# On the stability of gradient flow dynamics for a rank-one matrix approximation problem

Hesameddin Mohammadi, Meisam Razaviyayn, and Mihailo R. Jovanović

*Abstract*— In this paper, we examine the global stability of gradient flow dynamics associated with the problem of finding the best rank-one approximation of a given matrix. We partition the state-space into an infinite family of invariant manifolds over which the dynamics reduce to the special case of approximating a symmetric matrix. This allows us to employ a Lyapunov-based argument to explicitly characterize the region of attraction for the stable equilibrium points. This characterization proves an almost everywhere convergence for the gradient flow dynamics to the minimizers of the corresponding rank-one approximation problem.

## I. INTRODUCTION

Many modern inference problems, such as matrix completion [1], training of neural networks [2], and phase retrieval [3] require solving large-scale non-convex optimization problems. Although many of these optimization problems are known to be NP-hard in general [2], [4], *typical* problem instances can be solved in a polynomial time [1], [3]. For example, despite NP-hardness of the neural networks training problem [2], it was recently numerically observed that the gradient descent algorithm converges to the set of global optima for most initializations [5]. To some extent, these observations have been justified theoretically under various simplifying assumptions by showing either the existence of no spurious local optima or the exponentially vanishing number of spurious local optima [6], [7]. More specifically, in the matrix completion and the training of *linear* neural networks problems, it was recently established that all local optima are globally optimal [1], [6], [8]; consequently, the gradient descent method converges to the global optima despite non-convexity of the optimization problems [9]. Furthermore, for nonlinear neural networks, the recent work [7] shows that *most* of the local optima are globally optimal.

For non-convex learning problems, not all global optima result in the same statistical performance [10]. Hence, the choice of optimization algorithm and the initialization play a crucial role in biasing toward a specific global optima. This bias, which is also known as implicit regularization, is central to the understanding of various algorithms in non-convex problems with multiple global/local optima [5], [10]–[13].

As a first step toward better understanding of implicit regularization on learning problems, we consider the rank-one approximation problem of a given matrix under the gradient flow dynamics. Although this problem may appear much less involved than the aforementioned optimization problems, its complete understanding is crucial because it is the major building block for many non-convex learning problems (including matrix completion and training of neural networks). This problem is equivalent to the matrix completion under the full observation and rank-one restriction [1], and so to the problem of training linear neural networks with one hidden unit [6]. For this low rank approximation problem, all local optima are known to be globally optimal despite non-convexity of the optimization problem [14]. However, the behavior of the gradient descent algorithm, which is central to many learning tasks, is not yet fully understood. In this work, we examine the behavior of the gradient flow dynamics associated with this non-convex optimization problem.

The solution to the rank-one approximation problem under $\ell_2$ distance is closely related to the principal eigenspace problem; see [15, Chapter 1] for a brief survey on the existing algorithms. Among various procedures for finding the principal eigenspace, the power method, Rayleigh quotient procedure, and the Oja's method attracted significant attention because of their simplicity and scalability [15]–[17]. In particular, the Oja's flow [15], [18], [19], which is based on the non-normalized version of the Rayleigh quotient flow, has gained popularity in recent years for online (i.e., streaming) principal component analysis [17]. It is also worth noting that both Rayleigh quotient and the Oja's flows can be modified by regularizing the $\ell_2$ distance of the estimated rank-one matrix with a certain scaled version of the identity mapping [20]–[22]. This simple, yet insightful, modification allows for tracking of both the major and minor components of a given symmetric matrix [22], [23]. Although this modified gradient flow dynamics enjoy a globally stable set of equilibrium points, the modification requires *a priori* knowledge of the eigenvalues of the matrix, which is not always available.

In this paper, we consider the gradient flow dynamics for the rank-one approximation of a given matrix. By partitioning the state-space via the introduction of novel invariant manifolds, we reduce the problem to the simpler problem of rank-one approximation of a *symmetric* matrix. Building upon this connection and through the introduction of maximal Lyapunov functions, we completely characterize the regions of attraction of the stable equilibrium points.

Our presentation is organized as follows. In Section II, we formulate the problem in the general form and also discuss the important special case of symmetric matrices. In Section III, we use maximal Lyapunov functions to characterize

H. Mohammadi and M. R. Jovanović are with the Ming Hsieh Department of Electrical Engineering and M. Razaviyayn is with the Epstein Department of Industrial and Systems Engineering, University of Southern California, Los Angeles, CA 90089. E-mails: hesamedm@usc.edu, razaviya@usc.edu, mihailo@usc.edu.

the regions of attraction of the stable equilibrium points for the symmetric problem. In Section IV, we introduce a set of invariant manifolds that partition the state-space of the general problem. Over each of these manifolds, we further demonstrate that the general problem reduces to the special case of symmetric matrix approximation. This reduction leads to a complete characterization for the region of attraction of the stable equilibrium points.

## II. PROBLEM FORMULATION

In this section, we formulate the problem of finding the best rank-one approximation of a given matrix. Even though the solution is well known and determined by the singular value decomposition, the properties of the gradient flow dynamics associated with this nonconvex optimization problem are not well-understood. We first discuss the problem of approximating an $m \times n$ matrix and then introduce a special case (i.e., the best rank-one approximation of a given symmetric matrix) which plays a central role in the analysis of the general problem.

### A. Rank-one matrix approximation problem

Consider the rank-one approximation problem

$$\underset{x, y}{\text{minimize}} \quad \frac{1}{2} \|xy^T - M\|_F^2 \qquad \text{(P1)}$$

where $M \in \mathbb{R}^{m \times n}$ is a given matrix, $x \in \mathbb{R}^m$ and $y \in \mathbb{R}^n$ are the optimization variables, and $\| \cdot \|_F$ is the Frobenius norm. The solution to this nonconvex optimization problem is determined by the principal singular vectors of the matrix $M$. In this paper, we examine the gradient flow dynamics associated with (P1),

$$\begin{bmatrix} \dot{x} \\ \dot{y} \end{bmatrix} = \begin{bmatrix} My - (y^T y)\, x \\ M^T x - (x^T x)\, y \end{bmatrix} \qquad (1)$$

and analyze the stability properties of the equilibrium points.

Let $M = U\Sigma V^T$ be a singular value decomposition of $M$, where $U := [\, u_1 \cdots u_m \,]$ and $V := [\, v_1 \cdots v_n \,]$ are unitary matrices, and $\Sigma \in \mathbb{R}^{m \times n}$ is the matrix of singular values. In what follows, we restrict our analysis to the case where the principal singular value $\sigma_1$ is strictly larger than the other singular values.

*Assumption 1:* The singular values of the rank $r$ matrix $M$ satisfy $\sigma_1 > \sigma_2 \geq \cdots \geq \sigma_r$.

If Assumption 1 holds, the minimizers of (P1) are

$$(x^\star, y^\star) = (\frac{\sigma_1}{c}\, u_1, c\, v_1)$$

where $c$ is a nonzero number. However, the minimizers of (P1) are not the only equilibrium points of (1). In fact, any pair

$$(\bar{x}, \bar{y}) := \begin{cases} (\frac{\sigma_i}{c}\, u_i, c\, v_i), & i \in \{1, \ldots, r\} \\ (0, v), & v \in \mathcal{N}(M) \\ (u, 0), & u \in \mathcal{N}(M^T) \end{cases} \qquad (3)$$

is an equilibrium point of (1), where $\mathcal{N}(\cdot)$ denotes the null space of a given matrix.

The change of variables $x := U\xi$, $y := V\eta$ brings gradient flow dynamics (1) into the following form

$$\begin{bmatrix} \dot{\xi} \\ \dot{\eta} \end{bmatrix} = \begin{bmatrix} \Sigma \eta - (\eta^T \eta)\, \xi \\ \Sigma^T \xi - (\xi^T \xi)\, \eta \end{bmatrix}. \qquad \text{(GD)}$$

Let $\mathrm{e}_i$ and $\hat{\mathrm{e}}_j$ denote the unit vectors in the canonical basis of $\mathbb{R}^m$ and $\mathbb{R}^n$, respectively. Under this change of variable, the minimizers of (P1) are given by

$$(\xi^\star, \eta^\star) = (\frac{\sigma_1}{c}\, \mathrm{e}_1, c\, \hat{\mathrm{e}}_1) \qquad (2)$$

and any pair

$$(\bar{\xi}, \bar{\eta}) := \begin{cases} (\frac{\sigma_1}{c}\, \mathrm{e}_i, c\, \hat{\mathrm{e}}_i), & i \in \{1, \ldots, r\} \\ (0, \hat{\mathrm{e}}), & \hat{\mathrm{e}} \in \mathcal{N}(\Sigma) \\ (\mathrm{e}, 0), & \mathrm{e} \in \mathcal{N}(\Sigma^T) \end{cases} \qquad (3)$$

is an equilibrium point of (GD).

### B. Symmetric case

A special instance of problem (P1) is given by

$$\underset{z}{\text{minimize}} \quad \frac{1}{4} \|zz^T - W\|_F^2 \qquad \text{(P2)}$$

where $W \in \mathbb{R}^{n \times n}$ is a symmetric matrix and $z \in \mathbb{R}^n$ is the optimization variable. The corresponding gradient flow dynamics simplifies to

$$\dot{z} = (W - (z^T z)\, I_n)\, z \qquad (4)$$

where $I_n$ is the $n \times n$ identity matrix. As we demonstrate in Section IV, the analysis of (4) allows us to characterize the trajectories of general gradient flow dynamics (GD).

Let $W = U\Lambda U^T$ be an eigenvalue decomposition of $W$ where $\Lambda$ is the diagonal matrix of eigenvalues and $U := [\, u_1 \cdots u_n \,]$ is a unitary matrix of eigenvectors. If the largest eigenvalue $\lambda_1$ of $W$ is not positive, the only equilibrium point of (4) is $\bar{z}_0 = 0$. To avoid this trivial case, we assume $\lambda_1 > 0$. We also further restrict our analysis to the situation where $\lambda_1$ is strictly larger than the other eigenvalues.

*Assumption 2:* The eigenvalues of the matrix $W$ satisfy

$$\lambda_1 > \lambda_2 \geq \cdots \geq \lambda_n \quad \text{and} \quad \lambda_1 > 0.$$

Under Assumption 2, the minimizers of (P2) are given by $z^\star = \pm\sqrt{\lambda_1}\, u_1$. Finally, in addition to the origin, for any positive eigenvalue $\lambda_i > 0$,

$$\bar{z}_i := \pm\sqrt{\lambda_i}\, u_i, \qquad (5)$$

is an equilibrium point of (4). Note that (4) may have infinitely many equilibrium points because any $u \in \Theta_i$ such that $u^T u = \sqrt{\lambda_i}$, is an equilibrium point of (4), where $\Theta_i$ is the eigenspace corresponding to the eigenvalue $\lambda_i > 0$. However, since the choice of $U$ is arbitrary, our stability analysis of equilibrium points (5) covers all the equilibrium points of (4).

## III. GRADIENT FLOW DYNAMICS: SYMMETRIC CASE

In this section, we examine the symmetric rank-one approximation problem (P2) and employ a Lyapunov-based

approach to studying the behavior of gradient flow dynamics (4). We first consider the linearized version of (4) and demonstrate that the global minimizers of (P2) are the only locally stable equilibrium points of (4). We then continue our analysis by explicitly characterizing the region of attraction for each of the locally stable equilibrium points.

### A. Linearization around the equilibrium points

The linearization of (4) around the equilibrium point $\bar{z}_i$ is given by

$$\dot{z} = \left(W - (\bar{z}_i^T \bar{z}_i) I_n - 2\bar{z}_i \bar{z}_i^T\right) z =: A_i z.$$

The eigenvalues of $A_i$ are determined by the eigenvalues of the matrix $W$,

$$\text{eig}(A_i) = \{-2\lambda_i; \ \lambda_k - \lambda_i, \ k \in \{1, \ldots, n\}\backslash\{i\}\}. \quad (6)$$

Under Assumption 2, all eigenvalues of the matrix $A_1$ are negative, thereby implying local asymptotic stability of $\bar{z}_1 = \pm\sqrt{\lambda_1} u_1$. In contrast, for any other positive eigenvalue $\lambda_i \neq \lambda_1$ of $W$, the matrix $A_i$ has a positive eigenvalue and $\bar{z}_i$ is unstable. Finally, it is easy to show that $\bar{z}_0 = 0$ is also unstable if Assumption 2 holds.

Next, we establish a global convergence result by showing that the regions of attraction of the locally stable equilibrium points ($\sqrt{\lambda_1} u_1$ and $-\sqrt{\lambda_1} u_1$) are the two open half spaces corresponding to an invariant hyperplane.

### B. Regions of attraction of stable equilibrium points

We first establish the existence of an invariant hyperplane $\mathcal{H}$ with respect to gradient flow dynamics (4) which separates the two locally stable equilibrium points $\bar{z}_1 = \pm\sqrt{\lambda_1} u_1$. We then use a Lyapunov-based argument to show that the two open half spaces corresponding to $\mathcal{H}$ are the regions of attraction for $\sqrt{\lambda_1} u_1$ and $-\sqrt{\lambda_1} u_1$, respectively.

*Definition 1:* For any $f \colon \mathbb{R}^n \to \mathbb{R}^n$, $\mathcal{S} \subset \mathbb{R}^n$ is an invariant set with respect to

$$\dot{z} = f(z) \quad (7)$$

if, for any initial condition $z(0) \in \mathcal{S}$, $z(t) \in \mathcal{S}$ for all $t \geq 0$. Moreover, the region of attraction of an equilibrium point $\bar{z}$ is the set of all initial conditions for which $\lim_{t \to \infty} z(t) = \bar{z}$.

In Lemma 1, we characterize a family of hyperplanes that are invariant sets with respect to (4).

*Lemma 1:* For any eigenvector $u$ of a symmetric matrix $W$, the hyperplane $\mathcal{H} := \{z \,|\, z^T u = 0\}$ is an invariant set with respect to (4).

*Proof:* The inner product of the eigenvector $u$ with the the vector field in (4) is given by

$$u^T f(z) = u^T \left(W - (z^T z) I_n\right) z = u^T z \left(\lambda - z^T z\right)$$

where $\lambda$ is the corresponding eigenvalue of $W$. Thus, if $u^T z = 0$, then $u^T f(z) = 0$. This proves that $\mathcal{H}$ is an invariant hyperplane with respect to (4). ∎

To characterize the regions of attraction of (4) around the stable equilibrium points, we employ the notion of maximal Lyapunov functions introduced in [24].

*Definition 2:* Let $\mathcal{X} \subset \mathbb{R}^n$ be an open set that contains an equilibrium point $\bar{z}$ of (7) and let $\partial\mathcal{X}$ denote the boundary of $\mathcal{X}$. A differentiable function $V \colon \mathcal{X} \to \mathbb{R}$ is a maximal Lyapunov function for (7) if it satisfies

1) $V(\bar{z}) = 0$, $V(z) > 0$, $\forall z \in \mathcal{X}\backslash\{\bar{z}\}$;
2) $V(z) \to \infty$ as $z \to \infty$ and/or $z \to \partial\mathcal{X}$;
3) $\dot{V}(\bar{z}) = 0$, $\dot{V}(z) < 0$, $\forall z \in \mathcal{X}\backslash\{\bar{z}\}$.

The following theorem states that if there exists a maximal Lyapunov function for (7), then $\mathcal{X}$ is the region of attraction for the equilibrium point $\bar{z}$.

*Theorem 1 (Theorem 1 in [24]):* Let $\bar{z} \in \mathcal{X}$ be an equilibrium point of (7) and let $\mathcal{X} \subset \mathbb{R}^n$ be an open set. If there exists a maximal Lyapunov function $V \colon \mathcal{X} \to \mathbb{R}$ as in Definition 2, then $\mathcal{X}$ is the region of attraction of $\bar{z}$.

In Theorem 2, we show that there exists a hyperplane $\mathcal{H}$ for which the corresponding half spaces are the regions of attraction for the equilibrium points $\bar{z}_1 = \pm\sqrt{\lambda_1} u_1$ of (4).

*Theorem 2:* Let Assumption 2 hold. Then, the two open half spaces $\mathcal{H}_\pm := \{z \,|\, \pm z^T u_1 > 0\}$ are the regions of attraction for the equilibrium points $\pm\sqrt{\lambda_1} u_1$ of (4).

*Proof:* We find a maximal Lyapunov function to prove the result for $\sqrt{\lambda_1} u_1$. The proof for $-\sqrt{\lambda_1} u_1$ follows from similar arguments and is omitted for brevity.

Let $U \Lambda U^T$ be an eigenvalue decomposition of $W$ where $\Lambda := \text{diag}([\lambda_1 \cdots \lambda_n])$ and $U = [u_1 \cdots u_n]$ is the unitary matrix of eigenvectors. To simplify the presentation, we use the change of variable $z := U\zeta$ to rewrite (4) as

$$\dot{\zeta} = \left(\Lambda - (\zeta^T \zeta) I_n\right) \zeta. \quad (8)$$

This change of variable brings equilibrium point $\bar{z}_1 = \sqrt{\lambda_1} u_1$ to $\bar{\zeta}_1 = \sqrt{\lambda_1}\, e_1$, where $e_i$ is the $i$th unit vector in the canonical basis of $\mathbb{R}^n$. Similarly, the half-space $\mathcal{H}_+$ becomes $\hat{\mathcal{H}}_+ = \{\zeta \,|\, e_1^T \zeta > 0\}$.

Now, we employ Theorem 1 to establish that $\hat{\mathcal{H}}_+$ is the region of attraction for the equilibrium point $\bar{\zeta}_1$ of (8). In particular, we propose the maximal Lyapunov function candidate $V \colon \hat{\mathcal{H}}_+ \to \mathbb{R}$,

$$V(\zeta) := \frac{\|\zeta - \bar{\zeta}_1\|^2}{e_1^T \zeta}. \quad (9)$$

Clearly, $V(\zeta)$ satisfies conditions 1) and 2) in Definition 2. The remaining task is to show that 3) is also satisfied.

The derivate of $V$ along the solutions of (4) is given by

$$
\begin{aligned}
\dot{V}(\zeta) &= \frac{2\dot{\zeta}^T \left(\zeta - \bar{\zeta}_1\right)}{e_1^T \zeta} - e_1^T \dot{\zeta} \frac{\|\zeta - \bar{\zeta}_1\|^2}{\left(e_1^T \zeta\right)^2} \\
&= \frac{2\zeta^T \left(\Lambda - (\zeta^T \zeta) I_n\right)\left(\zeta - \bar{\zeta}_1\right)}{e_1^T \zeta} - \\
&\quad \left(\lambda_1 - \zeta^T \zeta\right) \frac{\|\zeta - \bar{\zeta}_1\|^2}{e_1^T \zeta}.
\end{aligned}
$$

Substitution of $\sqrt{\lambda_1}\,\mathrm{e}_1$ for $\bar{\zeta}_1$ into $\dot{V}$ yields

$$\dot{V}(\zeta) \;=\; -\frac{\left(\zeta^T\zeta\right)^2 - 2\,\zeta^T\Lambda\,\zeta \;+\; \lambda_1^2}{\mathrm{e}_1^T\zeta}$$

which implies $\dot{V}(\bar{\zeta}_1) = 0$. Note that $\mathrm{e}_1^T\zeta > 0$ for any $\zeta \in \hat{\mathcal{H}}_+$. Moreover, $\zeta^T\Lambda\zeta \leq \lambda_1\zeta^T\zeta$ and the equality holds only when $\mathrm{e}_i^T\zeta = 0$ for all $i > 1$. Therefore,

$$\dot{V}(\zeta) \;\leq\; -\frac{\left(\zeta^T\zeta\right)^2 - 2\,\lambda_1\zeta^T\zeta \;+\; \lambda_1^2}{\mathrm{e}_1^T\zeta} \;\leq\; 0$$

where at least one of the above inequalities is strict when $\zeta \neq \bar{\zeta}_1$. Thus, $\dot{V}(\zeta) < 0$ for every $\zeta \in \hat{\mathcal{H}}_+\backslash\{\bar{\zeta}_1\}$, and Theorem 1 implies that $\hat{\mathcal{H}}_+$ is the region of attraction for the equilibrium point $\bar{\zeta}_1$ of (8) or, equivalently, that $\mathcal{H}_+$ is the region of attraction for the equilibrium point $\bar{z}_1$ of (4). ∎

Figure 1 illustrates the convergence of trajectories of (4) to the global minimizers of (P2), for

$$W \;=\; \begin{bmatrix} 1.36 & 0.48 \\ 0.48 & 1.64 \end{bmatrix}. \tag{10}$$

The thick red line marks the boundary of the half spaces $\mathcal{H}_\pm$ (Theorem 2); the black dash-dotted curves mark the level sets of $V$ in (9), and the filled and hollow circles mark the global minimizers $\bar{z}_1$ and the unstable eq. points $\bar{z}_2$, respectively.
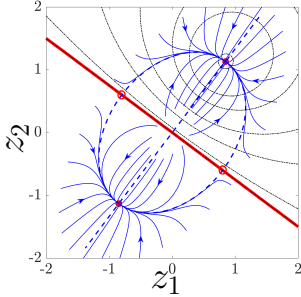


Fig. 1: Trajectories of (4) for the matrix $W$ given by (10) and the level-sets of a maximal Lyapunov function.

## IV. GRADIENT FLOW DYNAMICS: GENERAL CASE

In this section, we study gradient flow dynamics (GD) associated with the rank-one approximation of a matrix $M \in \mathbb{R}^{m \times n}$. We first identify unstable equilibrium points of (GD) by analyzing stability properties of the corresponding linearized systems. Next, we partition the state-space into invariant manifolds $\mathcal{S}_b := \{(\xi,\eta)|\,\xi^T\xi - \eta^T\eta = b\}$ parameterized by a real number $b$. This partitioning allows us to make a connection between (GD) and (4). In particular, we demonstrate that for any $b \in \mathbb{R}$, there exists a symmetric $W \in \mathbb{R}^{(m+n) \times (m+n)}$ such that for any initial conditions $(\xi_0, \eta_0) \in \mathcal{S}_b$ and $z_0 := [\,\xi_0^T\;\eta_0^T\,]^T$, the trajectories of (GD) and (4) are two different parameterizations of the same curve. We then utilize this relation to establish that, under Assumption 1, the trajectory of (GD) converges to a global

minimizer of (P1) for any initial condition apart from a set of measure zero.

### A. Linearization around the equilibrium points

In this subsection, we use linearization to establish that if the conditions in Assumption 1 hold then, except for the global minimizers of (P1), all other equilibrium points of gradient flow dynamics (GD) are unstable.

Let $\mathrm{e}_j \in \mathbb{R}^m$ and $\hat{\mathrm{e}}_j \in \mathbb{R}^n$ be the unit vectors in the canonical bases and let $\Gamma_i := \Sigma - 2\sigma_i\mathrm{e}_i\hat{\mathrm{e}}_i^T$. Linearization of (GD) around equilibrium points (3) is given by

$$\begin{bmatrix} \dot{\xi} \\ \dot{\eta} \end{bmatrix} \;=\; A \begin{bmatrix} \xi \\ \eta \end{bmatrix} \tag{11}$$

where, depending on the type of $(\bar{\xi}, \bar{\eta})$ in (3), the linearized dynamical generators are respectively given by

$$A_i \;:=\; \begin{bmatrix} -c^2 I_m & \Gamma_i \\ \Gamma_i^T & -\frac{\sigma_i^2}{c^2} I_n \end{bmatrix} \tag{12}$$

$$A_\infty \;:=\; \begin{bmatrix} -(\bar{\eta}^T\bar{\eta})\,I_m & \Sigma \\ \Sigma^T & 0 \end{bmatrix} \tag{13}$$

$$A_0 \;:=\; \begin{bmatrix} 0 & \Sigma \\ \Sigma^T & -(\bar{\xi}^T\bar{\xi})\,I_n \end{bmatrix}. \tag{14}$$

Under Assumption 1, it is straightforward to show that the matrices $A_\infty$, $A_0$, as well as $A_i$ for $i \neq 1$, have at least one positive eigenvalue. Moreover, all eigenvalues of $A_1$ are negative apart from a single eigenvalue at zero, for all $c \neq 0$. Thus, $(\bar{\xi}, \bar{\eta}) := (\frac{\sigma_1}{c}\,\mathrm{e}_1, c\,\hat{\mathrm{e}}_1)$ are the only candidates for stable equilibrium points of (GD). The presence of the zero eigenvalue makes the use of linearization inconclusive; stability of these equilibrium points is established in Subsection IV-D using a Lyapunov-based argument.

Note that these equilibrium points are exactly the global minimizers of (P1). However, since (P1) has infinitely many global minimizers, the stability analysis of (GD) is more subtle than that of the special case (4). In the next subsection, we characterize a family of invariant sets with respect to (GD) which allows us to carry out the analysis.

### B. Invariant manifolds and partitioning of the state-space

Herein, we partition the state-space of (GD) into a family of invariant manifolds. For any $b \in \mathbb{R}$, let us define

$$\mathcal{S}_b \;:=\; \{(\xi,\eta)|\,\xi^T\xi - \eta^T\eta = b,\; \xi \in \mathbb{R}^m,\; \eta \in \mathbb{R}^n\}. \tag{15}$$

Lemma 2 shows an invariance property of (GD) over $\mathcal{S}_b$.

*Lemma 2:* For any $b \in \mathbb{R}$, $\mathcal{S}_b$ is an invariant set with respect to (GD).

*Proof:* Follows from the fact that $\dot{\xi}$ and $\dot{\eta}$ given by (GD) satisfy $\mathrm{d}/\mathrm{d}t\left(\xi^T\xi - \eta^T\eta\right) = 0$. ∎

Lemma 2 implies that given an initial condition $(\xi_0, \eta_0)$, the trajectory of (GD) stays in $\mathcal{S}_{b_0}$, where $b_0 := \xi_0^T\xi_0 - \eta_0^T\eta_0$. Moreover, since the collection of sets $\{\mathcal{S}_b\}_{b \in \mathbb{R}}$ partitions the state-space of (GD), we can restrict the stability analysis of (GD) to $\mathcal{S}_b$ for arbitrary $b \in \mathbb{R}$. In the next subsection, we demonstrate that for any trajectory $\tilde{z} := [\xi^T\eta^T]^T$ of (GD),

there exists a symmetric matrix $W$ for which (4) has a trajectory $z$ that traverses the same curve as $\tilde{z}$.

## C. Reducing the asymmetric problem to the symmetric form

The following theorem establishes a connection between gradient flow dynamics (GD) and (4). In particular, we show that for any initialization of (GD), there exists a symmetric matrix $W \in \mathbb{R}^{(m+n)\times(m+n)}$ such that (4) has a trajectory that traverses the same curve as that of (GD). Even though the symmetric reformulation of (P1) has been proposed in the literature [13], to the best of our knowledge the connection between the two gradient flow dynamics is novel.

*Theorem 3:* Let $\tilde{z}^T(t) := \begin{bmatrix} \xi^T(t) & \eta^T(t) \end{bmatrix}$ where $(\xi, \eta)$ is the trajectory of (GD) with initial condition $(\xi_0, \eta_0)$. Define

$$W := \begin{bmatrix} (\xi_0^T \xi_0 - \eta_0^T \eta_0)I_m & 2\Sigma \\ 2\Sigma^T & (\eta_0^T \eta_0 - \xi_0^T \xi_0)I_n \end{bmatrix}. \quad (16)$$

The trajectory of (4) with initial condition $z(0) = \tilde{z}(0)$ satisfies $z(t) = \tilde{z}(2t)$ for all $t \geq 0$.

*Proof:* Let $b_0 := \xi_0^T \xi_0 - \eta_0^T \eta_0$. Lemma 2 implies that $\xi^T \xi - \eta^T \eta = b_0$ for all $t \geq 0$. Furthemore, since $\tilde{z}^T \tilde{z} = \xi^T \xi + \eta^T \eta$, we have

$$\xi^T(t)\xi(t) = \tfrac{1}{2} \left( \tilde{z}^T(t)\tilde{z}(t) + b_0 \right), \quad (17a)$$
$$\eta^T(t)\eta(t) = \tfrac{1}{2} \left( \tilde{z}^T(t)\tilde{z}(t) - b_0 \right). \quad (17b)$$

If we substitute (17) into (GD), we obtain

$$\dot{\xi} = \Sigma \eta - \tfrac{1}{2} \left( \tilde{z}^T \tilde{z} - b_0 \right) \xi$$
$$\dot{\eta} = \Sigma^T \xi - \tfrac{1}{2} \left( \tilde{z}^T \tilde{z} + b_0 \right) \eta$$

or, equivalently, $\dot{\tilde{z}} = \tfrac{1}{2} \left( W - (\tilde{z}^T \tilde{z}) I_{m+n} \right) \tilde{z}$. ∎

The relation between the trajectories of (GD) and (4) is important because it allows us to explicitly characterize the regions of attraction of (GD).

## D. Regions of attraction and global convergence

In this subsection, we use the relationship between (GD) and (4) to completely characterize the regions of attraction of (GD) for the minimizers of (P1). In the following Lemma, we demonstrate that for any $b \in \mathbb{R}$, the invariant set $\mathcal{S}_b$ contains exactly two of the global minimizers of (P1).

*Lemma 3:* Under Assumption 1, for any $b \in \mathbb{R}$, the invariant set $\mathcal{S}_b$ contains exactly two of the global minimizers of (P1). Moreover, these points are centrally symmetric.

*Proof:* The minimizers of (P1) are given by (2), $(\xi^\star, \eta^\star) = \left( \tfrac{\sigma_1}{c} e_1, c \hat{e}_1 \right)$. Now, for any $b \in \mathbb{R}$, the condition $\xi^{\star T} \xi^\star - \eta^{\star T} \eta^\star = b$ yields a quadratic equation in $c^2$ which has two real solutions for $c$ with opposite signs. ∎

Let $(\xi_b^\star, \eta_b^\star) \in \mathcal{S}_b$ be the global minimizer of (P1) corresponding to the positive $c$ introduced in the proof of Lemma 3. The next lemma introduces an invariant set $\mathcal{H}_b \subset \mathcal{S}_b$ for which $(\xi_b^\star, \eta_b^\star) \notin \mathcal{H}_b$ and $-(\xi_b^\star, \eta_b^\star) \notin \mathcal{H}_b$.

*Lemma 4:* Under Assumption 1, for any $b \in \mathbb{R}$, the set

$$\mathcal{H}_b := \mathcal{S}_b \cap \{(\xi,\eta) | \xi^T \xi_b^\star + \eta^T \eta_b^\star = 0\} \quad (19)$$

is an invariant set with respect to (GD). Moreover, neither $(\xi_b^\star, \eta_b^\star)$ nor $-(\xi_b^\star, \eta_b^\star)$ belongs to $\mathcal{H}_b$.

*Proof:* Since $\pm(\xi_b^\star, \eta_b^\star) \neq 0$, it is clear that $\pm(\xi_b^\star, \eta_b^\star) \notin \mathcal{H}_b$. Now, since $\mathcal{S}_b$ is an invariant set, it suffices to show that for any $(\xi, \eta) \in \mathcal{H}_b$, $\mathrm{d}/\mathrm{d}t \left( \xi^T \xi_b^\star + \eta^T \eta_b^\star \right) = 0$. Defining $\tilde{z}_b^\star := [\,\xi_b^{\star T} \ \eta_b^{\star T}\,]^T$ and $\tilde{z}(t) := [\,\xi^T(t) \ \eta^T(t)\,]^T$, we have

$$\frac{\mathrm{d}}{\mathrm{d}t} \left( \xi^T \xi_b^\star + \eta^T \eta_b^\star \right) = \frac{\mathrm{d}}{\mathrm{d}t} \left( \tilde{z}^T \tilde{z}_b^\star \right).$$

Now, from Theorem 3 it follows that

$$\frac{\mathrm{d}}{\mathrm{d}t} \left( \tilde{z}^T \tilde{z}_b^\star \right) = \frac{1}{2} \tilde{z}^T (W - (\tilde{z}^T \tilde{z}) I_{m+n}) \tilde{z}_b^\star$$

where $W := \begin{bmatrix} b\,I_m & 2\Sigma \\ 2\Sigma^T & -b\,I_n \end{bmatrix}$. Therefore, it suffices to show that for any $\tilde{z} = [\,\xi^T \ \eta^T\,]^T$ with $\tilde{z}^T \tilde{z}_b^\star = 0$ and $\xi^T \xi - \eta^T \eta = b$, we have

$$\tilde{z}^T (W - (\tilde{z}^T \tilde{z}) I_{m+n}) \tilde{z}_b^\star = 0. \quad (20)$$

Let $\alpha$ and $\beta$ be the first entries of $\xi_b^\star$ and $\eta_b^\star$, respectively. Since $\alpha\beta = \sigma_1$ and $\alpha^2 - \beta^2 = b$, we have

$$W^T \tilde{z}_b^\star = W \tilde{z}_b^\star = W[\alpha e_1^T \ \beta \hat{e}_1^T]^T = \left( \tilde{z}_b^{\star T} \tilde{z}_b^\star \right) \tilde{z}_b^\star$$

and (20) follows from the fact that $\tilde{z}^T \tilde{z}_b^\star = 0$. ∎

In Theorem 4, we combine Theorem 2, Theorem 3, and Lemma 4 to establish a global convergence result for (GD).

*Theorem 4:* Under Assumption 1, for any $b \in \mathbb{R}$, the sets

$$\mathcal{H}_b^+ := \mathcal{S}_b \cap \{(\xi,\eta) | \xi^T \xi_b^\star + \eta^T \eta_b^\star > 0\}$$
$$\mathcal{H}_b^- := \mathcal{S}_b \cap \{(\xi,\eta) | \xi^T \xi_b^\star + \eta^T \eta_b^\star < 0\}$$

are the respective regions of attraction for the equilibrium points $(\xi_b^\star, \eta_b^\star)$ and $-(\xi_b^\star, \eta_b^\star)$ of (GD).

*Proof:* Let $\mathcal{R}_b^+$ and $\mathcal{R}_b^-$ be the regions of attraction of $(\xi_b^\star, \eta_b^\star)$ and $-(\xi_b^\star, \eta_b^\star)$, respectively. The collection of invariant manifolds $\{\mathcal{S}_a\}_{a \in \mathbb{R}}$ partitions the state-space and hence, for any given $b \in \mathbb{R}$, $\mathcal{S}_b$ and $\mathbb{R}^{m+n}\backslash \mathcal{S}_b$ are both invariant with respect to (GD). Thus, $\mathcal{R}_b^\pm \subseteq \mathcal{S}_b$. Moreover, $\mathcal{R}_b^\pm \cap \mathcal{H}_b = \emptyset$ because of Lemma 4. Since, the sets $\mathcal{H}_b$, $\mathcal{H}_b^+$, and $\mathcal{H}_b^-$ partition $\mathcal{S}_b$, it suffices to show that $\mathcal{H}_b^+ \subseteq \mathcal{R}_b^+$ and $\mathcal{H}_b^- \subseteq \mathcal{R}_b^-$. We prove $\mathcal{H}_b^+ \subseteq \mathcal{R}_b^+$. The proof for $\mathcal{H}_b^- \subseteq \mathcal{R}_b^+$ follows from similar arguments and is omitted for brevity.

In order to show that $\mathcal{H}_b^+$ is a subset of $\mathcal{S}_b^+$, we use the relation between the trajectories of (GD) and (4). In particular, let $\tilde{z}_b^\star := [\xi_b^{\star T} \eta_b^{\star T}]^T$. Let $\tilde{z} := [\xi^T \eta^T]^T$ and $W$ be defined as in Theorem 3, where $(\xi, \eta)$ is the trajectory of (GD) with some initial condition $(\xi_0, \eta_0) \in \mathcal{H}_b^+$. Let $\lambda_1$ be the largest eigenvalue of $W$ with the corresponding unit eigenvector $u_1$. It is not hard to show that $\lambda_1 = \sqrt{b^2 + 4\sigma_1^2}$ is the strict largest eigenvalue of $W$ and therefore Assumption 2 holds. If $u_1^T \tilde{z}(0) > 0$, Theorem 3 in conjunction with Theorem 2 imply that $\tilde{z}$ converges to $\sqrt{\lambda_1} u_1$.

Hence, in order to establish the convergence of $\tilde{z}$ to $\tilde{z}_b^\star$, it suffices to show that $\tilde{z}_b^\star = \sqrt{\lambda_1} u_1$ and $u_1^T \tilde{z}(0) > 0$.

Based on the proof of Lemma 4, $\tilde{z}_b^\star$ is an eigenvector of $W$ with the corresponding eigenvalue $\|\tilde{z}_b^\star\|^2 = \alpha^2 + \beta^2$, where $\alpha$ and $\beta$ satisfy $\alpha\beta = \sigma_1$ and $\alpha^2 - \beta^2 = b$. Hence, we

can write

$$\lambda_1 = \sqrt{b^2 + 4\sigma_1^2} = \alpha^2 + \beta^2 = \|\tilde{z}_b^\star\|^2.$$

Therefore, $\tilde{z}_b^\star = \sqrt{\lambda_1} u_1$ holds. Moreover, since $(\xi_0, \eta_0) \in \mathcal{H}_b^+$, it follows that $(\tilde{z}_b^\star)^T \tilde{z}(0) > 0$ or equivalently $u_1^T \tilde{z}(0) > 0$, which completes the proof. ∎

Figure 2 illustrates the relationship between the trajectories of (1) and (4), over the invariant set $\mathcal{S}_{b=2}$. We consider $M = 1$ in (1) and $W = \begin{bmatrix} 2 & 2 \\ 2 & -2 \end{bmatrix}$ in (4). Notice that $W$ is the corresponding symmetric matrix of $M$ as defined in Theorem 3. The red thin dashed curves in this figure represent the set of non-zero equilibrium points of (1) given by $\{(x,y) \mid xy = 1\}$. The trajectories of (4) and (1) are marked by the blue dash-dotted and black curves, respectively, and the set $\mathcal{S}_b$ is marked by the blue curves. The red thick line represents the invariant set $\mathcal{H}$ with respect to (4) as established in Theorem 2. The red bullet points are the global minimizers $\pm(\xi_b^\star, \eta_b^\star)$. We observe that the trajectories of (4) and (1) traverse the same curve for any initial condition $z_0 = [\, x_0^T \ y_0^T \,]^T$ with $(x_0, y_0) \in \mathcal{S}_b$.
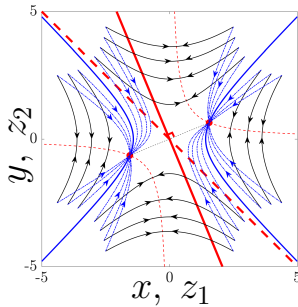


Fig. 2: An illustration of the trajectories of (1) (black curves) and (4) (blue dash-dotted curves) for a scalar $M$ and the corresponding matrix $W$ as established in Theorem 3.

## V. CONCLUDING REMARKS

The gradient flow dynamics associated with the problem of finding the best rank-one approximation of a given matrix has infinitely many stable and unstable equilibrium points. We partition the state-space into a family of invariant manifolds that each contains two stable equilibrium points. We show that, over each of these manifolds, the gradient flow simplifies to the gradient flow for the symmetric rank-one approximation problem. For each stable equilibrium point, this allows us to explicitly characterize the region of attraction as the domain of a radially unbounded Lyapunov function that also goes to infinity at the boundary. Our result establishes the almost everywhere convergence of gradient flow dynamics to the minimizers of the corresponding rank-one approximation problem. Our ongoing work focuses on extending the analysis to higher rank matrix approximation problems and problems with incomplete data.

## REFERENCES

[1] R. Ge, J. D. Lee, and T. Ma, "Matrix completion has no spurious local minimum," in *Advances in Neural Information Processing Systems*, 2016, pp. 2973–2981.

[2] A. Blum and R. L. Rivest, "Training a 3-node neural network is NP-complete," in *Advances in Neural Information Processing Systems*, 1989, pp. 494–501.

[3] E. J. Candes, X. Li, and M. Soltanolkotabi, "Phase retrieval via Wirtinger flow: Theory and algorithms," *IEEE Transactions on Information Theory*, vol. 61, no. 4, pp. 1985–2007, 2015.

[4] H. Sahinoglou and S. D. Cabrera, "On phase retrieval of finite-length sequences using the initial time sample," *IEEE Transactions on Circuits and Systems*, vol. 38, no. 8, pp. 954–958, 1991.

[5] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," *arXiv preprint arXiv:1611.03530*, 2016.

[6] K. Kawaguchi, "Deep learning without poor local minima," in *Advances in Neural Information Processing Systems*, 2016, pp. 586–594.

[7] D. Soudry and E. Hoffer, "Exponentially vanishing sub-optimal local minima in multilayer neural networks," *arXiv preprint arXiv:1702.05777*, 2017.

[8] H. Lu and K. Kawaguchi, "Depth creates no bad local minima," *arXiv preprint arXiv:1702.08580*, 2017.

[9] J. D. Lee, M. Simchowitz, M. I. Jordan, and B. Recht, "Gradient descent only converges to minimizers," in *Conference on Learning Theory*, 2016, pp. 1246–1257.

[10] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang, "On large-batch training for deep learning: Generalization gap and sharp minima," *arXiv preprint arXiv:1609.04836*, 2016.

[11] S. Bhojanapalli, B. Neyshabur, and N. Srebro, "Global optimality of local search for low rank matrix recovery," in *Advances in Neural Information Processing Systems*, 2016, pp. 3873–3881.

[12] B. Neyshabur, R. Tomioka, and N. Srebro, "In search of the real inductive bias: On the role of implicit regularization in deep learning," *arXiv preprint arXiv:1412.6614*, 2014.

[13] S. Gunasekar, B. Woodworth, S. Bhojanapalli, B. Neyshabur, and N. Srebro, "Implicit regularization in matrix factorization," *arXiv preprint arXiv:1705.09280*, 2017.

[14] N. Srebro and T. Jaakkola, "Weighted low-rank approximations," in *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, 2003, pp. 720–727.

[15] U. Helmke and J. B. Moore, *Optimization and dynamical systems*. Springer Science & Business Media, 2012.

[16] G. H. Golub and C. F. Van Loan, *Matrix computations*. JHU Press, 2012, vol. 3.

[17] P. Jain, C. Jin, S. M. Kakade, P. Netrapalli, and A. Sidford, "Streaming PCA: Matching matrix Bernstein and near-optimal finite sample guarantees for Oja's algorithm," in *Conference on Learning Theory*, 2016, pp. 1147–1164.

[18] E. Oja, "Simplified neuron model as a principal component analyzer," *Journal of Mathematical Biology*, vol. 15, no. 3, pp. 267–273, 1982.

[19] E. Oja, "Neural networks, principal components, and subspaces," *International Journal of Neural Systems*, vol. 1, no. 1, pp. 61–68, 1989.

[20] J. H. Manton, U. Helmke, and I. M. Mareels, "A dual purpose principal and minor component flow," *Syst. Control Lett.*, vol. 54, no. 8, pp. 759–769, 2005.

[21] C. Webers, J. H. Manton, *et al.*, "A generalisation of the Oja subspace flow," in *17th International Symposium on Mathematical Theory of Networks and Systems*, 2006, pp. 24–28.

[22] X. Kong, C. Hu, and C. Han, "A dual purpose principal and minor subspace gradient flow," *IEEE Transactions on Signal Processing*, vol. 60, no. 1, pp. 197–210, 2012.

[23] L. Liu, R. Ge, J. Meng, and G. You, "Dual subspace learning via geodesic search on Stiefel manifold," *International Journal of Machine Learning and Cybernetics*, vol. 5, no. 5, pp. 753–759, 2014.

[24] A. Vannelli and M. Vidyasagar, "Maximal Lyapunov functions and domains of attraction for autonomous nonlinear systems," *Automatica*, vol. 21, no. 1, pp. 69 – 80, 1985.