \in [M]} \sum_{j \in [L]} \left(Y_{ij} - \sum_{k \in [N]} A_{ik} g_k^2 V_{kj} \right)^2 - \sum_{i \in [M]} \sum_{j \in [L]} \left(Y_{ij} - \sum_{k \in [N]} A_{ik} \tilde{g}_k^2 \tilde{V}_{kj} \right)^2 \right| \\ &\leq \sum_{i \in [M]} \sum_{j \in [L]} \left| \left(\sum_{k \in [N]} A_{ik} (\tilde{g}_k^2 \tilde{V}_{kj} - g_k^2 V_{kj}) \right) \left(2Y_{ij} - \sum_{k \in [N]} A_{ik} (g_k^2 V_{kj} + \tilde{g}_k^2 \tilde{V}_{kj}) \right) \right| \\ &\leq 2C_1 \sum_{i \in [M]} \sum_{j \in [L]} \left| \sum_{k \in [N]} (\tilde{g}_k^2 \tilde{V}_{kj} - g_k^2 V_{kj}) \right|, \end{aligned} \quad (25)$$

where $C_1 = B_Y + N\mu B_g^2 B_V$. Last line follows by noting that $g_i, V_{ij} \in \mathcal{D}$, and that \mathbf{A} is μ -coherent with ℓ_2 -normalized columns i.e., $\left| \sum_{b \in [M]} A_{bn} A_{bc} \right| \leq \mu \leq 1$. Let $B_Y = \max_{i,j} \{Y_{ij}\}$:

$$\begin{aligned} \left| \tilde{g}_k^2 \tilde{V}_{kj} - g_k^2 V_{kj} \right| &= \left| \tilde{g}_k^2 \tilde{V}_{kj} - \tilde{g}_k^2 V_{kj} + \tilde{g}_k^2 V_{kj} - g_k^2 V_{kj} \right| \\ &= \left| \tilde{g}_k^2 (\tilde{V}_{kj} - V_{kj}) + V_{kj} (\tilde{g}_k^2 - g_k^2) \right|. \end{aligned}$$

Substituting the parameter bounds from the domain \mathcal{D}

$$\begin{aligned} &\leq B_g^2 |\tilde{V}_{kj} - V_{kj}| + 2B_V B_g |\tilde{g}_k - g_k| \\ &\leq C_2 \left(|\tilde{V}_{kj} - V_{kj}| + |\tilde{g}_k - g_k| \right). \end{aligned} \quad (26)$$

where $C_2 = \max \{1, B_g^2, 2B_V B_g\}$. Substituting (26) back into (25), we can finally write:

$$\left| \mathcal{L}(\mathbf{g}, \mathbf{V}) - \mathcal{L}(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \right| \leq 2C_1 C_2 M \sum_{j \in [L]} \sum_{k \in [N]} \left(|\tilde{V}_{kj} - V_{kj}| + |\tilde{g}_k - g_k| \right).$$

Using the norm equivalence $\sum_{i=1}^M \sum_{j=1}^N |a_{ij}| \leq \sqrt{MN \sum_{i=1}^M \sum_{j=1}^N |a_{ij}|^2}$:

$$\begin{aligned} \left| \mathcal{L}(\mathbf{g}, \mathbf{V}) - \mathcal{L}(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \right| &\leq 2C_1 C_2 M \sqrt{NL} \sqrt{\sum_{j \in [L]} \sum_{k \in [N]} \left((\tilde{V}_{kj} - V_{kj})^2 + (\tilde{g}_k - g_k)^2 \right)} \\ &\leq 2C_1 C_2 M \sqrt{NL} \left| (\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) - (\mathbf{g}, \mathbf{V}) \right|. \end{aligned}$$

Thus, $\mathcal{L}(\mathbf{g}, \mathbf{V})$ is locally Lipchitz continuous with respect to $g_l \forall l \in [N]$ and $V_{lm} \forall l \in [N], m \in [L]$. Now, we proceed to prove the uniform continuity of $\mathcal{L}(\mathbf{g}(t), \mathbf{V}(t))$ with respect to t .

$$\frac{d}{dt} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) = \sum_{n \in [N]} \nabla_{g_n} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) \frac{d}{dt} g_n(t) + \sum_{n \in [N]} \sum_{l \in [L]} \nabla_{V_{nl}} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) \frac{d}{dt} V_{nl}(t).$$

By substituting the gradient flow relations $\frac{d}{dt} g_i(t) = -\nabla_{g_i(t)} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t))$ and $\frac{d}{dt} V_{il}(t) = -\nabla_{V_{il}(t)} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t))$, this becomes:

$$\frac{d}{dt} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) = - \sum_{n \in [N]} (\nabla_{g_n} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)))^2 - \sum_{n \in [N]} \sum_{l \in [L]} (\nabla_{V_{nl}} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)))^2.$$

Substituting equation (7) and (8) in the above equation, we have:

$$\frac{d}{dt} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) = \sum_{n \in [N]} \left(\left(4g_n(t) \sum_{l \in [L]} \Lambda_{nl}(t) V_{nl}(t) \right)^2 + \sum_{l \in [L]} 4g_n^4(t) \Lambda_{nl}^2(t) \right). \quad (27)$$

To complete the proof, we must show that the magnitude of this time derivative is bounded by a constant. First, we establish an upper bound for $|\Lambda_{nl}(t)|$:

$$\begin{aligned} |\Lambda_{nl}(t)| &= |(\mathbf{A}^\top (\mathbf{Y} - \mathbf{A}\mathbf{X}(t)))_{nl}| \\ &= \left| \sum_{b \in [M]} \left(A_{bn} \left(Y_{bl} - \sum_{c \in [N]} A_{bc} X_{cl}(t) \right) \right) \right| \\ &\leq \left| \sum_{b \in [M]} A_{bn} Y_{bl} \right| + \left| \sum_{b \in [M]} A_{bn} \sum_{c \in [N]} A_{bc} X_{cl}(t) \right|. \end{aligned}$$

Using \mathbf{A} is ℓ_2 -normalized columns and $X_{cl}(t) = g_c^2(t) V_{cl}(t)$, we have:

$$|\Lambda_{nl}| \leq \left| \sum_{b \in [M]} A_{bn} Y_{bl} \right| + \sum_{c \in [N]} \left| \sum_{b \in [M]} A_{bn} A_{bc} \right| \left| g_c^2(t) V_{cl}(t) \right|.$$

Substituting the parameter bounds from the domain of interest \mathcal{D} in the above equation, we have:

$$\begin{aligned} |\Lambda_{nl}| &\leq MB_Y + N\mu B_g^2 B_V \\ &\leq M (B_Y + N\mu B_g^2 B_V). \end{aligned}$$

Substituting this bound along with the parameter bounds form \mathcal{D} in (27), we obtain:

$$\begin{aligned} \left| \frac{d}{dt} \mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) \right| &\leq \sum_n (16L^2 B_g^2 B_V^2 M^2 C_1^2 + 4LB_g^4 M^2 C_1^2) \\ &\leq 4NM^2 LC_2^2 C_1^2 (L+1). \end{aligned}$$

Since $M, N, L, B_g, B_V, C_1, C_2$ are all finite constants, the entire expression is bounded by a constant. Thus, $\mathcal{L}(\cdot, \cdot)$ is uniform with respect to t . \square

Lemma B.6. *For the system evolving under continuous gradient flow, as described in Lemma B.5:*

Case I: *If $\frac{1}{\sqrt{2}}|g_i(0)| = \|V_{i:}(0)\| = 0$, then $\frac{1}{\sqrt{2}}|g_i(t)| = \|V_{i:}(t)\| = 0$ for all $t > 0$.*

Case II: *If $\frac{1}{\sqrt{2}}|g_i(0)| = \|V_{i:}(0)\| > 0$, then $\frac{1}{\sqrt{2}}|g_i(t)| = \|V_{i:}(t)\| > 0$ for all $t > 0$.*

Proof. Case I: $\frac{1}{\sqrt{2}}|g_i(0)| = \|V_{i:}(0)\| = 0$; We also have $V_{il}(0) = 0 \forall l \in [L]$.

Recall that gradient flow is related to gradient as:

$$\frac{d}{dt}g_i(t) = -\nabla_{g_i(t)}\mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) = 4g_i \sum_j \Lambda_{ij}V_{ij},$$

and

$$\frac{d}{dt}V_{il}(t) = -\nabla_{V_{il}(t)}\mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) = 2g_i^2\Lambda_{il}.$$

This system forms an initial value problem in both $g_i(t)$ and $V_{il}(t)$. Because the loss function gradient is locally Lipschitz continuous with respect to the parameters $(\mathbf{g}(t), \mathbf{V}(t))$ and uniformly continuous with respect to time t , Lemma B.5 and Teschl (2012, Theorem 2.2)¹, guarantees a unique solution for $g_i(t)$ and $V_{il}(t)$. At $t = 0$,

$$\left. \frac{d}{dt}g_i(t) \right|_{t=0} = 0; \quad \left. \frac{d}{dt}V_{il}(t) \right|_{t=0} = 0.$$

The constant functions, $g_i(t) = 0, V_{il}(t) = 0$ form a valid solution to the initial value problem. Since the solution is unique, this must be the only solution.

In our setting $g_i(t)$ and $V_{i:}(t)$ are continuous, $\frac{1}{\sqrt{2}}|g_i(0)| = \|V_{i:}(0)\| > 0$, and $\langle \boldsymbol{\lambda}_i(t), \hat{\mathbf{x}}_i(t) \rangle$, is bounded because parameters lie in bounded domain. To complete the proof, we now show that the differential equations governing these norms follow the format required by the cited lemma:

$$\frac{d}{dt}|g_i(t)| = \frac{1}{2|g_i(t)|} \frac{d}{dt}g_i^2(t); \quad \frac{d}{dt}\|V_{i:}(t)\| = \frac{1}{2\|V_{i:}\|} \frac{d}{dt}\|V_{i:}(t)\|^2.$$

Recalling the result from Lemma B.2, we have the dynamics of the norms in terms of $(\boldsymbol{\Lambda}\mathbf{X}^\top)_{ii}$. By applying the innerproduct identity (equation (16)), we have:

$$\frac{d}{dt}|g_i(t)| = \frac{4g_i^2(t)\|V_{i:}(t)\|}{|g_i(t)|} \langle \boldsymbol{\lambda}_i(t), \hat{\mathbf{x}}_i(t) \rangle. \quad (28)$$

Similarly, we have

$$\frac{d}{dt}\|V_{i:}(t)\| = 2\frac{g_i^2(t)\|V_{i:}(t)\|}{\|V_{i:}(t)\|} \langle \boldsymbol{\lambda}_i(t), \hat{\mathbf{x}}_i(t) \rangle. \quad (29)$$

Applying the perfect row-balancedness condition (namely, $\frac{1}{2}g_i^2(t) = \sum_j V_{ij}^2 = \|\mathbf{V}_{i:}(t)\|^2$) to equations (28) and (29) simplify to:

$$\frac{d}{dt}|g_i(t)| = 2\sqrt{2}g_i^2(t) \langle \boldsymbol{\lambda}_i(t), \hat{\mathbf{x}}_i(t) \rangle,$$

¹(Teschl, 2012, Theorem 2.2) (Picard-Lindelöf): Suppose $f \in C(U, \mathbb{R}^n)$, where U is an open subset of \mathbb{R}^{n+1} , and $(t_0, x_0) \in U$. If f is locally Lipschitz continuous in the second argument, uniformly with respect to the first, then there exists a unique local solution $\bar{x}(t) \in C^1(I)$ of the initial value problem, where I is some interval around t_0 .

and

$$\frac{d}{dt}\|V_{i:}(t)\| = 4\|V_{i:}(t)\|^2 \langle \boldsymbol{\lambda}_i(t), \hat{\mathbf{x}}_i(t) \rangle.$$

As all the necessary conditions of Razin et al. (2021, Lemma 5) are met, we have $g_i(t) > 0$, $\|V_{i:}(t)\| > 0$. \square

Now that we have characterized the fundamental dynamics of the gradient flow, in the next section, we develop the essential theoretical tools needed to prove our main result.

C Supporting Lemmas and the Proof of the Main Theorem

We begin by defining the initialization schemes for both the estimated trajectory $\mathbf{X}(t)$, and some reference rank- K trajectory $\tilde{\mathbf{X}}(t)$.

Initialization of estimated trajectory

At time $t = 0$, the parameters are initialized uniformly irrespective of the sparsity support \mathcal{S}_K , i.e.,

$$g_i(0) = \alpha_g; \quad V_{ij}(0) = \alpha_V \quad \forall i \in [N], j \in [L].$$

To ensure perfect row-balancedness (i.e., $\frac{1}{2}g_i^2(t) = \sum_j V_{ij}^2(t)$), at $t = 0$, we keep $\frac{1}{2}\alpha_g^2 = L\alpha_V^2$. This initialization results in the following results, at $t = 0$

$$X_{ij}(0) = g_i^2(0)V_{ij}(0) = \alpha_g^2\alpha_V = 2L\alpha_V^3 = \frac{1}{\sqrt{2L}}\alpha_g^3$$

As a result the Euclidean norm of i -th row of the estimated trajectory at time $t = 0$ is:

$$\|\mathbf{X}_{i:}(0)\|_2 = g_i^2(0)\|\mathbf{V}_{i:}(0)\|_2 = \alpha_g^2\sqrt{L}\alpha_V = \frac{1}{\sqrt{2}}\alpha_g^3 = 2L\sqrt{L}\alpha_V^3.$$

Initialization of the reference trajectory

Next, we construct a reference rank- K trajectory, denoted as $\tilde{\mathbf{X}}(t)$. For a given support set \mathcal{S}_K , we initialize its corresponding parameters so that the rows outside the support are zero:

$$\tilde{g}_i(0) = \begin{cases} \sqrt{2}\rho^{1/3} & \forall i \in \mathcal{S}_K \\ 0 & \forall i \notin \mathcal{S}_K \end{cases}; \quad \tilde{V}_{i,j}(0) = \begin{cases} \rho^{1/3} \frac{V_{ij}(0)}{\|\mathbf{V}_{i:}(0)\|} & \forall i \in \mathcal{S}_K, \forall j \in [N] \\ 0 & \forall i \notin \mathcal{S}_K \end{cases},$$

which results in

$$\|\tilde{\mathbf{V}}_{i:}(0)\|_2 = \begin{cases} \rho^{1/3} & \forall i \in \mathcal{S}_K, \\ 0 & \forall i \notin \mathcal{S}_K \end{cases}$$

Here $\rho > 0$ is a very small positive constant. This construction ensures that the reference trajectory is perfectly row-balanced (i.e., $\frac{1}{2}\tilde{g}_i^2(t) = \sum_j \tilde{V}_{ij}^2(t)$), at initialization and rank- K at initialization. Note that since $X_{ij}(t) = g_i^2(t)V_{ij}(t)$ the above initializations result in the following expression:

$$\|\tilde{\mathbf{X}}_{i:}(0)\|_2 = \begin{cases} 2\rho & \forall i \in \mathcal{S}_K, \\ 0 & \forall i \notin \mathcal{S}_K \end{cases} \quad (30)$$

Note that this construction is purely for this theoretical proof; the proposed algorithm does not assume any knowledge of the support set or sparsity.

With the initializations established, the following lemma formalizes the dynamics of the constructed reference rank- K trajectory.

Lemma C.1. *Let a K -component system be defined by parameters $\bar{\mathbf{g}}(t) \in \mathbb{R}^K$ and $\bar{\mathbf{V}}(t) \in \mathbb{R}^{K \times L}$. Let $\mathcal{S}_K \subseteq [N]$ be a support set of size $K = |\mathcal{S}_K| \leq N$. Let us construct an N -component trajectory $\tilde{\mathbf{g}}(t) \in \mathbb{R}^N$ and $\tilde{\mathbf{V}}(t) \in \mathbb{R}^{N \times L}$ by embedding the K -component system into a larger N -dimensional space and padding the remaining components with zeros. To formally relate the K -component system to its N -component embedding, we create an ordered correspondence between the indices of the small system, $\{1, \dots, K\}$, and the ordered indices of the support set, $\mathcal{S}_K = \{s_1, s_2, \dots, s_K\}$. The embedding is then defined as:*

$$\tilde{g}_i(t) := \begin{cases} \bar{g}_k(t) & \text{if } i = s_k \text{ for } k \in [K] \\ 0 & \text{if } i \notin \mathcal{S}_K \end{cases} \quad (31)$$

$$\tilde{V}_{il}(t) := \begin{cases} \bar{V}_{kl}(t) & \text{if } i = s_k \text{ for } k \in [K] \\ 0 & \text{if } i \notin \mathcal{S}_K \end{cases} \quad (32)$$

Let $\bar{\mathbf{A}}$ be the sub-matrix of \mathbf{A} such that $\bar{\mathbf{A}}$ contains only columns of \mathbf{A} that correspond to the support indices \mathcal{S}_K . Let the dynamics of $\bar{\mathbf{g}}(t)$ and $\bar{\mathbf{V}}(t)$ be derived from the loss function $\mathcal{L}_K(\bar{\mathbf{g}}(t), \bar{\mathbf{V}}(t)) = \|\mathbf{Y} - \bar{\mathbf{A}}(\bar{\mathbf{g}}(t)^{\odot 2} \mathbf{1}_L \odot \bar{\mathbf{V}}(t))\|_F^2$ where $\tilde{\mathbf{X}}(t) = (\bar{\mathbf{g}}(t)^{\odot 2} \mathbf{1}_L) \odot \bar{\mathbf{V}}(t)$ follow gradient flow path. Then, the dynamics of $\tilde{\mathbf{g}}(t)$ and $\tilde{\mathbf{V}}(t)$, derived from the loss function $\mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)) = \|\mathbf{Y} - \mathbf{A}(\tilde{\mathbf{g}}(t)^{\odot 2} \mathbf{1}_L \odot \tilde{\mathbf{V}}(t))\|_F^2$ where $\tilde{\mathbf{X}}(t) = (\tilde{\mathbf{g}}(t)^{\odot 2} \mathbf{1}_L) \odot \tilde{\mathbf{V}}(t)$, also follow the gradient flow path. That is:

$$\begin{aligned} \frac{d}{dt} \tilde{g}_i(t) &= -\nabla_{\tilde{g}_i} \mathcal{L}(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)), \\ \frac{d}{dt} \tilde{V}_{il}(t) &= -\nabla_{\tilde{V}_{il}} \mathcal{L}(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)). \end{aligned}$$

Proof. Let the system parameters $(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$ be derived from the loss function $\mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)) := \|\mathbf{Y} - \mathbf{A}(\tilde{\mathbf{g}}^{\odot 2}(t) \mathbf{1}_L \odot \tilde{\mathbf{V}}(t))\|$.

The objective of this proof is to demonstrate that the N -component trajectory $(\tilde{g}(t), \tilde{V}(t))$ (which is constructed by embedding a K -component system $(\bar{\mathbf{g}}(t), \bar{\mathbf{V}}(t))$ and padding it with zeros) follows the gradient flow dynamics derived from the N -component loss function, $\mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$.

The proof is divided into two parts, considering the components within the support set ($i \in \mathcal{S}_K$) and those outside of it ($i \notin \mathcal{S}_K$).

For the components within the support set ($i \in \mathcal{S}_K$): Let k be the index such that i is the k -th element in the ordered support set \mathcal{S}_K . Recall that the gradient flow for $\bar{g}_k(t)$ is related to the gradient as:

$$\frac{d}{dt} \bar{g}_k(t) = -\nabla_{\bar{g}_k(t)} \mathcal{L}_K(\bar{\mathbf{g}}(t), \bar{\mathbf{V}}(t)).$$

Substituting the definition of $\mathcal{L}_K(\cdot, \cdot)$:

$$\frac{d}{dt} \bar{g}_k(t) = -\nabla_{\bar{g}_k(t)} \left(\sum_m \sum_l \left((Y_{ml} - \sum_{j \in [K]} \bar{A}_{mj} \bar{g}_j^2(t) \bar{V}_{jl}(t)) \right)^2 \right).$$

Taking the derivative:

$$\frac{d}{dt} \bar{g}_k(t) = 4\bar{g}_k(t) \sum_m \sum_l \bar{A}_{mk} \left(Y_{ml} - \sum_{j \in [K]} \bar{A}_{mj} \bar{g}_j^2(t) \bar{V}_{jl}(t) \right) \bar{V}_{kl}(t),$$

By the definition of parameters, we can replace the K -component parameters $\bar{g}_k(t)$, $\bar{V}_{kl}(t)$ and the columns of $\bar{\mathbf{A}}$ by their corresponding N -component parameters $\tilde{g}_i(t)$, $\tilde{V}_{il}(t)$ and columns of \mathbf{A} .

$$\frac{d}{dt}\bar{g}_k(t) = 4\tilde{g}_i(t) \sum_m \sum_l A_{mi} \left(Y_{ml} - \sum_{j \in \mathcal{S}_K} A_{mj} \tilde{g}_j^2(t) \tilde{V}_{jl}(t) \right) \tilde{V}_{il}(t).$$

For any $j \notin \mathcal{S}_K$, we have $\tilde{g}_j(t) = 0$, $\tilde{V}_{jl}(t) = 0$, hence we can also write it as

$$\frac{d}{dt}\bar{g}_k(t) = 4\tilde{g}_i(t) \sum_m \sum_l A_{mi} \left(Y_{ml} - \sum_{j \in [N]} A_{mj} \tilde{g}_j^2(t) \tilde{V}_{jl}(t) \right) \tilde{V}_{il}(t),$$

This expression is the negative partial derivative of the loss function ($\mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$) with respect to $\tilde{g}_i(t)$ i.e.,

$$\begin{aligned} \frac{d}{dt}\bar{g}_k(t) &= -\nabla_{\tilde{g}_i(t)} \sum_m \sum_l \left(Y_{ml} - \sum_{j \in [N]} A_{mj} \tilde{g}_j^2(t) \tilde{V}_{jl}(t) \right)^2, \\ &= -\nabla_{\tilde{g}_i(t)} \mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)). \end{aligned} \quad (33)$$

By the definition (31), we have:

$$\frac{d}{dt}\tilde{g}_i(t) = \frac{d}{dt}\bar{g}_k(t). \quad (34)$$

Thus combining the results from (33) and (34), we have:

$$\frac{d}{dt}\tilde{g}_i(t) = -\nabla_{\tilde{g}_i} \mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)). \quad (35)$$

Following the same line of reasoning, we begin with the gradient flow for $\bar{V}_{kl}(t)$

$$\begin{aligned} \frac{d}{dt}\bar{V}_{kl}(t) &= -\nabla_{\bar{V}_{kl}(t)} \mathcal{L}_K(\bar{\mathbf{g}}(t), \bar{\mathbf{V}}(t)), \\ &= -\nabla_{\bar{V}_{kl}} \left(\sum_m \sum_l \left(Y_{ml} - \sum_{j \in [K]} \bar{A}_{mj} \bar{g}_j^2(t) \bar{V}_{jl}(t) \right)^2 \right) \\ &= 2\bar{g}_i^2(t) \sum_m \bar{A}_{mk} \left(Y_{ml} - \sum_{j \in [K]} \bar{A}_{mj} \bar{g}_j^2(t) \bar{V}_{jl}(t) \right) \\ &= 2\tilde{g}_i^2(t) \sum_m A_{mi} \left(Y_{ml} - \sum_{j \in \mathcal{S}_K} A_{mj} \tilde{g}_j^2(t) \tilde{V}_{jl}(t) \right). \end{aligned}$$

The last line is obtained by replacing the K -component parameters by their N -component counterparts.

For any $j \notin \mathcal{S}_K$, we have $\tilde{g}_j(t) = 0$, $\tilde{V}_{jl}(t) = 0$, hence we can also write the above as

$$\frac{d}{dt}\bar{V}_{kl}(t) = 2\tilde{g}_i^2(t) \sum_m A_{mi} \left(Y_{ml} - \sum_{j \in [N]} A_{mj} \tilde{g}_j^2(t) \tilde{V}_{jl}(t) \right),$$

This expression is the negative partial derivative of the loss function ($\mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$) with respect to $\tilde{V}_{il}(t)$ i.e.,

$$\begin{aligned} \frac{d}{dt}\bar{V}_{kl}(t) &= -\nabla_{\tilde{V}_{il}(t)} \sum_m \sum_l \left(Y_{ml} - \sum_{j \in [N]} A_{mj} \tilde{g}_j^2(t) \tilde{V}_{jl}(t) \right)^2, \\ &= -\nabla_{\tilde{V}_{il}} \mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)). \end{aligned} \quad (36)$$

By the definition (31), we have

$$\frac{d}{dt} \tilde{V}_{il}(t) = \frac{d}{dt} \bar{V}_{kl}(t). \quad (37)$$

Combining the results from (36) and (37) we have:

$$\frac{d}{dt} \tilde{V}_{kl}(t) = -\nabla_{\tilde{V}_{il}(t)} \mathcal{L}_N(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}). \quad (38)$$

For the components outside the support set ($i \notin \mathcal{S}_K$):

By definitions (31) and (32), $\tilde{g}_i(t) = 0$ and $\tilde{V}_{il}(t) = 0$ for all $t \geq 0$. Also, note that $\nabla_{\tilde{V}_{il}(t)} \mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)) \propto \tilde{g}_i^2(t)$ and $\nabla_{\tilde{g}_i(t)} \mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)) \propto \tilde{g}_i(t)$. Combining we arrive at:

$$\frac{d}{dt} \tilde{g}_i(t) = 0 = -\nabla_{\tilde{g}_i(t)} \mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)), \quad (39)$$

and

$$\frac{d}{dt} \tilde{V}_{il}(t) = 0 = -\nabla_{\tilde{V}_{il}(t)} \mathcal{L}_N(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)). \quad (40)$$

Thus, from equations (35), (38), (39), and (40) we have that $\tilde{\mathbf{g}}(t)$ and $\tilde{\mathbf{V}}(t)$ both follow gradient flow path on N component. \square

Lemma C.2. *Let the domain of interest be $\mathcal{D} := \{(\mathbf{g}(t), \mathbf{V}(t)) : |g_n(t)| \leq B_g, |V_{nl}(t)| \leq B_V \forall n \in [N], l \in [L]\}$ for some constants $B_g > 0$ and $B_V > 0$. For the system evolving under continuous gradient flow, the loss function $\mathcal{L}(\mathbf{g}(t), \mathbf{V}(t)) = \|\mathbf{Y} - \mathbf{A}(\mathbf{g}(t)^{\odot 2} \mathbf{1}_L \odot \mathbf{V}(t))\|_F^2$ is β -smooth over \mathcal{D} , where \mathbf{A} is a μ -coherent matrix with l_2 -normalized columns. The Lipschitz smoothness constant β is given by:*

$$\beta = 16NL^{3/2} ((N+1)\mu C^4 + MB_Y C),$$

where $C \geq \max\{1, B_g, 2B_V\}$ and $B_Y = \max_{n \in [N], l \in [L]} \{Y_{nl}\}$.

Proof. For brevity, we omit the explicit time dependency (t) for variables $\mathbf{g}, \mathbf{V}, \mathbf{A}, \mathbf{X}; g_n, V_{nl}, \Lambda_{nl}$ and X_{nl} correspond to the respective value at time t unless specified otherwise.

To establish that $\mathcal{L}(\mathbf{g}, \mathbf{V})$ is β -smooth, we must show that its gradient is Lipschitz continuous. Let (\mathbf{g}, \mathbf{V}) and $(\tilde{\mathbf{g}}, \tilde{\mathbf{V}})$ be two distinct points within the domain of interest \mathcal{D} corresponding to the signal matrices $\mathbf{X} = \mathbf{g}^{\odot 2} \mathbf{1}_L \odot \mathbf{V}$ and $\tilde{\mathbf{X}} = \tilde{\mathbf{g}}^{\odot 2} \mathbf{1}_L \odot \tilde{\mathbf{V}}$. Note that for all $i \in [N]$ and $j \in [L]$, $X_{ij} = g_i^2 V_{ij}$, $\tilde{X}_{ij} = \tilde{g}_i^2 \tilde{V}_{ij}$. We will now analyze the squared distance between the gradients at these two points:

$$\begin{aligned} E_{\mathcal{L}}^2 &= \|\nabla \mathcal{L}(\mathbf{g}, \mathbf{V}) - \nabla \mathcal{L}(\tilde{\mathbf{g}}, \tilde{\mathbf{V}})\|^2 \\ &= \sum_{n \in [N]} \left| \nabla_{g_n} \mathcal{L}(\mathbf{g}, \mathbf{V}) - \nabla_{\tilde{g}_n} \mathcal{L}(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \right|^2 \\ &\quad + \sum_{l \in [L]} \sum_{n \in [N]} \left| \nabla_{V_{nl}} \mathcal{L}(\mathbf{g}, \mathbf{V}) - \nabla_{\tilde{V}_{nl}} \mathcal{L}(\tilde{\mathbf{g}}, \tilde{\mathbf{V}}) \right|^2. \end{aligned}$$

Substituting the gradient expressions from (7) and (8) into the above equation yields:

$$E_{\mathcal{L}}^2 = \sum_{n \in [N]} \left| 4g_n \sum_{l \in [L]} \Lambda_{nl} V_{nl} - 4\tilde{g}_n \sum_{l \in [L]} \tilde{\Lambda}_{nl} \tilde{V}_{nl} \right|^2 + \sum_{l \in [L]} \sum_{n \in [N]} \left| 2g_n^2 \Lambda_{nl} - 2\tilde{g}_n^2 \tilde{\Lambda}_{nl} \right|^2.$$

Here $\tilde{\Lambda} = \mathbf{A}^\top (\mathbf{Y} - \mathbf{A}(\tilde{\mathbf{g}}^{\odot 2} \mathbf{1}_L \odot \tilde{\mathbf{V}}))$. Hence,

$$\begin{aligned} E_{\mathcal{L}}^2 &\leq L \sum_{n \in [N]} \sum_{l \in [L]} \left| 4g_n \Lambda_{nl} V_{nl} - 4\tilde{g}_n \tilde{\Lambda}_{nl} \tilde{V}_{nl} \right|^2 + \sum_{l \in [L]} \sum_{n \in [N]} \left| 2g_n^2 \Lambda_{nl} - 2\tilde{g}_n^2 \tilde{\Lambda}_{nl} \right|^2, \\ &\leq L \sum_{n \in [N]} \sum_{l \in [L]} \left| 4g_n \Lambda_{nl} V_{nl} - 4\tilde{g}_n \tilde{\Lambda}_{nl} \tilde{V}_{nl} \right|^2 + L \sum_{l \in [L]} \sum_{n \in [N]} \left| 2g_n^2 \Lambda_{nl} - 2\tilde{g}_n^2 \tilde{\Lambda}_{nl} \right|^2. \end{aligned} \quad (41)$$

Now, we try to bound the terms in each sum.

Substituting $X_{cl}(t) = g_c^2(t) V_{cl}(t)$ and using parameter bounds of domain of interest \mathcal{D} in above equation yields

$$\begin{aligned} \left| \tilde{X}_{cl} - X_{cl} \right| &= \left| \tilde{g}_c^2 \tilde{V}_{cl} - g_c^2 V_{cl} \right| \\ &= \left| \tilde{g}_c^2 \tilde{V}_{cl} - g_c^2 \tilde{V}_{cl} + g_c^2 \tilde{V}_{cl} - g_c^2 V_{cl} \right| \\ &= \left| (\tilde{g}_c - g_c)(\tilde{g}_c + g_c) \tilde{V}_{cl} + g_c^2 (\tilde{V}_{cl} - V_{cl}) \right| \end{aligned}$$

Substituting the parameter bounds from the domain \mathcal{D} ,

$$\begin{aligned} \left| \tilde{X}_{cl} - X_{cl} \right| &\leq 2B_g B_V |\tilde{g}_c - g_c| + B_g^2 |\tilde{V}_{cl} - V_{cl}| \\ &\leq C^2 \left(|\tilde{g}_c - g_c| + |\tilde{V}_{cl} - V_{cl}| \right), \end{aligned} \quad (42)$$

where, $C \geq \max\{1, B_g, 2B_V\}$.

$$\begin{aligned} |\Lambda_{nl} - \tilde{\Lambda}_{nl}| &= \left| (\mathbf{A}^\top (\mathbf{Y} - \mathbf{A}\mathbf{X}))_{nl} - (\mathbf{A}^\top (\mathbf{Y} - \mathbf{A}\tilde{\mathbf{X}}))_{nl} \right| \\ &= \left| \sum_{b \in [M]} \left(A_{bn} \left(Y_{bl} - \sum_{c \in [N]} A_{bc} X_{cl} \right) \right) - \sum_{b \in [M]} \left(A_{bn} \left(Y_{bl} - \sum_{c \in [N]} A_{bc} \tilde{X}_{cl} \right) \right) \right| \\ &= \left| \sum_{b \in [M]} \left(A_{bn} \sum_{c \in [N]} A_{bc} \tilde{X}_{cl} \right) - \sum_{b \in [M]} \left(A_{bn} \sum_{c \in [N]} A_{bc} X_{cl} \right) \right| \\ &\leq \sum_{c \in [N]} \left(\left| \sum_{b \in [M]} A_{bn} A_{bc} \right| \left| \tilde{X}_{cl} - X_{cl} \right| \right). \end{aligned}$$

Since we have \mathbf{A} is μ -coherent with ℓ_2 -normalized columns, we have $\left| \sum_{b \in [M]} A_{bn} A_{bc} \right| \leq \mu \leq 1$. Substituting, we have:

$$|\Lambda_{nl} - \tilde{\Lambda}_{nl}| \leq \mu \sum_{c \in [N]} \left| \tilde{X}_{cl} - X_{cl} \right|.$$

Substituting (42):

$$|\Lambda_{nl} - \tilde{\Lambda}_{nl}| \leq \mu C^2 \sum_{c \in [N]} \left(|\tilde{g}_c - g_c| + |\tilde{V}_{cl} - V_{cl}| \right).$$

We can use the following looser inequality:

$$|\Lambda_{nl} - \tilde{\Lambda}_{nl}| \leq \mu C^2 \sum_{l \in [L]} \sum_{c \in [N]} \left(|\tilde{g}_c - g_c| + |\tilde{V}_{cl} - V_{cl}| \right).$$

Let $S = \sum_{c \in [N]} \sum_{l \in [L]} \left(|\tilde{g}_c - g_c| + |\tilde{V}_{cl} - V_{cl}| \right)$. The difference is then given by:

$$|\Lambda_{nl} - \tilde{\Lambda}_{nl}| \leq \mu C^2 S. \quad (43)$$

Similarly, we bound the term $|g_n V_{nl} - \tilde{g}_n \tilde{V}_{nl}|$ as:

$$|g_n V_{nl} - \tilde{g}_n \tilde{V}_{nl}| \leq |g_n - \tilde{g}_n| |\tilde{V}_{nl}| + |V_{nl} - \tilde{V}_{nl}| |g_n|.$$

By applying the parameter bounds from the domain \mathcal{D} :

$$\begin{aligned} |g_n V_{nl} - \tilde{g}_n \tilde{V}_{nl}| &\leq B_V |g_n - \tilde{g}_n| + B_g |V_{nl} - \tilde{V}_{nl}| \\ &\leq \max \{1, B_g, B_V\} \left(|\tilde{g}_n - g_n| + |\tilde{V}_{nl} - V_{nl}| \right) \\ &\leq C \left(|\tilde{g}_n - g_n| + |\tilde{V}_{nl} - V_{nl}| \right), \end{aligned}$$

where $C = \max \{1, B_g, 2B_V\}$. Since the non-negative term, $|\tilde{g}_n - g_n| + |\tilde{V}_{nl} - V_{nl}|$, is smaller than or equal to the sum over all terms, we use the following looser upper bound for simplicity:

$$\begin{aligned} |g_n V_{nl} - \tilde{g}_n \tilde{V}_{nl}| &\leq C \sum_{l \in [L]} \sum_{c \in [N]} \left(|\tilde{g}_c - g_c| + |\tilde{V}_{cl} - V_{cl}| \right) \\ |g_n V_{nl} - \tilde{g}_n \tilde{V}_{nl}| &\leq CS. \end{aligned} \quad (44)$$

Next, we find an upper bound for the magnitude of the term $|\tilde{\Lambda}_{nl}|$.

$$\begin{aligned} |\tilde{\Lambda}_{nl}| &\leq \left| \left(\mathbf{A}^\top (\mathbf{Y} - \mathbf{A} \tilde{\mathbf{X}}) \right)_{nl} \right| \\ &\leq \left| \sum_{b \in [M]} \left(A_{bn} \left(Y_{bl} - \sum_{c \in [N]} A_{bc} \tilde{X}_{cl} \right) \right) \right|. \end{aligned}$$

By applying the triangle inequality again and substituting $\tilde{X}_{cl} = \tilde{g}_c^2 \tilde{V}_{cl}$:

$$\leq \left| \sum_{b \in [M]} A_{bn} Y_{bl} \right| + \left| \sum_{b \in [M]} A_{bn} \sum_{c \in [N]} A_{bc} \tilde{X}_{cl} \right|.$$

Finally, applying the parameter bounds from the domain \mathcal{D} , and the fact that \mathbf{A} is μ -coherent with ℓ_2 -normalized columns ($|A_{ij}| \leq 1$):

$$|\tilde{\Lambda}_{nl}| \leq \left| \sum_{b \in [M]} A_{bn} Y_{bl} \right| + \sum_{c \in [N]} \left| \sum_{b \in [M]} A_{bn} A_{bc} \right| \left| \tilde{g}_c^2 \tilde{V}_{cl} \right|.$$

Using $\left| \tilde{g}_c^2 \tilde{V}_{cl} \right| \leq B_g^2 B_V \leq C^3$, and $Y_{nl} \leq B_Y$, the above equation can be written as

$$|\tilde{\Lambda}_{nl}| \leq M B_Y + N \mu C^3. \quad (45)$$

$$\begin{aligned} \left| 4g_n \Lambda_{nl} V_{nl} - 4\tilde{g}_n \tilde{\Lambda}_{nl} \tilde{V}_{nl} \right| &= 4|(\Lambda_{nl} - \tilde{\Lambda}_{nl})g_n V_{nl} + (g_n V_{nl} - \tilde{g}_n \tilde{V}_{nl})\tilde{\Lambda}_{nl}| \\ &\leq 4|(\Lambda_{nl} - \tilde{\Lambda}_{nl})g_n V_{nl} + (g_n V_{nl} - \tilde{g}_n \tilde{V}_{nl})\tilde{\Lambda}_{nl}| \\ &\leq 4|\Lambda_{nl} - \tilde{\Lambda}_{nl}| |g_n V_{nl}| + 4|g_n V_{nl} - \tilde{g}_n \tilde{V}_{nl}| |\tilde{\Lambda}_{nl}|. \end{aligned}$$

Substituting $\left| \tilde{g}_c \tilde{V}_{cl} \right| \leq B_g B_V \leq C^2$, (43) and (44) in the previous equation, we have:

$$\begin{aligned} \left| 4g_n \Lambda_{nl} V_{nl} - 4\tilde{g}_n \tilde{\Lambda}_{nl} \tilde{V}_{nl} \right| &\leq 4\mu C^4 S + 4(MB_Y + N\mu C^3)CS \\ &\leq 4((N+1)\mu C^4 + MB_Y C) S. \end{aligned} \quad (46)$$

Next, we find an upper bound for the magnitude of the term $\left| 2g_n^2 \Lambda_{nl} - 2\tilde{g}_n^2 \tilde{\Lambda}_{nl} \right|$.

$$\begin{aligned} \left| 2g_n^2 \Lambda_{nl} - 2\tilde{g}_n^2 \tilde{\Lambda}_{nl} \right| &\leq 2|g_n^2 - \tilde{g}_n^2| \left| \tilde{\Lambda}_{nl} \right| + 2|g_n^2| \left| \Lambda_{nl} - \tilde{\Lambda}_{nl} \right| \\ &\leq 2|g_n - \tilde{g}_n| |g_n + \tilde{g}_n| \left| \tilde{\Lambda}_{nl} \right| + 2|g_n^2| \left| \Lambda_{nl} - \tilde{\Lambda}_{nl} \right| \\ &\leq 4B_g |g_n - \tilde{g}_n| \left| \tilde{\Lambda}_{nl} \right| + 2B_g^2 \left| \Lambda_{nl} - \tilde{\Lambda}_{nl} \right| \end{aligned}$$

Substituting (43) and (45) in the above equation yields:

$$\left| 2g_n^2 \Lambda_{nl} - 2\tilde{g}_n^2 \tilde{\Lambda}_{nl} \right| \leq 4B_g (MB_Y + N\mu C^3) |g_n - \tilde{g}_n| + 2B_g^2 \mu C^2 S.$$

Note that $|g_n - \tilde{g}_n| \leq \sum_{l \in [L]} \sum_{c \in [N]} \left(|\tilde{g}_c - g_c| + |\tilde{V}_{cl} - V_{cl}| \right) = S$ and $B_g \leq C$. Substituting this in the previous equation, we have:

$$\begin{aligned} \left| 2g_n^2 \Lambda_{nl} - 2\tilde{g}_n^2 \tilde{\Lambda}_{nl} \right| &\leq 4B_g (MB_Y + N\mu C^3) S + 2B_g^2 \mu C^2 S \\ &\leq 4(C(MB_Y + N\mu C^3) + \mu C^4) S \\ &\leq 4((N+1)\mu C^4 + MB_Y C) S. \end{aligned} \quad (47)$$

Substituting the derived upper bounds ((46) and (47)) back into (41), we have:

$$\begin{aligned} E_{\mathcal{L}}^2 &\leq 2L \sum_{n' \in [N]} \sum_{l' \in [L]} \left(4((N+1)\mu C^4 + MB_Y C) \right)^2 S^2 \\ &\leq 2NL^2 \left(4((N+1)\mu C^4 + MB_Y C) \right)^2 S^2 \\ &\leq 2NL^2 \left(4((N+1)\mu C^4 + MB_Y C) \right)^2 \left(\sum_{l \in [L]} \sum_{c \in [N]} \left(|\tilde{g}_c - g_c| + |\tilde{V}_{cl} - V_{cl}| \right) \right)^2. \end{aligned}$$

Since $\sum_{i=1}^M \sum_{j=1}^N |a_{ij}| \leq \sqrt{MN \sum_{i=1}^M \sum_{j=1}^N |a_{ij}|^2}$:

$$\begin{aligned} E_{\mathcal{L}}^2 &\leq 4N^2 L^3 \left(4((N+1)\mu C^4 + MB_Y C) \right)^2 \sum_{l \in [L]} \sum_{c \in [N]} \left(|\tilde{g}_c - g_c| + |\tilde{V}_{cl} - V_{cl}| \right)^2 \\ E_{\mathcal{L}} &\leq 8NL^{3/2} ((N+1)\mu C^4 + MB_Y C) \sqrt{\sum_{l \in [L]} \sum_{c \in [N]} \left(|\tilde{g}_c - g_c| + |\tilde{V}_{cl} - V_{cl}| \right)^2} \\ E_{\mathcal{L}} &\leq 16NL^{3/2} ((N+1)\mu C^4 + MB_Y C) \sqrt{\sum_{l \in [L]} \sum_{c \in [N]} \left(|\tilde{g}_c - g_c|^2 + |\tilde{V}_{cl} - V_{cl}|^2 \right)}. \end{aligned}$$

Thus, the loss function $\mathcal{L}(\mathbf{g}, \mathbf{V})$ is β -smooth over \mathcal{D} where $\beta = 16NL^{3/2} ((N+1)\mu C^4 + MB_Y C)$ for $g_n, V_{nl} \in \mathcal{D} \forall n \in [N], l \in [L]$. \square

Lemma C.3 (Lemma 5.4 restated). *Let the estimated trajectory, $\mathbf{X}(t)$ be parametrized by $(\mathbf{g}(t), \mathbf{V}(t))$ and a reference trajectory $\tilde{\mathbf{X}}(t)$, be parameterized by $(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$ both of which evolve under continuous gradient flow, as described in Lemma C.2. Let the parameters $(\mathbf{g}(t), \mathbf{V}(t))$ and $(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$ be at a distance D from the origin. At any time $t \in [0, T]$, the following inequality holds:*

$$\left\| \mathbf{X}(t) - \tilde{\mathbf{X}}(t) \right\|_F^2 \leq 8D^4 \sum_{n \in [N]} \sum_{l \in [L]} \left(\left(V_{nl}(t) - \tilde{V}_{nl}(t) \right)^2 + (g_n(t) - \tilde{g}_n(t))^2 \right).$$

Proof.

$$\begin{aligned} \left\| \mathbf{X}(t) - \tilde{\mathbf{X}}(t) \right\|_F^2 &= \sum_{n \in [N]} \left\| g_n^2(t) \mathbf{V}_{n:}(t) - \tilde{g}_n^2(t) \tilde{\mathbf{V}}_{n:}(t) \right\|_2^2 \\ &= \sum_{n \in [N]} \left\| g_n^2(t) \mathbf{V}_{n:}(t) + \tilde{g}_n^2(t) \tilde{\mathbf{V}}_{n:}(t) - g_n^2(t) \tilde{\mathbf{V}}_{n:}(t) - \tilde{g}_n^2(t) \mathbf{V}_{n:}(t) \right\|_2^2 \\ &= \sum_{n \in [N]} \left\| g_n^2(t) (\mathbf{V}_{n:}(t) - \tilde{\mathbf{V}}_{n:}(t)) + \tilde{\mathbf{V}}_{n:}(t) (g_n^2(t) - \tilde{g}_n^2(t)) \right\|_2^2 \\ &\leq \sum_{n \in [N]} \left(g_n^2(t) \left\| \mathbf{V}_{n:}(t) - \tilde{\mathbf{V}}_{n:}(t) \right\|_2 + \left\| \tilde{\mathbf{V}}_{n:}(t) \right\|_2 |g_n^2(t) - \tilde{g}_n^2(t)| \right)^2, \end{aligned}$$

using $(a + b)^2 \leq 2a^2 + 2b^2$,

$$\begin{aligned} \left\| \mathbf{X}(t) - \tilde{\mathbf{X}}(t) \right\|_F^2 &\leq 2 \sum_{n \in [N]} \left((g_n^2(t))^2 \left\| \mathbf{V}_{n:}(t) - \tilde{\mathbf{V}}_{n:}(t) \right\|_2^2 + \left\| \tilde{\mathbf{V}}_{n:}(t) \right\|_2^2 (g_n^2(t) - \tilde{g}_n^2(t))^2 \right) \\ &\leq 2 \sum_{n \in [N]} \left(g_n^4(t) \sum_{l \in [L]} \left(V_{nl}(t) - \tilde{V}_{nl}(t) \right)^2 + \sum_{l \in [L]} \tilde{V}_{nl}(t)^2 ((g_n(t) - \tilde{g}_n(t))(g_n(t) + \tilde{g}_n(t)))^2 \right) \end{aligned}$$

Using the inequalities $\{|g_n(t)|, \|\mathbf{V}_{n:}(t)\|_2\} \leq \|(\mathbf{g}(t), \mathbf{V}(t))\| \leq D$ and $\{|\tilde{g}_n(t)|, \|\tilde{\mathbf{V}}_{n:}(t)\|_2\} \leq \|(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))\| \leq D$, we can write the above equation as:

$$\begin{aligned} \left\| \mathbf{X}(t) - \tilde{\mathbf{X}}(t) \right\|_F^2 &\leq 2 \sum_{n \in [N]} \sum_{l \in [L]} \left(D^4 \left(V_{nl}(t) - \tilde{V}_{nl}(t) \right)^2 + 4D^4 (g_n(t) - \tilde{g}_n(t))^2 \right) \\ &\leq 8D^4 \sum_{n \in [N]} \sum_{l \in [L]} \left(\left(V_{nl}(t) - \tilde{V}_{nl}(t) \right)^2 + (g_n(t) - \tilde{g}_n(t))^2 \right). \end{aligned}$$

□

C.1 Proof of Theorem V.1

To prove the theorem, our strategy is to first construct an idealized, reference low-rank trajectory, denoted as $\tilde{\mathbf{X}}(t)$, which will serve as a theoretical benchmark. This trajectory is designed to be perfectly rank- K by initializing it such that only the rows within a specific support set, \mathcal{S}_K , are non-zero. The initial magnitude of these non-zero rows is controlled by a small positive scalar, ρ .

The core of our proof will now be to show that the estimated trajectory found by the algorithm, namely, $\mathbf{X}(t)$, remains close to this reference trajectory, namely, $\tilde{\mathbf{X}}(t)$, provided the initialization is sufficiently small. It is crucial to note that estimated trajectory, $\mathbf{X}(t)$, does not assume any knowledge of the support set \mathcal{S}_K or sparsity level K .

We begin the proof using the idealized rank- K trajectory, $\tilde{\mathbf{X}}(t)$ as in (30). This initialization along with Lemmas 5.3 and B.6 ensures that the rows outside the support \mathcal{S}_K are initialized at zero and remain zero under the gradient flow dynamics. Thus, $\tilde{\mathbf{X}}(t)$ is guaranteed to be an at-most rank- K matrix for all $t \geq 0$.

To quantify the proximity between the estimated parameters $(\mathbf{g}(t), \mathbf{V}(t))$ and reference parameters $(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))$, we use the squared distance defined in (24), namely $\left\|(\mathbf{g}(t), \mathbf{V}(t)) - (\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))\right\|^2$.

Bounding the Initial Distance

We first analyze the initial distance between the parameters for any row i that is outside the support set ($i \notin \mathcal{S}_K$). For these rows, we recall that the estimated parameters are initialized as $g_i(0) = \alpha_g$ and $V_{il}(0) = \alpha_V$, while the reference parameters are $\tilde{g}_i(0) = 0$ and $\tilde{V}_{il}(0) = 0 \forall l \in [L]$.

The distance contributed by the non-support points ($i \notin \mathcal{S}_K$) is:

$$\begin{aligned} \sum_{l \in [L]} \left(|g_i(0) - \tilde{g}_i(0)|^2 + |V_{il}(0) - \tilde{V}_{il}(0)|^2 \right) \\ = L\alpha_g^2 + L\alpha_V^2. \end{aligned}$$

Using the balancedness condition at $t = 0$ $\left(\frac{1}{2}g_i^2(t) = \sum_j V_{ij}^2(t) \right)$ and $\alpha_g^2 = 2L\alpha_V^2$, this simplifies to

$$\begin{aligned} \sum_{l \in [L]} \left(|g_i(0) - \tilde{g}_i(0)|^2 + |V_{il}(0) - \tilde{V}_{il}(0)|^2 \right) \\ = (2L + 1) L\alpha_V^2. \end{aligned} \quad (48)$$

To control how this distance evolves over time, we use Razin et al. (2021, Lemma 6). This lemma requires us to set the initial distance to be bounded by a term related to the final approximation error ϵ and the smoothness constant β , which we ensure by imposing a condition on the initialization scales, which is:

$$\sum_{l \in [L]} \left(|g_i(0) - \tilde{g}_i(0)|^2 + |V_{il}(0) - \tilde{V}_{il}(0)|^2 \right) \leq \frac{1}{N} \exp(-2\beta T) \quad (49)$$

Then we have $(2L + 1) L\alpha_V^2 \leq \frac{1}{N} \exp(-2\beta T)$, rearranging this inequality gives us the upper bound on the initialization as:

$$\alpha_V \leq \frac{\exp(-\beta T)}{\sqrt{LN(2L + 1)}}. \quad (50)$$

The distance contributed by the support set ($i \in \mathcal{S}_K$) at initialization $t = 0$ is:

$$\begin{aligned} \sum_{l \in [L]} \left(|g_i(0) - \tilde{g}_i(0)|^2 + |V_{il}(0) - \tilde{V}_{il}(0)|^2 \right) &= L(\alpha_g^2 + \rho^{2/3} - 2\alpha_g\rho^{1/3} + \alpha_V^2 + \rho^{2/3} - 2\alpha_V\rho^{1/3}) \\ &= L((2L + 1)\alpha_V^2 + 2\rho^{2/3} - 2(\sqrt{2L} + 1)\alpha_V\rho^{1/3}). \end{aligned} \quad (51)$$

Let $0 < \rho \leq \alpha_V^3$. By applying the bound from (49) we have:

$$\begin{aligned} \left(L + \frac{1}{2} \right) \alpha_V^2 + \rho^{2/3} - (\sqrt{2L} + 1) \alpha_V \rho^{1/3} &\leq \frac{1}{2LN} \exp(-2\beta T) \\ \rho^{2/3} - (\sqrt{2L} + 1) \alpha_V \rho^{1/3} + \left(\sqrt{\frac{L}{2}} + \frac{1}{2} \right)^2 \alpha_V^2 &\leq \frac{1}{2LN} \exp(-2\beta T) - \left(L + \frac{1}{2} \right) \alpha_V^2 + \left(\sqrt{\frac{L}{2}} + \frac{1}{2} \right)^2 \alpha_V^2 \\ \left(\rho^{1/3} - \left(\sqrt{\frac{L}{2}} + \frac{1}{2} \right) \alpha_V \right)^2 &\leq \frac{1}{2LN} \exp(-2\beta T) - \frac{1}{4} \alpha_V^2 (\sqrt{2L} - 1)^2. \end{aligned} \quad (52)$$

For a real solution for $\rho^{1/3}$ to exist, the right-hand side of this inequality must be non-negative. Then, let α_V be designed such that $\frac{1}{2LN}\exp(-2\beta T) - \frac{1}{4}\alpha_V^2 \left(\sqrt{2L} - 1\right)^2 \geq 0$ holds. Thus,

$$\alpha_V \leq \frac{\sqrt{2}}{\sqrt{LN}(\sqrt{2L} - 1)}\exp(-\beta T).$$

To ensure that the initialization conditions for rows both inside and outside the support set are met, α_V must satisfy the stricter of the two derived upper bounds ((50) and (50)). This leads to the initialization condition:

$$\alpha_V \leq \frac{1}{\sqrt{LN}}\exp(-\beta T)\min\left\{\frac{\sqrt{2}}{\sqrt{2L} - 1}, \frac{1}{\sqrt{(2L+1)}}\right\}.$$

Since $\frac{\sqrt{2}}{\sqrt{2L}-1} > \frac{1}{\sqrt{(2L+1)}}$ for all $L \geq 1$. Note that this is still Equation (50).

Equation (52) provides a valid interval for $\rho^{1/3}$:

$$\rho^{1/3} \in \left[\left(\sqrt{\frac{L}{2}} + \frac{1}{2} \right) \alpha_V - \epsilon_\alpha, \left(\sqrt{\frac{L}{2}} + \frac{1}{2} \right) \alpha_V + \epsilon_\alpha \right], \quad (53)$$

where $\epsilon_\alpha = \left(\frac{1}{2LN}\exp(-2\beta T) - \frac{1}{4}\alpha_V^2 \left(\sqrt{2L} - 1\right)^2 \right)^{1/2}$. But we have $\left(\sqrt{\frac{L}{2}} + \frac{1}{2} \right) \alpha_V + \epsilon_\alpha \geq \alpha_V$ in order to keep $\epsilon_\alpha \geq 0$. Hence the set becomes:

$$\rho^{1/3} \in \left[\left(\sqrt{\frac{L}{2}} + \frac{1}{2} \right) \alpha_V - \epsilon_\alpha, \alpha_V \right].$$

For this set to be non-empty, we need lower bound to be less than or equal to the upper bound. This gives us:

$$\left(\sqrt{\frac{L}{2}} + \frac{1}{2} \right) \alpha_V - \epsilon_\alpha \leq \alpha_V.$$

Rearranging, we have:

$$\left(\sqrt{\frac{L}{2}} - \frac{1}{2} \right) \alpha_V \leq \epsilon_\alpha.$$

Rearranging both sides and substituting for ϵ_α :

$$\left(\sqrt{\frac{L}{2}} - \frac{1}{2} \right)^2 \alpha_V^2 \leq \frac{1}{2LN}\exp(-2\beta T) - \frac{1}{4}\alpha_V^2 \left(\sqrt{2L} - 1\right)^2.$$

Reorganizing,

$$\begin{aligned} \alpha_V^2 \left(\sqrt{2L} - 1\right)^2 &\leq \frac{1}{LN}\exp(-2\beta T), \\ \alpha_V &\leq \frac{1}{\sqrt{LN}(\sqrt{2L} - 1)}\exp(-\beta T). \end{aligned}$$

This confirms that the interval for ρ (equation (53)) is consistent with Equation (50).

With the initialization conditions for α_V and ρ established, we can bound the total initial squared distance between the parameters. A time $t = 0$:

$$\left\| (\mathbf{g}(0), \mathbf{V}(0)) - (\tilde{\mathbf{g}}(0), \tilde{\mathbf{V}}(0)) \right\|^2 = \sum_{i \in [N]} \sum_{l \in [L]} \left(|g_i(0) - \tilde{g}_i(0)|^2 + |V_{il}(0) - \tilde{V}_{il}(0)|^2 \right),$$

Splitting the contributions of support and out-side support parameters,

$$\begin{aligned} \left\| (\mathbf{g}(0), \mathbf{V}(0)) - (\tilde{\mathbf{g}}(0), \tilde{\mathbf{V}}(0)) \right\|^2 &= \sum_{i \in \mathcal{S}_K} \sum_{l \in [L]} \left(|g_i(0) - \tilde{g}_i(0)|^2 + |V_{il}(0) - \tilde{V}_{il}(0)|^2 \right) \\ &\quad + \sum_{i \notin \mathcal{S}_K} \sum_{l \in [L]} \left(|g_i(0) - \tilde{g}_i(0)|^2 + |V_{il}(0) - \tilde{V}_{il}(0)|^2 \right), \end{aligned}$$

Substituting back the equations (51) and (48),

$$\begin{aligned} \left\| (\mathbf{g}(0), \mathbf{V}(0)) - (\tilde{\mathbf{g}}(0), \tilde{\mathbf{V}}(0)) \right\|^2 &\leq KL((2L+1)\alpha_V^2 + 2\rho^{2/3} - 2(\sqrt{2L}+1)\alpha_V\rho^{1/3}) + (N-K)((2L+1)L\alpha_V^2), \\ &\leq KL(2L+1)\alpha_V^2 + 2KL\rho^{2/3} - 2KL(\sqrt{2L}+1)\alpha_V\rho^{1/3} + (N-K)((2L+1)L\alpha_V^2), \\ &\leq NL(2L+1)\alpha_V^2 + 2KL\rho^{2/3} - 2KL(\sqrt{2L}+1)\alpha_V\rho^{1/3}, \end{aligned}$$

to find a simpler upper bound, we can drop the negative term.

$$\left\| (\mathbf{g}(0), \mathbf{V}(0)) - (\tilde{\mathbf{g}}(0), \tilde{\mathbf{V}}(0)) \right\|^2 \leq NL(2L+1)\alpha_V^2 + 2KL\rho^{2/3}.$$

We have $0 < \rho \leq \alpha_V^3$, which implies $\rho^{2/3} \leq \alpha_V^2$. Hence,

$$\begin{aligned} \left\| (\mathbf{g}(0), \mathbf{V}(0)) - (\tilde{\mathbf{g}}(0), \tilde{\mathbf{V}}(0)) \right\|^2 &\leq NL(2L+1)\alpha_V^2 + 2K\alpha_V^2 \\ &\leq (N(2L+1) + 2K)L\alpha_V^2. \end{aligned}$$

Note that the substituting the upper bound on the value α_V from equation (50) yields:

$$\left\| (\mathbf{g}(0), \mathbf{V}(0)) - (\tilde{\mathbf{g}}(0), \tilde{\mathbf{V}}(0)) \right\|^2 \leq \exp(-2\beta T).$$

Bounding the Evolved Distance

With the initial distance bounded, we can now use Razin et al. (2021, Lemma 6)², to control how the distance

$\left\| (\mathbf{g}(t), \mathbf{V}(t)) - (\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)) \right\|^2$ evolves over time. In other words:

$$\begin{aligned} \left\| (\mathbf{g}(t), \mathbf{V}(t)) - (\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t)) \right\|^2 &\leq \left\| (\mathbf{g}(0), \mathbf{V}(0)) - (\tilde{\mathbf{g}}(0), \tilde{\mathbf{V}}(0)) \right\|^2 \exp(2\beta t) \\ &\leq (N(2L+1) + 2K)\exp(2\beta t) L\alpha_V^2 \end{aligned} \tag{54}$$

holds until $t \in [0, T]$ or $\|(\tilde{\mathbf{g}}(t), \tilde{\mathbf{V}}(t))\| \geq \tilde{D}$ where, $\tilde{D} = N^{\frac{1}{3}}\sqrt{2L+1} \left(\frac{D}{2} + 0.5\right)^{\frac{1}{3}}$ and $D \geq 2\rho\sqrt{K}$.

²(Razin et al., 2021, Lemma 6): Let $\theta, \theta' : [0, T] \rightarrow \mathbb{R}^d$, where $T > 0$, be two curves born from gradient flow over a continuously differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$:

$$\begin{aligned} \theta(0) &= \theta_0 \in \mathbb{R}^d, \quad \frac{d}{dt}\theta(t) = -\nabla f(\theta(t)), t \in [0, T], \\ \theta'(0) &= \theta'_0 \in \mathbb{R}^d, \quad \frac{d}{dt}\theta'(t) = -\nabla f(\theta'(t)), t \in [0, T]. \end{aligned}$$

Let $D > 0$, and suppose that $f(\cdot)$ is β -smooth over \mathcal{D}_{D+1} for some $\beta \geq 0$,¹⁶ where $\mathcal{D}_{D+1} := \{\theta \in \mathbb{R}^d : \|\theta\| \leq D+1\}$. Then, if $\|\theta(0) - \theta'(0)\| < \exp(-\beta T)$, it holds that:

$$\|\theta(t) - \theta'(t)\| \leq \|\theta(0) - \theta'(0)\| \exp(\beta t)$$

at least until $t \geq T$ or $\|\theta'(t)\| \geq D$. That is, Equation (13) holds for all $t \in [0, \min\{T, T_D\}]$, where $T_D := \inf\{t \geq 0 : \|\theta'(t)\| \geq D\}$.

Note that,

$$\begin{aligned} (N(2L+1) + 2K) \exp(2\beta T) L \alpha_V^2 &\stackrel{a}{\leq} (N(2L+1) + 2K) \frac{1}{N(2L+1)} \\ &\leq 1 + \frac{2K}{N(2L+1)} \leq 1 + \frac{2}{2L+1} \leq 2. \end{aligned} \quad (55)$$

The inequality (a) is obtained by substituting the upper bound on the value α_V from equation (50). The following inequalities follow by noting that $K \leq N$, and $L \geq 1$. Applying the reverse triangle inequality to (54) and using (55) yields $\|(\mathbf{g}(t), \mathbf{V}(t))\| < \tilde{D} + 2$.

Using Lemma C.3, we can translate the bound on the parameters to a bound on the matrices. This yields:

$$\begin{aligned} \|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\|_F &< 2(\tilde{D} + 2)^2 \sqrt{2(N(2L+1) + 2K)} \exp(\beta t) \alpha_V \\ &< 2(\tilde{D} + 2)^2 \sqrt{2NL(2L+3)} \exp(\beta T) \alpha_V, \end{aligned} \quad (56)$$

where $\beta = 16NL^{3/2}((N+1)\mu C^4 + MB_Y C)$ as derived in C.2. Let the degree of approximation be denoted as $\epsilon_{app} \in (0, 1)$. To ensure this in (56) is less than a desired approximation error ϵ_{app} , we introduce a stricter condition on the initialization as:

$$\alpha_V \leq \frac{\epsilon_{app}}{2(\tilde{D} + 2)^2} \frac{\exp(-\beta T)}{\sqrt{2LN(2L+3)}}.$$

Thus for $t \in [0, T]$ we have

$$\|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\|_F < \epsilon_{app}. \quad (57)$$

The proof is complete if we can show that the condition $\|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\|_F < \epsilon_{app}$ holds for all $t \in [0, \min\{T, T_D\}]$, where $T_D := \inf\{t \geq 0 \mid \|\mathbf{X}(t)\|_F \geq D\}$. Let $t' \in [0, T]$ be the first instance at which

$$\|(\tilde{\mathbf{g}}(t'), \tilde{\mathbf{V}}(t'))\|^2 = \tilde{D}^2$$

Using the definition of distance equation, we have:

$$L\|\tilde{\mathbf{g}}(t')\|_2^2 + \|\tilde{\mathbf{V}}(t')\|_F^2 = \tilde{D}^2. \quad (58)$$

The proof is complete if we can show that for $t \in [0, t']$ the bound (57) holds, at $t = t'$, $\|\mathbf{X}(t')\| \geq D$.

Assuming perfect balancedness condition $(\frac{1}{2}\|\mathbf{g}(t')\|_2^2 = \|\mathbf{V}(t')\|_F^2)$ in (58), we have:

$$\|(\tilde{\mathbf{g}}(t'), \tilde{\mathbf{V}}(t'))\|_2^2 = 2\tilde{D}^2/(2L+1).$$

Recall that norm of the reference matrix $\tilde{\mathbf{X}}(t)$ is:

$$\|\tilde{\mathbf{X}}(t')\|_F^2 = \sum_{n \in [N]} \tilde{g}_n^4(t') \sum_{l \in [L]} \tilde{V}_{lj}^2(t').$$

From row-balancedness condition $\left(\frac{1}{2}\tilde{g}_i^2(t) = \sum_{l \in [L]} \tilde{V}_{lj}^2(t)\right)$, we can write it as,

$$\|\tilde{\mathbf{X}}(t')\|_F^2 = \frac{1}{2} \sum_{n \in [N]} \tilde{g}_n^6(t'), \quad (59)$$

Using the generalized mean inequality, namely, for $p < q$, $\left(\frac{1}{N} \sum_{i=1}^N x^p\right)^{1/p} \leq \left(\frac{1}{N} \sum_{i=1}^N x^q\right)^{1/q}$.

Let $p = 1$ and $q = 3$, $x = \tilde{g}_n^2(t')$.

$$\begin{aligned} \left(\frac{1}{N} \sum_i^N \tilde{g}_n^2(t')\right) &\leq \left(\frac{1}{N} \sum_i^N (\tilde{g}_n^2(t'))^3\right)^{1/3} \\ \left(\frac{1}{N} \sum_i^N \tilde{g}_n^2(t')\right)^3 &\leq \left(\frac{1}{N} \sum_i^N (\tilde{g}_n^2(t'))^3\right) \end{aligned}$$

We can now substitute this into (59):

$$\begin{aligned} \|\tilde{\mathbf{X}}(t')\|_F^2 &\geq \frac{1}{2N^2} \left(\sum_{n \in [N]} \tilde{g}_n^2(t')\right)^3 \\ &\geq \frac{1}{2N^2} \left(\frac{2\tilde{D}^2}{2L+1}\right)^3 \\ &\geq \frac{4\tilde{D}^6}{((2L+1)^3 N^2)}. \end{aligned} \tag{60}$$

Substituting \tilde{D} in (60) gives:

$$\|\tilde{\mathbf{X}}(t')\|_F \geq \frac{2\tilde{D}^3}{N(2L+1)^{3/2}} \geq D+1. \tag{61}$$

But at time $t = 0$, by the choice of our initialization we have $\|\tilde{\mathbf{X}}(0)\|_F = 2\rho\sqrt{K} < D+1$. This implies that $t' > 0$. Now, we focus on the time interval $t \in [0, t')$. Note that until $t < t'$, $\|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\|_F \leq \epsilon_{app} < 1$ holds.

Using continuity with respect to time, at $t = t'$ we have

$$\|\mathbf{X}(t') - \tilde{\mathbf{X}}(t')\|_F \leq \epsilon_{app} < 1.$$

Together with equation (61), and applying reverse triangular inequality,

$$\begin{aligned} \|\mathbf{X}(t')\|_F &\geq \|\tilde{\mathbf{X}}(t')\|_F - \epsilon_{app} \\ &> \|\tilde{\mathbf{X}}(t')\|_F - 1 \\ &> D. \end{aligned}$$

This shows that at time $t = t'$, $\|\mathbf{X}(t')\|_F > D$. This guarantees that

$$\|\mathbf{X}(t') - \tilde{\mathbf{X}}(t')\|_F \leq \epsilon_{app}$$

until time T or time t' at which $\|\mathbf{X}(t')\|_F \geq D$. This completes the proof.

C.2 Proof of Corollary 1

Corollary C.0.1 (Corollary 5.1.1 restated). *Assume the conditions of Theorem 5.1, and in addition assume that all reference trajectories $\tilde{\mathbf{X}}(t)$ converge to a solution $\mathbf{X}^* \in \mathbb{R}^{N \times L}$. This convergence is uniform in the sense that the trajectories are confined to a bounded domain, and for any $\epsilon_{app} > 0$ there exists a time $T_c \leq T$ after which they are all within a distance ϵ_{app} from \mathbf{X}^* . Then for any $\epsilon_{app} > 0$, if the initialization scales α_V, α_g are sufficiently small, for any time $t \in [T_c, T]$ it holds that $\|\mathbf{X}(t) - \mathbf{X}^*\|_F \leq 2\epsilon_{app}$*

Proof. For $\epsilon_{app} > 0$, there exists a time $T_c > 0$ after which all the reference trajectories are within a distance of ϵ_{app} from the solution \mathbf{X}^* . i.e.,

$$\|\tilde{\mathbf{X}}(t) - \mathbf{X}^*\| \leq \epsilon_{app}.$$

where $t \geq T_c$. These reference trajectories are also confined within a ball of radius D . i.e., $\|\tilde{\mathbf{X}}(t)\| \leq D$ for all $t \geq 0$. From Theorem 5.1, if initialization scales α_V , α_g are sufficiently small, then

$$\|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\| \leq \epsilon_{app} < 1.$$

This bound holds atleast until $t \geq T$ or $\|\mathbf{X}(t)\| \geq D + 1$. In line with the analysis done by Razin et al. (2021), by the way of contradiction, we show that the second condition cannot hold, namely $\|\mathbf{X}(t)\| \leq D + 1$ cannot be true for any $t \in [0, T]$. Let $t' \in [0, T]$ be an initial time at which $\|\mathbf{X}(t')\| \geq D + 1$. Since $t' \leq T$, the guarantee from Theorem 5.1 is still valid, meaning, $\|\mathbf{X}(t') - \tilde{\mathbf{X}}(t')\| < 1$, using the triangle inequality we have:

$$\begin{aligned} \|\tilde{\mathbf{X}}(t')\| &\geq \|\mathbf{X}(t')\| - \|\mathbf{X}(t') - \tilde{\mathbf{X}}(t')\| > (D + 1) - 1 \\ &> D \end{aligned}$$

This is a contradiction to the assumption that $\tilde{\mathbf{X}}(t)$ is confined to a ball of radius D for $t = t'$. Thus, for all $t \in [0, T]$, the guarantee $\|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\| \leq \epsilon_{app}$ holds.

For any time $t \in [T_c, T]$ we can apply the triangle inequality:

$$\|\mathbf{X}(t) - \mathbf{X}^*\| \leq \|\mathbf{X}(t) - \tilde{\mathbf{X}}(t)\| + \|\tilde{\mathbf{X}}(t) - \mathbf{X}^*\| \leq 2\epsilon_{app}$$

This completes the proof. \square

D Result on MNIST

To demonstrate the efficacy of the proposed algorithm on real-world data, we conduct a recovery experiment on a set of images from the MNIST dataset. The experiment uses the first 500 images from the MNIST dataset, each reshaped into a vector of 784 pixels. The images are normalized to have pixel values between 0 and 1. A random Gaussian sensing matrix \mathbf{A} of size 1024×784 with ℓ_2 -normalized columns is used to generate the measurements. The measurement matrix \mathbf{Y} is then computed as $\mathbf{Y} = \mathbf{A}\mathbf{X}$, where \mathbf{X} represents the matrix of vectorized MNIST images. The proposed IR-MMV algorithm is used to recover the images from the measurements \mathbf{Y} and the sensing matrix \mathbf{A} . For comparison, we also show the results from the M-OMP, M-SP and M-FOCUSS algorithm. For M-OMP and M-SP, which require prior knowledge of the row sparsity level, is set to $K = 18$. The simulations were performed on the CPU of a desktop computer equipped with an Intel Core i7-12700 processor, 16GB of DDR5 RAM. The MNIST dataset was obtained under the Creative Commons Attribution-Share Alike 3.0 license.

As shown in Figure 4, the proposed method successfully recovers the digit images, achieving a visual quality comparable to that of the other approaches without requiring any prior information about the signal sparsity.

E Code and Reporducability

The source code, and instructions needed to reproduce all experimental results are provided in the supplementary material as a ZIP file. All experiments were conducted in a Python 3 environment with dependencies listed in ‘requirements.txt’. To reproduce all results, please follow the detailed instructions in the README.md file included in the ZIP file.

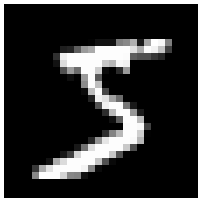
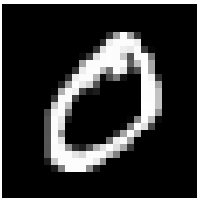
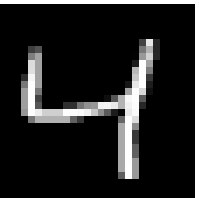
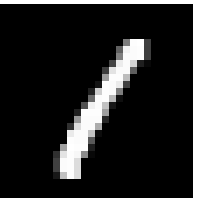
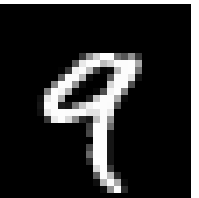
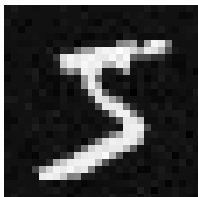
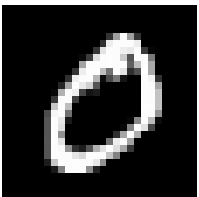
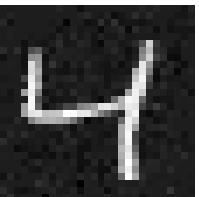
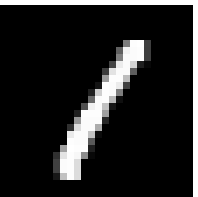
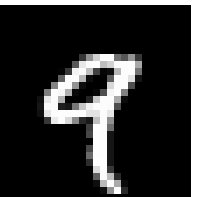
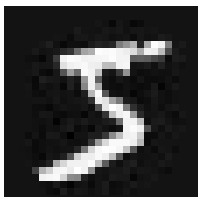
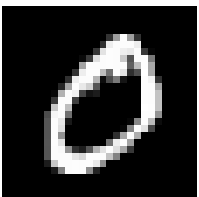
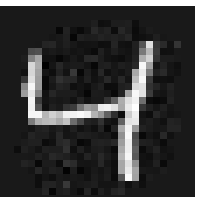
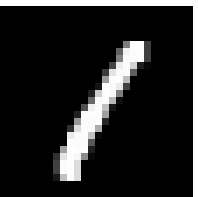
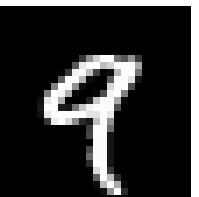
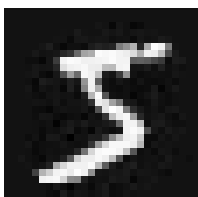
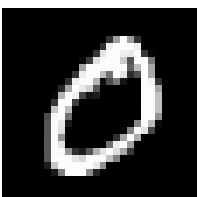
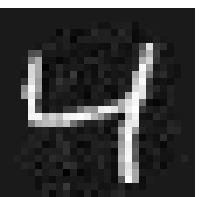
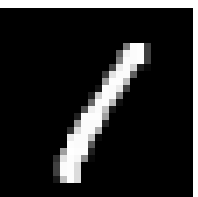
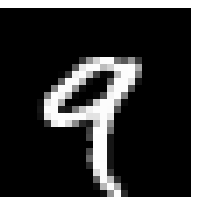
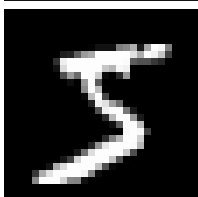
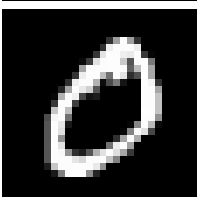
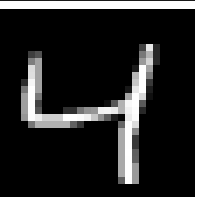
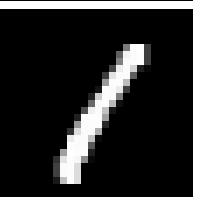
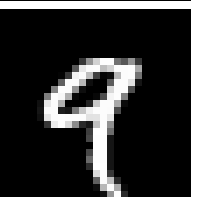
Method	Digit 5	Digit 0	Digit 4	Digit 1	Digit 9
Original					
Proposed					
M-OMP					
M-SP					
M-FOCUSS					

Figure 4: Comparison of MNIST digit reconstruction using different methods. The first row shows the original images.