

---

# On Constrained Nonconvex Stochastic Optimization: A Case Study for Generalized Eigenvalue Decomposition

---

Zhehui Chen<sup>1</sup>  
Georgia Tech<sup>1</sup>

Xingguo Li<sup>2</sup>

Lin Yang<sup>2</sup>  
Princeton University<sup>2</sup>

Jarvis Haupt<sup>3</sup>  
University of Minnesota<sup>3</sup>

Tuo Zhao<sup>1</sup>

## Abstract

We study constrained nonconvex optimization problems in machine learning and signal processing. It is well-known that these problems can be rewritten to a min-max problem in a Lagrangian form. However, due to the lack of convexity, their landscape is not well understood and how to find the stable equilibria of the Lagrangian function is still unknown. To bridge the gap, we study the landscape of the Lagrangian function. Further, we define a special class of Lagrangian functions. They enjoy the following two properties: 1. Equilibria are either stable or unstable (Formal definition in Section 2); 2. Stable equilibria correspond to the global optima of the original problem. We show that a generalized eigenvalue (GEV) problem, including canonical correlation analysis and other problems as special examples, belongs to the class. Specifically, we characterize its stable and unstable equilibria by leveraging an invariant group and symmetric property (more details in Section 3). Motivated by these neat geometric structures, we propose a simple, efficient, and stochastic primal-dual algorithm solving the online GEV problem. Theoretically, under sufficient conditions, we establish an asymptotic rate of convergence and obtain the first sample complexity result for the online GEV problem by diffusion approximations, which are widely used in applied probability. Numerical results are also provided to support our theory.

---

Proceedings of the 22<sup>nd</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2019, Naha, Okinawa, Japan. PMLR: Volume 89. Copyright 2019 by the author(s).

## 1 Introduction

We often encounter the following optimization problem in machine learning and signal processing:

$$\min_X f(X) \quad \text{subject to} \quad X \in \Omega, \quad (1)$$

where  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is a loss function,  $\Omega \triangleq \{X \in \mathbb{R}^d : g_i(X) = 0, i = 1, 2, \dots, m\}$  denotes a feasible set,  $m$  is the number of constraints, and  $g_i : \mathbb{R}^d \rightarrow \mathbb{R}$ 's are the differentiable functions that impose constraints into model parameters. For notational simplicity, we define  $\mathcal{G}(X) = [g_1(X), \dots, g_m(X)]^\top$  and  $\Omega = \{X \in \mathbb{R}^d : \mathcal{G}(X) = 0\}$ . Principal component analysis (PCA), canonical correlation analysis (CCA), matrix factorization/sensing/completion, phase retrieval, and many other problems (Friedman et al., 2001; Sun et al., 2016; Bhojanapalli et al., 2016; Li et al., 2016b; Ge et al., 2016b; Chen et al., 2017; Zhu et al., 2017) can be viewed as special examples of (1). Many algorithms have been proposed to solve (1). For the unconstrained ( $\Omega = \mathbb{R}^d$ ) or a simple constraint  $\mathcal{G}(X)$ , e.g., the spherical constraint,  $\mathcal{G}(X) := \|X\|_2 - 1$ , we can apply simple first order algorithms such as the projected gradient descent algorithm (Luenberger et al., 1984).

However, when  $\mathcal{G}(X)$  is complicated, the aforementioned algorithms are often not applicable or inefficient. This is because the projection to  $\Omega$  does not admit a closed form expression and can be computationally expensive in each iteration. To address this issue, we convert (1) to a min-max problem using the Lagrangian multiplier method. Specifically, instead of solving (1), we solve the following:

$$\min_{X \in \mathbb{R}^d} \max_{Y \in \mathbb{R}^m} \mathcal{L}(X, Y) := f(X) + Y^\top \mathcal{G}(X), \quad (2)$$

where  $Y \in \mathbb{R}^m$  is the Lagrangian multiplier.  $\mathcal{L}(X, Y)$  is often referred as the Lagrangian function in literature (Boyd and Vandenberghe, 2004). The literature on optimization also refers to  $X$  as

the primal variable and  $Y$  as the dual variable. Accordingly, (1) is called the primal problem. From the perspective of game theory, it can be viewed as two players competing with each other and eventually achieving some equilibrium. When  $f(X)$  is convex and  $\Omega$  is convex or the boundary of a convex set, the optimization landscape of (2) is essentially convex-concave, i.e., for any fixed  $Y$ ,  $\mathcal{L}(X, Y)$  is convex in  $X$ , and for any fixed  $X$ ,  $\mathcal{L}(X, Y)$  is concave in  $Y$ . Such a landscape further implies that the equilibrium of (2) is a saddle point, whose primal variable is equivalent to the global optimum of (1) under strong duality conditions. To solve (2), we resort to primal-dual algorithms, which iterate over both  $X$  and  $Y$  (usually in an alternating manner). The global convergence rates to the equilibrium are also established accordingly for these algorithms (Lan et al., 2011; Chen et al., 2014; Iouditski and Nesterov, 2014).

When  $f(X)$  and  $\Omega$  are nonconvex, both (1) and (2) become much more computationally challenging, NP-Hard in general. Significant progress has been made toward solving primal problem (1). For example, Ge et al. (2015) show that when certain tensor factorization problem satisfies the so-called strict saddle properties, one can apply some first order algorithms, e.g., the projected gradient algorithm, and the global convergence in polynomial time can be guaranteed. Their results further motivate many follow-up works, proving that many problems can be formulated as strict saddle optimization problems, including PCA, multiview learning, phase retrieval, and matrix factorization/sensing/completion (Sun et al., 2016; Bhojanapalli et al., 2016; Li et al., 2016b; Ge et al., 2016b; Chen et al., 2017; Zhu et al., 2017; Liu et al., 2018). Note that these strict saddle optimization problems are either unconstrained or with a simple constraint. However, for many other nonconvex optimization problems,  $\Omega$  can be much more complicated. To the best of our knowledge, when  $\Omega$  is nonconvex and complicated, the applicable algorithms and convergence guarantees are still largely unknown.

To handle the complicated  $\Omega$ , this paper proposes to investigate min-max problem (2). Specifically, we first define a special class of Lagrangian functions, where the landscape of  $\mathcal{L}(X, Y)$  enjoys the following good properties: (1) *There exist only two types of equilibria – stable and unstable equilibria. At an unstable one,  $\mathcal{L}(X, Y)$  has negative curvature with respect to the primal variable  $X$ . More details in Section 2;* (2) *All stable equilibria correspond to the global optima of primal problem (1).*

Both properties are intuitive. On the one hand, the negative curvature in the first property enables the primal variable to escape from the unstable equilibria along some descent direction. On the other hand, the second property ensures that we do not get spurious local optima of (1), that is, all local minima must also be global optima.

We then study a generalized eigenvalue (GEV) problem, which includes CCA, Fisher discriminant analysis (FDA, Mika et al. (1999)), and sufficient dimension reduction (SDR, Cook and Ni (2005)) as special examples. Specifically, GEV solves

$$\begin{aligned} X^* &= \operatorname{argmin}_{X \in \mathbb{R}^{d \times r}} f(X) := -\operatorname{tr}(X^\top A X) \\ \text{s.t. } X &\in \mathcal{T}_B := \{X \in \mathbb{R}^{d \times r} : X^\top B X = I_r\}, \end{aligned} \quad (3)$$

where  $A, B \in \mathbb{R}^{d \times d}$  are symmetric,  $B$  is positive semidefinite. We rewrite (3) as a min-max problem,

$$\min_X \max_Y \mathcal{L}(X, Y) = -\operatorname{tr}(X^\top A X) + \langle Y, X^\top B X - I_r \rangle, \quad (4)$$

where  $Y \in \mathbb{R}^{r \times r}$  is the Lagrangian multiplier. Theoretically, we show that the Lagrangian function in (4) exactly belongs to our previously defined class. Motivated by our defined landscape structures, we then solve an online version of (4), where we can only access independent unbiased stochastic approximations of  $A$ ,  $B$  and directly accessing  $A$  and  $B$  is prohibited. Specifically, at the  $k$ -th iteration, we only obtain independent  $A^{(k)}$  and  $B^{(k)}$  satisfying

$$\mathbb{E}A^{(k)} = A \quad \text{and} \quad \mathbb{E}B^{(k)} = B.$$

Computationally, we propose a simple stochastic primal-dual algorithm, which is a stochastic variant of the generalized Hebbian algorithm (GHA, Gorrell (2006)). Theoretically, we establish its asymptotic rate of convergence to stable equilibria for our stochastic GHA (SGHA) based on the diffusion approximations (Kushner and Yin, 2003). Specifically, we show that, asymptotically, the solution trajectory of SGHA weakly converges to the solutions of stochastic differential equations (SDEs). By studying the analytical solutions of these SDEs, we further establish the asymptotic sample/iteration complexity of SGHA under certain regularity conditions (Harold et al., 1997; Li et al., 2016a; Chen et al., 2017). To the best of our knowledge, this is the first asymptotic sample/iteration complexity analysis of a stochastic optimization algorithm for solving the online version of GEV problem. Numerical experiments are also presented to justify our theory.

Our work is closely related to several recent results on solving GEV problems. For example, Ge et al.

(2016a) propose a multistage semi-stochastic optimization algorithm for solving the GEV problem with a finite sum structure. At each optimization stage, their algorithm needs to access the exact  $B$  matrix, and compute the approximate inverse of  $B$  by solving a quadratic program, which is forbidden in our setting. Similar matrix inversion approaches are also adopted by a few other recently proposed algorithms for solving the GEV problem (Allen-Zhu and Li, 2016; Arora et al., 2017). In contrast, our proposed SGHA is a fully stochastic algorithm, which does not require any matrix inversion.

Moreover, our work is also related to several more complicated min-max problems, such as Markov Decision Process with function approximation and Generative Adversarial Network (Sutton et al., 2000; Shapiro et al., 2009; Goodfellow et al., 2014). Many primal-dual algorithms have been proposed to solve these problems. However, most of these algorithms are not guaranteed to converge. As mentioned earlier, when the convex-concave structure is missing, the min-max problems go far beyond the existing theories. Moreover, both primal and dual iterations involve sophisticated stochastic approximations (equal or more difficult than our online GEV). This paper makes the attempt on understanding the optimization landscape of these challenging min-max problems. Taking our results as an initial start, we expect more sophisticated and stronger follow-up works that apply to these min-max problems.

**Notations.** Given an integer  $d$ , we denote  $I_d$  as a  $d \times d$  identity matrix,  $[d] = \{1, 2, \dots, d\}$ . Given an index set  $\mathcal{I} \subseteq [d]$  and a matrix  $X \in \mathbb{R}^{d \times r}$ , we denote  $\mathcal{I}^\perp = [d] \setminus \mathcal{I}$  as the complement set of  $\mathcal{I}$ ,  $X_{:,i}$  ( $X_{i,:}$ ) as the  $i$ -th column (row) of  $X$ ,  $X_{i,j}$  as the  $(i, j)$ -th entry of  $X$ ,  $X_{:, \mathcal{I}}$  ( $X_{\mathcal{I}, :}$ ) as the column (row) submatrix of  $X$  indexed by  $\mathcal{I}$ ,  $\text{vec}(X) \in \mathbb{R}^{dr}$  as the vectorization of  $X$ ,  $\text{Col/Null}(X)$  as the column/null space of  $X$ . Given a symmetric matrix  $X \in \mathbb{R}^{d \times d}$ , we denote  $\lambda_{\min/\max}(X)$  as its smallest/largest singular value, and denote the eigenvalue decomposition of  $X$  as  $X = O\Lambda O^\top$ , where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$  with  $\lambda_1 \geq \dots \geq \lambda_d$ , denote  $\|X\|_2$  as the spectral norm of  $X$ . Given two matrices  $X$  and  $Y$ ,  $X \otimes Y$  as the Kronecker product of  $X, Y$ .

## 2 Characterization of Equilibria

We start with characterizing the equilibria of (2). By KKT conditions, an equilibrium  $(X, Y)$  satisfies

$$\begin{aligned} \nabla_X \mathcal{L}(X, Y) &= \nabla_X f(X) + Y^\top \nabla_X \mathcal{G}(X) = 0 \\ \text{and } \nabla_Y \mathcal{L}(X, Y) &= \mathcal{G}(X) = 0, \end{aligned}$$

which only contains the first order information of  $\mathcal{L}(X, Y)$ . To further distinguish the difference among the equilibria, we define two types of equilibria by the second order information.

**Definition 1.** Given the Lagrangian function  $\mathcal{L}(X, Y)$  in (2), a point  $(X, Y)$  is called:

(1) An **equilibrium** of  $\mathcal{L}(X, Y)$ , if

$$\nabla \mathcal{L}(X, Y) = \begin{bmatrix} \nabla_X \mathcal{L}(X, Y) \\ \nabla_Y \mathcal{L}(X, Y) \end{bmatrix} = 0.$$

(2) An **equilibrium**  $(X, Y)$  is **unstable**, if  $(X, Y)$  is an equilibrium and  $\lambda_{\min}(\nabla_X^2 \mathcal{L}(X, Y)) < 0$ .

(3) An **equilibrium**  $(X, Y)$  is **stable**, if  $(X, Y)$  is an equilibrium,  $\nabla_X^2 \mathcal{L}(X, Y) \succeq 0$ , and  $\mathcal{L}(X, Y)$  is strongly convex over a restricted domain.

Note that (2) in Definition 1 has a similar strict saddle property over a manifold in Ge et al. (2015). The motivation behind Definition 1 is intuitive. When  $\mathcal{L}(X, Y)$  has negative curvature with respect to the primal variable  $X$  at an equilibrium, we can find a direction in  $X$  to further decrease  $\mathcal{L}(X, Y)$ . Therefore, a tiny perturbation can break this unstable equilibrium. An illustrative example is presented in Figure 1. Moreover, at a stable equilibrium  $(X^*, Y^*)$ , there is restricted strong convexity, which relates to several conditions, e.g., Polyak Lojasiewicz conditions (Polyak, 1963), i.e.,

$$\|\nabla_X \mathcal{L}(X, Y^*)\|^2 \geq \mu(\mathcal{L}(X, Y^*) - \mathcal{L}(X^*, Y^*)),$$

for  $X$  belonging to a small region near  $X^*$  and  $\mu > 0$  is a constant, or Error Bound conditions (Luo and Tseng, 1993). With this property, we cannot decrease  $\mathcal{L}(X, Y)$  along any direction with respect to  $X$ . Definition 1 excludes the high order unstable equilibrium, which may exist due to the degeneracy of  $\nabla_X^2 \mathcal{L}(X, Y)$ . Specifically, such a high order unstable equilibrium cannot be identified by the second order information, e.g.,

$$\mathcal{L}(x_1, x_2, y) = x_1^3 + x_2^2 + y \cdot (x_1 - x_2).$$

$(0, 0, 0)$  is an equilibrium with a positive semidefinite Hessian matrix. However, it is an unstable equilibrium, since a small perturbation to  $x_1$  can break this equilibrium. Such an equilibrium makes the landscape highly more complicated. Overall, we consider a specific class of Lagrangian functions throughout the rest of this paper. They enjoy the following properties: (1) All equilibria are either stable or unstable (i.e., no high order unstable equilibria); (2) All stable equilibria correspond to the global optima of the primal problem.

As mentioned earlier, the first property ensures that second order information can identify the type of

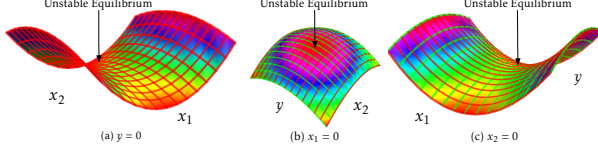


Figure 1: An illustration of an unstable equilibrium:  $\min_{x_1, x_2} \max_y \mathcal{L}(x_1, x_2, y) = x_1^2 - x_2^2 - y^2$ . Notice that  $(0, 0, 0)$  is an unstable equilibrium. For visualization, we show three views: (a)  $\mathcal{L}(x_1, x_2, 0)$ ; (b)  $\mathcal{L}(0, x_2, y)$ ; (c)  $\mathcal{L}(x_1, 0, y)$ . Red lines correspond to  $x_1$  and  $x_2$ , and the green one corresponds to the  $y$ .

equilibria. The second property guarantees that we do not get spurious optima for (1) as long as an algorithm attains a stable equilibrium. Several machine learning problems belong to this class, such as generalized eigenvalue problem.

### 3 Generalized Eigenvalue Problem

We consider the generalized eigenvalue (GEV) problem as a motivating example, which includes CCA, FDA, SDR, etc. as special examples. Recall its min-max formulation (4):

$$\min_{X \in \mathbb{R}^{d \times r}} \max_{Y \in \mathbb{R}^{r \times r}} \mathcal{L}(X, Y),$$

where  $\mathcal{L}(X, Y) = -\text{tr}(X^\top A X) + \langle Y, X^\top B X - I_r \rangle$ .

Before we proceed, we impose the following assumption on the problem.

**Assumption 1.** Given a symmetric matrix  $A \in \mathbb{R}^{d \times d}$  and a positive definite matrix  $B \in \mathbb{R}^{d \times d}$ , the eigenvalues of  $\tilde{A} = B^{-\frac{1}{2}} A B^{-\frac{1}{2}}$ , denoted by  $\lambda_1^{\tilde{A}}, \dots, \lambda_d^{\tilde{A}}$ , satisfy

$$\lambda_1^{\tilde{A}} \geq \dots \geq \lambda_r^{\tilde{A}} > \lambda_{r+1}^{\tilde{A}} \geq \dots \geq \lambda_d^{\tilde{A}}.$$

Such an eigengap assumption avoids the identifiability issue. The full rank assumption on  $B$  in Assumption 1 ensures that the original constrained optimization problem is bounded. This assumption can be further relaxed but require more involved analysis. We will discuss this in Appendix B.

To characterize all equilibria of GEV, we leverage the idea of an invariant group. Li et al. (2016b) use similar techniques for an unconstrained matrix factorization problem. However, it does not work for the Lagrangian function due to the more complicated landscape. Therefore, we consider a more general invariant group. Moreover, by analyzing the Hessian matrix of  $\mathcal{L}(X, Y)$  at the equilibria, we demonstrate that each equilibrium is either unstable or stable and the stable equilibria correspond to the

global optima of the primal problem (3). Therefore, GEV belongs to the class we defined earlier.

#### 3.1 Invariant Group and Symmetric Property

Denote the orthogonal group in dimension  $r$  as

$$O(r, \mathbb{R}) = \{ \Psi \in \mathbb{R}^{r \times r} : \Psi \Psi^\top = \Psi^\top \Psi = I_r \}.$$

Notice that for any  $\Psi \in O(r, \mathbb{R})$ ,  $\mathcal{L}(X, Y)$  in (4) has the same landscape as  $\mathcal{L}(X\Psi, \Psi^\top Y\Psi)$ . This further indicates that given an equilibrium  $(X, Y)$ ,  $(X\Psi, \Psi^\top Y\Psi)$  is also an equilibrium. This symmetric property motivates us to characterize the equilibria of  $\mathcal{L}(X, Y)$  with an invariant group.

We introduce several important definitions in group theory (Dummit and Foote, 2004).

**Definition 2.** Given a group  $\mathcal{H}$  and a set  $\mathcal{X}$ , a map  $\phi(\cdot, \cdot)$  from  $\mathcal{H} \times \mathcal{X}$  to  $\mathcal{X}$  is called the **group action** of  $\mathcal{H}$  on  $\mathcal{X}$  if  $\phi$  satisfies the following two properties: **Identity:**  $\phi(\mathbf{1}, x) = x \quad \forall x \in \mathcal{X}$ , where  $\mathbf{1}$  denotes the identity element of  $\mathcal{H}$ .

**Compatibility:**  $\phi(gh, x) = \phi(g, \phi(h, x)) \quad \forall g, h \in \mathcal{H}, x \in \mathcal{X}$ .

**Definition 3.** Given a function  $f(x, y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ , a group  $\mathcal{H}$  is a **stationary invariant group** of  $f$  with respect to two group actions of  $\mathcal{H}$ ,  $\phi_1$  on  $\mathcal{X}$  and  $\phi_2$  on  $\mathcal{Y}$ , if  $\mathcal{H}$  satisfies

$$f(x, y) = f(\phi_1(g, x), \phi_2(g, y)) \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}, g \in \mathcal{H}.$$

For notational simplicity, we denote  $\mathcal{G} = O(r, \mathbb{R})$ . Given the group  $\mathcal{G}$ , two sets  $\mathbb{R}^{d \times r}$  and  $\mathbb{R}^{r \times r}$ , we define a group action with  $\phi_1$  of  $\mathcal{G}$  on  $\mathbb{R}^{d \times r}$  and a group action  $\phi_2$  of  $\mathcal{G}$  on  $\mathbb{R}^{r \times r}$  as

$$\phi_1(\Psi, X) = X\Psi \quad \forall \Psi \in \mathcal{G}, X \in \mathbb{R}^{d \times r}$$

$$\text{and } \phi_2(g, Y) = \Psi^{-1} Y \Psi \quad \forall \Psi \in \mathcal{G}, Y \in \mathbb{R}^{r \times r}.$$

One can check that the orthogonal group  $\mathcal{G}$  is a stationary invariant group of  $\mathcal{L}(X, Y)$  with respect to two group actions of  $\mathcal{G}$ ,  $\phi_1$  on  $\mathbb{R}^{d \times r}$  and  $\phi_2$  on  $\mathbb{R}^{r \times r}$ . By this invariant group, we define the equivalence relation between  $(X_1, Y_1)$  and  $(X_2, Y_2)$ , if there exists a  $\Psi \in \mathcal{G}$  such that

$$(X_1, Y_1) = (X_2 \Psi, \Psi^{-1} Y_2 \Psi) = (X_2 \Psi, \Psi^\top Y_2 \Psi). \quad (5)$$

To find all equilibria of GEV, we examine the KKT conditions of (4):

$$\begin{aligned} 2BXY - 2AX = 0 \quad \text{and} \quad X^\top BX - I_r = 0 \\ \implies Y = X^\top AX =: \mathcal{D}(X). \end{aligned}$$

Given the eigendecomposition  $B = O^B \Lambda^B O^{B^\top}$ , we denote  $\tilde{A} = (\Lambda^B)^{-\frac{1}{2}} O^{B^\top} A O^B (\Lambda^B)^{-\frac{1}{2}}$  and  $\tilde{X} =$



$(\Lambda^B)^{\frac{1}{2}}O^{B\top}X$ . We then consider eigendecomposition  $\tilde{A} = O^{\tilde{A}}\Lambda^{\tilde{A}}O^{\tilde{A}\top}$ . The following theorem shows the connection between the equilibrium of  $\mathcal{L}(X, Y)$  and the column submatrix of  $O^{\tilde{A}}$ , denoted as  $O_{:, \mathcal{I}}^{\tilde{A}}$ , where

$$\mathcal{I} \in \mathcal{X}_d^r := \left\{ \{i_1, \dots, i_r\} : \{i_1, \dots, i_r\} \subseteq [d] \right\}$$

is the index set to determine a column submatrix.

**Theorem 4** (Symmetric Property). *Suppose Assumption 1 holds. Then  $(X, \mathcal{D}(X))$  is an equilibrium of  $\mathcal{L}(X, Y)$ , if and only if  $X$  can be written as*

$$X = (O^B(\Lambda^B)^{-\frac{1}{2}}O_{:, \mathcal{I}}^{\tilde{A}}) \cdot \Psi,$$

where index  $\mathcal{I} \in \mathcal{X}_d^r$  and  $\Psi \in \mathcal{G}$ .

The proof of Theorem 4 is provided in Appendix A.1. Theorem 4 implies that under the equivalence relation given in (5), there are  $\binom{d}{r}$  equilibria of  $\mathcal{L}(X, Y)$ . Each corresponds to an  $O_{:, \mathcal{I}}^{\tilde{A}}$ , where  $\mathcal{I} \in \mathcal{X}_d^r$  is an index set. Then whole equilibria set is generated by  $O_{:, \mathcal{I}}^{\tilde{A}}$ 's with the transformation matrix  $O^B(\Lambda^B)^{-\frac{1}{2}}$  and the invariant group action induced by  $\mathcal{G}$ .

### 3.2 Unstable and Stable Equilibria

We further identify the stable and unstable equilibria. Specifically, given  $(X, Y)$  as an equilibrium of  $\mathcal{L}(X, Y)$ , we denote the Hessian matrix of  $\mathcal{L}(X, Y)$  with respect to the primal variable  $X$  as

$$H_X \triangleq \nabla_X^2 \mathcal{L}(X, Y)|_{Y=\mathcal{D}(X)} \in \mathbb{R}^{dr \times dr}.$$

Then we calculate the eigenvalues of  $H_X$ . By Definition 1,  $(X, \mathcal{D}(X))$  is unstable if  $H_X$  has a negative eigenvalue; Otherwise, we analyze the local landscape at  $(X, \mathcal{D}(X))$  to determine whether it is stable or not. The following theorem shows that all equilibria are either stable or unstable and demonstrates how the choice of index set  $\mathcal{I}$  corresponds to the unstable and stable equilibria of  $\mathcal{L}(X, Y)$ .

**Theorem 5.** *Suppose Assumption 1 holds, and  $(X, \mathcal{D}(X))$  is an equilibrium in (4). By Theorem 4,  $X$  can be represented as  $X = (O^B(\Lambda^B)^{-\frac{1}{2}}O_{:, \mathcal{I}}^{\tilde{A}}) \cdot \Psi$  for some  $\Psi \in \mathcal{G}$  and  $\mathcal{I} \in \mathcal{X}_d^r$ . Then, if  $\mathcal{I} \neq [r]$ , equilibrium  $(X, \mathcal{D}(X))$  is unstable with*

$$\lambda_{\min}(H_X) \leq \frac{2(\lambda_{\max \mathcal{I}}^{\tilde{A}} - \lambda_{\min \mathcal{I}^\perp}^{\tilde{A}})}{\|X_{:, \min \mathcal{I}^\perp}\|_2^2} < 0,$$

where  $\lambda_{\max / \min \mathcal{I}}^{\tilde{A}} = \max / \min_{i \in \mathcal{I}} \lambda_i^{\tilde{A}}$ , and  $\lambda_i^{\tilde{A}}$  is the  $i$ -th leading eigenvalue of  $\tilde{A}$ ;

Otherwise, we have  $H_X \succeq 0$  and  $\text{rank}(H_X) = dr - r(r-1)/2$ . Moreover,  $(X, \mathcal{D}(X))$  is a stable equilibrium of problem (4).

The proof of Theorem 5 is provided in Appendix A.2. Theorem 5 indicates that when  $\tilde{X} = O_{:, [r]}^{\tilde{A}}$ , that is, the eigenvectors of  $\tilde{A}$  corresponding to the  $r$  largest eigenvalues,  $(X, \mathcal{D}(X))$  is a stable equilibrium of  $\mathcal{L}(X, Y)$ , where  $X = (O^B(\Lambda^B)^{-\frac{1}{2}}O_{:, \mathcal{I}}^{\tilde{A}}) \cdot \Psi$  for some  $\Psi \in \mathcal{G}$ . Although  $H_X$  is degenerate at this equilibrium, all directions in  $\text{Null}(H_X)$  essentially point to the primal variables of other stable equilibria. Excluding these directions, the rest all have positive curvature, which implies that this equilibrium is stable. Moreover, such an  $X$  corresponds to the optima of (3). When  $\mathcal{I} \neq [r]$ , due to the negative curvature, these equilibria are unstable. Therefore, all stable equilibria of  $\mathcal{L}(X, Y)$  correspond to the global optima in (3) and other equilibria are unstable, which further indicates that GEV belongs to the class we defined earlier.

## 4 Stochastic Optimization for GEV

For GEV, we propose a fully stochastic primal-dual algorithm to solve (4), which only requires access to the stochastic approximations (SA) of  $A$  and  $B$  matrices. This is very different from other existing semi-stochastic algorithms that require to access the exact  $B$  matrix (Ge et al., 2016a). Specifically, we propose a stochastic variant of the generalized Hebbian algorithm (GHA), also referred as Sanger's rule in existing literature (Sanger, 1989), to solve (4). For online setting, accessing the exact  $A$  and  $B$  is prohibitive and we only get  $A^{(k)} \in \mathbb{R}^{d \times d}$  and  $B^{(k)} \in \mathbb{R}^{d \times d}$  that are independently sampled from the distribution associated with  $A$  and  $B$  at the  $k$ -th iteration. Our proposed SGHA updates primal and dual variables as follows:

Primal Update:

$$X^{(k+1)} \leftarrow X^{(k)} - \eta \underbrace{\left( B^{(k)} X^{(k)} Y^{(k)} - A^{(k)} X^{(k)} \right)}_{\text{Stoc. Approx. of } \nabla_X \mathcal{L}(X^{(k)}, Y^{(k)})}, \quad (6)$$

Dual Update:

$$Y^{(k+1)} \leftarrow \underbrace{X^{(k)\top} A^{(k)} X^{(k)}}_{\text{Stoc. Approx. of } X^{(k)\top} A X^{(k)}}, \quad (7)$$

where  $\eta > 0$  is a step size parameter. Note that the primal update is a stochastic gradient descent step, while the dual update is motivated by the KKT conditions of (4). SGHA is simple and easy to implement. The constraint is handled by the dual update. Further, motivated by the landscape of GEV, we analyze the algorithm by diffusion approximations and obtain the asymptotical sample complexity.

## 4.1 Numerical Evaluations

We first provide numerical evaluations to illustrate the effectiveness of SGHA, and then provide an asymptotic convergence analysis of SGHA. We choose  $d = 500$  and select three different settings:

**Setting(1)** :  $\eta = 10^{-4}$ ,  $r = 1$ ,  $A_{ii} = 1/100 \ \forall i \in [d]$ ,  $A_{ij} = 0.5/10$  and  $B_{ij} = 0.5^{|i-j|}/3 \ \forall i \neq j$ ;

**Setting(2)** :  $\eta = 5 \times 10^{-5}$ ,  $r = 3$ , and randomly generate an orthogonal matrix  $U \in \mathbb{R}^{d \times d}$  such that  $A = U \cdot \text{diag}(1, 1, 1, 0.1, \dots, 0.1) \cdot U^\top$  and  $B = U \cdot \text{diag}(2, 2, 2, 1, \dots, 1) \cdot U^\top$ ;

**Setting(3)** :  $\eta = 2.5 \times 10^{-5}$ ,  $r = 3$ , and randomly generate two orthogonal matrices  $U, V \in \mathbb{R}^{d \times d}$  such that  $A = U \cdot \text{diag}(1, 1, 1, 0.1, \dots, 0.1) \cdot U^\top$  and  $B = V \cdot \text{diag}(2, 2, 2, 1, \dots, 1) \cdot V^\top$ .

At the  $k$ -th iteration of SGHA, we independently sample 40 random vectors from  $N(0, A)$  and  $N(0, B)$  respectively. Accordingly, we compute the sample covariance matrices  $A^{(k)}$  and  $B^{(k)}$  as the approximations of  $A$  and  $B$ . We repeat numerical simulations under each setting for 20 times using random data generations, and present all results in Figure 2. The horizontal and vertical axes correspond to the number of iterations and the optimization error

$$\|B^{1/2}X^{(t)}X^{(t)\top}B^{1/2} - B^{1/2}X^*X^{*\top}B^{1/2}\|_F$$

, respectively. Our experiments indicate that SGHA converges to a global optimum in all settings.

## 4.2 Analysis for Commutative $A$ and $B$

As a special case, we first prove the convergence of SGHA for GEV with  $r = 1$ , and commutative  $A$  and  $B$ . We will discuss more on noncommutative cases and  $r > 1$  in the next section. Before we proceed, we introduce our assumptions on the problem.

**Assumption 2.** *We assume that the following conditions hold:*

(a):  $A^{(k)}$ 's and  $B^{(k)}$ 's are independently sampled from two different distributions  $\mathcal{D}_A$  and  $\mathcal{D}_B$  respectively, where  $\mathbb{E}A^{(k)} = A$  and  $\mathbb{E}B^{(k)} = B \succ 0$ ;

(b):  $A$  and  $B$  are commutative, i.e., there exists an orthogonal matrix  $O$  such that  $A = O\Lambda^A O^\top$  and  $B = O\Lambda^B O^\top$ , where  $\Lambda^B = \text{diag}(\mu_1, \dots, \mu_d)$  and  $\Lambda^A = \text{diag}(\lambda_1, \dots, \lambda_d)$  with  $\lambda_j \neq 0$ ;

(c):  $A^{(k)}$  and  $B^{(k)}$  satisfy the moment conditions, that is, for some generic constants  $C_0$  and  $C_1$ ,  $\mathbb{E}\|A^{(k)}\|_2^2 \leq C_0$  and  $\mathbb{E}\|B^{(k)}\|_2^2 \leq C_1$ .

Note that (a) and (c) in Assumption 2 are mild, but (b) is stringent. For convenience of analysis, we combine (6) and (7) as

$$X^{(k+1)} \leftarrow X^{(k)} - \eta (B^{(k)} X^{(k)} X^{(k)\top} - I_d) A^{(k)} X^{(k)}. \quad (8)$$

We remark that (8) is very different from existing optimization algorithms over the generalized Stiefel manifold. Specifically, computing the gradient over the generalized Stiefel manifold requires  $B^{-1}$ , which is not allowed in our setting. For notational convenience, we further denote

$$\Lambda = (\Lambda^B)^{-\frac{1}{2}} \Lambda^A (\Lambda^B)^{-\frac{1}{2}} : \text{diag}(\beta_1, \dots, \beta_d).$$

Without loss of generality, we assume  $\beta_1 > \beta_2 \geq \beta_3 \geq \dots \geq \beta_d$ , and  $\beta_i \neq 0 \ \forall i \in [d]$ . However,  $\mu_i$  and  $\lambda_i$  are not necessarily to be monotonic. Denote  $\mu_{\min} = \min_{i \neq 1} \mu_i$ ,  $\mu_{\max} = \max_{i \neq 1} \mu_i$ , and  $\text{gap} = \beta_1 - \beta_2$ .

Moreover, we denote  $W^{(k)} = (\Lambda^B)^{\frac{1}{2}} O X^{(k)}$ . One can verify that (8) can be rewritten as follows:

$$W^{(k+1)} \leftarrow W^{(k)} - \eta \left( (\Lambda^B)^{\frac{1}{2}} \widehat{\Lambda}_B^{(k)} (\Lambda^B)^{-\frac{1}{2}} \cdot W^{(k)} W^{(k)\top} - \Lambda^B \right) \cdot \widetilde{\Lambda}^{(k)} W^{(k)}, \quad (9)$$

where  $\widehat{\Lambda}_B^{(k)} = O^\top B^{(k)} O$  and  $\widetilde{\Lambda}^{(k)} = O^\top B^{-\frac{1}{2}} \Lambda^{(k)} B^{-\frac{1}{2}} O$ . Note that  $W^* = [1, 0, 0, \dots, 0]^\top$  corresponds to the optimal solution of (3).

By diffusion approximation, we show that our algorithm converges through three Phases:

**Phase I:** Given an initial near a saddle point, we show that after rescaling of time properly, the algorithm can be characterized by a stochastic differential equation (SDE). Such an SDE further implies our algorithm can escape from the saddle fast;

**Phase II:** We show that away from the saddle, the trajectory of our algorithm can be approximated by an ordinary differential equation (ODE);

**Phase III:** We first show that after Phase II, the norm of solution converges to a constant. Then, the algorithm can be characterized by an SDE, like Phase I. By the SDE, we analyze the error fluctuation when the solution is within a small neighborhood of the global optimum.

Overall, we have an asymptotic sample complexity.

**ODE Characterization:** To demonstrate an ODE characterization for the trajectory of our algorithm, we introduce a continuous time random process

$$w^{(\eta)}(t) := W^{(k)},$$

where  $k = \lfloor \frac{t}{\eta} \rfloor$  and  $\eta$  is the step size in (8). For notational simplicity, we drop  $(t)$  when it is clear from the context. Instead of directly showing a global convergence of  $w^{(\eta)}$ , we construct a new quantity as

$$v_{i,j}^{(\eta)} = (w_i^{(\eta)})^{\mu_j} / (w_j^{(\eta)})^{\mu_i},$$

where  $w_i^{(\eta)}$  is the  $i$ -th component (coordinate) of  $w^{(\eta)}$ . We then show that  $v_{i,j}^{(\eta)}$  converges to an exponential decay function.

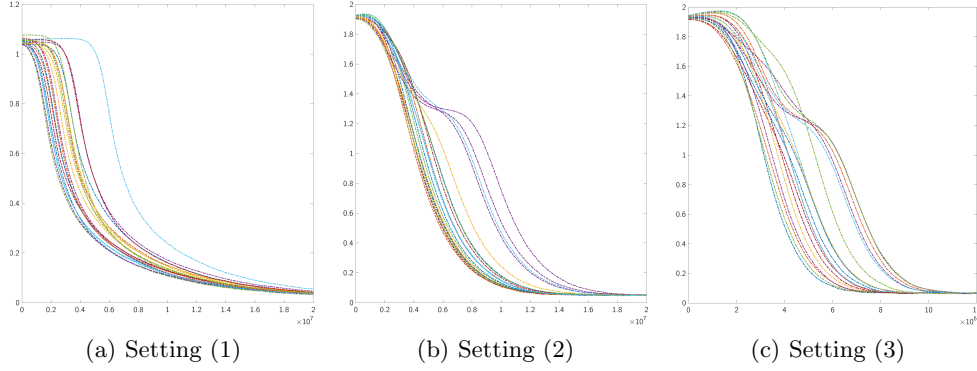


Figure 2: Plots of the optimization error  $\|B^{1/2}X^{(t)}X^{(t)\top}B^{1/2} - B^{1/2}X^*X^{*\top}B^{1/2}\|_F$  over SGHA iterations on synthetic data of 20 random data generations under different settings of parameters.

**Lemma 6.** *Suppose Assumption 2 holds and the initial is away from saddle points, i.e., given constants,  $\tau > 0$  and  $\delta < \frac{1}{2}$ , there exist  $i, j$  such that*

$$i \neq j, \quad |w_j^{(\eta)}| > \tau, \quad \text{and} \quad |w_i^{(\eta)}| > \eta^{\frac{1}{2} + \delta}.$$

As  $\eta \rightarrow 0$ ,  $v_{k,j}^{(\eta)}$  weakly converges to the solution of the following ODE:

$$dx_{k,j} = x_{k,j} \cdot (\mu_j \mu_k (\beta_k - \beta_j)) dt \quad \forall k \neq j. \quad (10)$$

The proof of Lemma 6 is provided in Appendix C.1. Lemma 6 essentially implies the global convergence of SGHA. Specifically, the solution of (10) is

$$x_{k,j}(t) = x_{k,j}(0) \cdot \exp(\mu_j \mu_k (\beta_k - \beta_j) t) \quad \forall k \neq j,$$

where  $x_{k,j}(0)$  is the initial value of  $v_{k,j}^{(\eta)}$ . In particular, we consider  $j = 1$ . Then, as  $t \rightarrow \infty$ , the dominating component of  $w$  will be  $w_1$ .

The ODE approximation of the algorithm implies that after long enough time, i.e.,  $t$  is large enough, the solution of the algorithm can be arbitrarily close to a global optimum. Nevertheless, to obtain the asymptotic ‘‘convergence rate’’, we need to study the variance of the trajectory at time  $t$ . Thus, we resort to the following SDE-based approach for a more precise characterization.

**SDE Characterization:** We notice that such a variance in the order of  $\mathcal{O}(\eta)$  vanishes as  $\eta \rightarrow 0$ . To characterize this variance, we rescale the updates by a factor of  $\eta^{-\frac{1}{2}}$ , i.e., by defining a new process as  $z^{(\eta)} = \eta^{-\frac{1}{2}} w^{(\eta)}$ . After rescaling, the variance of  $z^{(\eta)}$  is of the order of  $\mathcal{O}(1)$ . The following lemma characterizes how the algorithm escapes from the saddle, i.e.,  $w^{(\eta)}(0) \approx e_i$ , for  $i \neq 1$ , in Phase I.

**Lemma 7.** *Suppose Assumption 2 holds and the initial is close to a saddle point, i.e., given constants  $\delta < \frac{1}{2}$  and  $D$ , there exists an  $i \in [d] \setminus \{1\}$  such that*

$$|w_i^{(\eta)} - 1| \leq D\eta^{\frac{1}{2} + \delta} \quad \text{and} \quad |w_j^{(\eta)}| \leq D\eta^{\frac{1}{2} + \delta} \quad \forall j \neq i. \quad 7$$

As  $\eta \rightarrow 0$ , then for  $\forall j \neq i$   $z_j^{(\eta)}$  weakly converges to the solution of the following SDE:

$$dz_j(t) = (-\beta_j \mu_i \cdot z_i + \lambda_i z_i) dt + \sqrt{G_{j,i}} dB(t), \quad (11)$$

where  $G_{j,i} = \mathbb{E}\left(\left(\widehat{\Lambda}_B^{(k)}\right)_{j,i} \cdot \sqrt{\mu_j / \mu_i} \cdot \widetilde{\Lambda}_{i,i} - \mu_j \widetilde{\Lambda}_{j,i}\right)^2$  and  $B(t)$  is a standard Brownian motion.

The proof of Lemma 7 is provided in Appendix C.2. Note that (11) is a Fokker-Plank equation, whose solution is an Ornstein-Uhlenbeck (O-U) process (Doob, 1942) as follows:

$$z_j(t) = \exp[-\mu_j (\beta_i - \beta_j) t] \cdot \underbrace{\left[ z_j(0) + \sqrt{G_{j,i}} \int_0^t \exp[\mu_j (\beta_i - \beta_j) s] dB(s) \right]}_{Q_1}. \quad (12)$$

We consider  $j = 1$ . Note that  $Q_1$  is essentially a random variable with mean  $z_j(0)$  and variance smaller than  $\frac{G_{1,i} \mu_1}{2(\beta_1 - \beta_i)}$ . However, the larger  $t$  is, the closer its variance gets to this upper bound. Moreover, the term  $\exp[\mu_1 (\beta_1 - \beta_i) t]$  essentially amplifies  $Q_1$  by a factor exponentially increasing in  $t$ . This tremendous amplification forces  $z_1(t)$  to quickly get away from 0, as  $t$  increases, which indicates that the algorithm will escape from the saddle. Further, the following lemma characterizes the local behavior of the algorithm near the optimal.

**Lemma 8.** *Suppose that Assumption 2 holds and the initial solution is close to an optimal solution, that is, given pre-specified constants  $\kappa$  and  $\delta < \frac{1}{2}$ , we have  $\frac{|w_1^{(\eta)}|^2}{\|w^{(\eta)}\|_2^2} > 1 - \kappa\eta^{1+2\delta}$ . As  $\eta \rightarrow 0$ , then we have  $\|w^{(\eta)}(t)\|_2 \xrightarrow{t \rightarrow \infty} 1$  and for  $\forall i \neq 1$ ,  $z_i^{(\eta)}$  weakly converges to the solution of the following SDE:*

$$dz_i(t) = (-\beta_1 \cdot \mu_i z_i + \lambda_i z_i) dt + \sqrt{G_{i,1}} dB(t), \quad (13)$$

where  $G_{i,1} = \mathbb{E}((\widehat{\Lambda}_B)_{i,1} \cdot \sqrt{\mu_i/\mu_1} \cdot \widetilde{\Lambda}_{1,1} - \mu_i \Lambda_{i,1})^2$ , and  $B(t)$  is a standard Brownian motion.

The proof of Lemma 8 is provided in Appendix C.3. The solution of (13) is as follows:

$$z_i(t) = \sqrt{G_{i,1}} \int_0^t \exp[\mu_i(\beta_1 - \beta_i)(s-t)] dB(s) + z_i(0) \cdot \exp[-\mu_i(\beta_1 - \beta_i)t]. \quad (14)$$

Note the second term of the right hand side in (14) decays to 0, as time  $t \rightarrow \infty$ . The rest is a pure random walk. Thus, the fluctuation of  $z_i(t)$  is essentially the error fluctuation of the algorithm after sufficiently long time.

By Lemma 6, 7, and 8, we have the next theorem.

**Theorem 9.** *Suppose Assumption 2 holds. Given a sufficiently small error  $\epsilon > 0$ ,  $\phi = \sum_{i=1}^d G_{i,1}$ , and*

$$\eta \asymp \frac{\epsilon \cdot \mu_{\min} \cdot \text{gap}}{\phi},$$

we need

$$T \asymp \frac{\mu_{\max}/\mu_{\min}}{\mu_1 \cdot \text{gap}} \log(\eta^{-1}) \quad (15)$$

such that with probability at least  $\frac{5}{8}$ ,  $\|w(T) - W^*\|_2^2 \leq \epsilon$ , where  $W^*$  is the optima of (3).

The proof of Theorem 9 is provided in Appendix C.4. Theorem 9 implies that asymptotically, our algorithm yields an iterations of complexity:

$$N \asymp \frac{T}{\eta} \asymp \frac{\phi \cdot \mu_{\max}/\mu_{\min}}{\epsilon \cdot \mu_1 \cdot \mu_{\min} \cdot \text{gap}^2} \log\left(\frac{\phi}{\epsilon \cdot \mu_{\min} \cdot \text{gap}}\right),$$

which depends on the gap, i.e.,  $\beta_1 - \beta_2$ , but also depends on  $\frac{\mu_{\max}}{\mu_{\min}}$ , which is the condition number of  $B$  in the worst case. As can be seen, for an ill-conditioned  $B$ , problem (3) is more difficult to solve.

### 4.3 When $A$ and $B$ are Noncommutative?

Unfortunately, when  $A$  and  $B$  are noncommutative, the analysis is more difficult, even for  $r = 1$ . Recall that the optimization landscape of the Lagrangian function in (4) enjoys a nice geometric property: At an unstable equilibrium, the negative curvature with respect to the primal variable encourages the algorithm to escape. Specifically, suppose the algorithm is initialized at an unstable equilibrium  $(X^{(0)}, Y^{(0)})$ , the descent direction for  $X^{(0)}$  is determined by the eigenvectors of

$$H_{X^{(0)}} = A + Y^{(0)}B$$

associated with the negative eigenvalues. After one iteration, we obtain  $(X^{(1)}, Y^{(1)})$ . The Hessian matrix becomes

$$H_{X^{(1)}} = A + Y^{(1)}B.$$

Since  $Y^{(1)} = X^{(0)\top} A^{(0)} X^{(0)}$  is a stochastic approximation, the random noise can make  $Y^{(1)}$  significantly different from  $Y^{(0)}$ . Thus, the eigenvectors of  $H_{X^{(1)}}$  associated with negative eigenvalues can be also very different from those of  $H_{X^{(0)}}$ . This phenomenon can seriously confuse the algorithm about the descent direction of the primal variable. We remark that such an issue does not appear under the commutative assumption. We suspect this is very likely an artifact of our proof technique, since our numerical experiments have provided some empirical evidences of the convergence of SGHA.

## 5 Discussion

Here we briefly discuss a few related works:

- [Li et al. \(2016b\)](#) propose a framework for characterizing the stationary points in the unconstrained nonconvex matrix factorization problem, while our GEV problem is constrained. Different from their analysis, we analyze the optimization landscape of the corresponding Lagrangian function. When characterize the stationary points, we need to take both primal and dual variables into consideration, which is technically more challenging.
- [Ge et al. \(2016a\)](#) consider the (off-line) generalized eigenvalue problem in a finite sum form. Unlike online setting, they access exact  $A$  and  $B$  in each iteration. Specifically, they need to access exact  $A$  and  $B$  to compute an approximate inverse of  $B$  to find the descent direction. Meanwhile, they also need a modified Gram Schmidt process, which also requires accessing exact  $B$ , to maintain the solution on the generalized Stiefel manifold (defined by  $X^\top B X = I_r$  via exact  $B$ , [Mishra and Sepulchre \(2016\)](#)). Our proposed stochastic search, however, is a full stochastic primal-dual algorithm, which neither require accessing exact  $A$  and  $B$ , nor enforcing the the primal variables to stay on the manifold.

### Acknowledgements

This research was partially supported by DARPA Young Faculty Award N66001-14-1-4047.



## References

- ALLEN-ZHU, Z. and LI, Y. (2016). Doubly accelerated methods for faster CCA and generalized eigendecomposition. *arXiv preprint arXiv:1607.06017* .
- ARORA, R., MARINOV, T. V., MIANJY, P. and SREBRO, N. (2017). Stochastic approximation for canonical correlation analysis. In *Advances in Neural Information Processing Systems*.
- BOJANAPALLI, S., NEYSHABUR, B. and SREBRO, N. (2016). Global optimality of local search for low rank matrix recovery. *arXiv preprint arXiv:1605.07221* .
- BOYD, S. and VANDENBERGHE, L. (2004). *Convex optimization*. Cambridge university press.
- CHEN, Y., LAN, G. and OUYANG, Y. (2014). Optimal primal-dual methods for a class of saddle point problems. *SIAM Journal on Optimization* **24** 1779–1814.
- CHEN, Z., YANG, F. L., LI, C. J. and ZHAO, T. (2017). Online multiview representation learning: Dropping convexity for better efficiency. *arXiv preprint arXiv:1702.08134* .
- COOK, R. D. and NI, L. (2005). Sufficient dimension reduction via inverse regression: A minimum discrepancy approach. *Journal of the American Statistical Association* **100** 410–428.
- DOOB, J. L. (1942). The brownian movement and stochastic equations. *Annals of Mathematics* 351–369.
- DUMMIT, D. S. and FOOTE, R. M. (2004). *Abstract algebra*, vol. 3. Wiley Hoboken.
- ETHIER, S. N. and KURTZ, T. G. (2009). *Markov processes: characterization and convergence*, vol. 282. John Wiley & Sons.
- FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2001). *The elements of statistical learning*, vol. 1. Springer series in statistics New York.
- GE, R., HUANG, F., JIN, C. and YUAN, Y. (2015). Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Conference on Learning Theory*.
- GE, R., JIN, C., NETRAPALLI, P., SIDFORD, A. ET AL. (2016a). Efficient algorithms for large-scale generalized eigenvector computation and canonical correlation analysis. In *International Conference on Machine Learning*.
- GE, R., LEE, J. D. and MA, T. (2016b). Matrix completion has no spurious local minimum. In *Advances in Neural Information Processing Systems*.
- GOODFELLOW, I., POUGET-ABADIE, J., MIRZA, M., XU, B., WARDE-FARLEY, D., OZAIR, S., COURVILLE, A. and BENGIO, Y. (2014). Generative adversarial nets. In *Advances in neural information processing systems*.
- GORRELL, G. (2006). Generalized hebbian algorithm for incremental singular value decomposition in natural language processing. In *EACL*, vol. 6. Citeseer.
- HAROLD, J., KUSHNER, G. and YIN, G. (1997). Stochastic approximation and recursive algorithm and applications. *Application of Mathematics* **35**.
- IOUDITSKI, A. and NESTEROV, Y. (2014). Primal-dual subgradient methods for minimizing uniformly convex functions. *arXiv preprint arXiv:1401.1792* .
- KUSHNER, H. and YIN, G. G. (2003). *Stochastic approximation and recursive algorithms and applications*, vol. 35. Springer Science & Business Media.
- LAN, G., LU, Z. and MONTEIRO, R. D. (2011). Primal-dual first-order methods with  $\{O\}(1/\epsilon)$  iteration-complexity for cone programming. *Mathematical Programming* **126** 1–29.
- LI, C. J., WANG, M., LIU, H. and ZHANG, T. (2016a). Near-optimal stochastic approximation for online principal component estimation. *arXiv preprint arXiv:1603.05305* .
- LI, X., WANG, Z., LU, J., ARORA, R., HAUPT, J., LIU, H. and ZHAO, T. (2016b). Symmetry, saddle points, and global geometry of nonconvex matrix factorization. *arXiv preprint arXiv:1612.09296* .
- LIU, T., CHEN, Z., ZHOU, E. and ZHAO, T. (2018). Toward deeper understanding of non-convex stochastic optimization with momentum using diffusion approximations. *arXiv preprint arXiv:1802.05155* .
- LUENBERGER, D. G., YE, Y. ET AL. (1984). *Linear and nonlinear programming*, vol. 2. Springer.
- LUO, Z.-Q. and TSENG, P. (1993). Error bounds and convergence analysis of feasible descent methods: a general approach. *Annals of Operations Research* **46** 157–178.
- MIKA, S., RATSCH, G., WESTON, J., SCHOLKOPF, B. and MULLERS, K.-R. (1999). Fisher discriminant analysis with kernels. In *Neural networks for signal processing IX, 1999. Proceedings of the 1999 IEEE signal processing society workshop*. Ieee.

- MISHRA, B. and SEPULCHRE, R. (2016). Riemannian preconditioning. *SIAM Journal on Optimization* **26** 635–660.
- POLYAK, B. T. (1963). Gradient methods for minimizing functionals. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki* **3** 643–653.
- SANGER, T. D. (1989). Optimal unsupervised learning in a single-layer linear feedforward neural network. *Neural networks* **2** 459–473.
- SHAPIRO, A., DENTCHEVA, D. and RUSZCZYŃSKI, A. (2009). *Lectures on stochastic programming: modeling and theory*. SIAM.
- SUN, J., QU, Q. and WRIGHT, J. (2016). A geometric analysis of phase retrieval. In *Information Theory (ISIT), 2016 IEEE International Symposium on*. IEEE.
- SUTTON, R. S., MCALLESTER, D. A., SINGH, S. P. and MANSOUR, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*.
- ZHU, Z., LI, Q., TANG, G. and WAKIN, M. B. (2017). The global optimization geometry of non-symmetric matrix factorization and sensing. *arXiv preprint arXiv:1703.01256* .