

Nonlinear Dependent Component Analysis: Identifiability and Algorithm

Qi Lyu

School of EECS, Oregon State University
Corvallis, OR 97331, USA
lyuqi@oregonstate.edu

Xiao Fu

School of EECS, Oregon State University
Corvallis, OR 97331, USA
xiao.fu@oregonstate.edu

Abstract—This work studies the model identification problem of a class of post-nonlinear mixture models in the presence of dependent latent components. Particularly, our interest lies in latent components that are nonnegative and sum-to-one. This problem is motivated by applications such as hyperspectral unmixing under nonlinear distortion effects. Many prior works tackled nonlinear mixture analysis using statistical independence among the latent components, which is not applicable in our case. A recent work by Yang *et al.* put forth a solution for this problem leveraging functional equations. However, the identifiability conditions derived there are somewhat restrictive. The associated implementation also has difficulties—the function approximator used in their work may not be able to represent general nonlinear distortions and the formulated constrained neural network optimization problem may be challenging to handle. In this work, we advance both the theoretical and practical aspects of the problem of interest. On the theory side, we offer a new identifiability condition that circumvents a series of stringent assumptions in Yang *et al.*'s work. On the algorithm side, we propose an easy-to-implement unconstrained neural network-based algorithm—without sacrificing function approximation capabilities. Numerical experiments are employed to support our design.

Index Terms—post-nonlinear mixture, dependent component analysis, identifiability, neural networks, nonnegative matrix factorization

I. INTRODUCTION

Latent component analysis has been an essential tool for a large variety of applications in signal processing (SP) and machine learning (ML). Many component analysis tools have been proposed, e.g., principal component analysis (PCA) [1], independent component analysis (ICA) [2], [3], nonnegative matrix factorization (NMF) [4], [5], dictionary learning/sparse coding [6], and tensor decomposition models [7], just to name a few.

One of the most important theoretical aspects pertaining to component analysis is *model identifiability*—since these tools are oftentimes associated with unsupervised learning and blind signal processing tasks, e.g., topic model learning [8], community detection [9], and blind source separation [2]. With model identifiability, the latent components of interest can be identified (often up to trivial ambiguities such as scaling and permutation) through learning the model parameters of the employed component analysis models from the observed data.

This work is supported in part by the National Science Foundation (NSF) under projects NSF ECCS 1608961, CNS 2003082 and ECCS 1808159.

Establishing model identifiability is a nontrivial task. In a nutshell, many component analysis models can be understood as matrix factorization models—which is in general non-unique, thereby lacking identifiability. A common practice to circumvent this issue is introducing structural information as constraints, e.g., statistical independence in ICA, nonnegativity in NMF, and sparsity in dictionary learning. The identifiability analyses for these classic component analysis models are elegant, and the model uniqueness results have improved performance of many core tasks in SP and ML; see [2], [4], [6], [7], [10].

On the other hand, most classic component analysis models can be understood as variants of matrix and tensor factorization models. These models essentially assume that all the data vectors are generated from a *linear* subspace (a Khatri-Rao subspace for tensors). This is often over-simplified for reality—since nonlinear distortions happen ubiquitously. Starting from the 1980s, efforts have been made towards incorporating nonlinear distortions into component analysis [11]. One notable line of work is the so-called *nonlinear independent component analysis* (nICA) [12]–[19]. The nICA framework considers *unknown* nonlinear distortions on top of the ICA model. One take-home point learned is that general nonlinear distortion is not identifiable under the framework of ICA [12]. To circumvent this, one may exploit certain structures of nonlinear distortions—e.g., under the so-called post-nonlinear mixture model [15]–[18] that is considered realistic for a number of sensing problems in radar, wireless communications and bioinformatics. One may also utilize more prior information from the data to remove nonlinear distortions, e.g., temporal correlations; see [13], [14], [19].

The model identifiability results under the nonlinear ICA frameworks are encouraging—showing that nonlinear distortions may be provably removable. However, assuming statistical independence among the latent components may be stringent. To relax this assumption, the recent work in [20] addresses the problem of *dependent component* identification under the post-nonlinear mixture model. To be specific, the authors of [20] considered a model where the latent components are nonnegative and sum-to-one—which is often considered in weighted mixture models such as soft clustering [21] and hyperspectral imaging [22]. Working from there, and combining insights from NMF identifiability, latent component

identifiability was shown.

However, some challenges remain. First, the work in [20] assumes that the model parameters are all nonnegative, which may restrict the applicability in some cases. Second, the nonlinear model identifiability hinges on a special assumption that the composition of the learned nonlinearity-compensating function and the nonlinear distortion is convex or concave—which is hard to verify or control, thereby being restrictive. Third, the work in [20] utilizes a neural network (NN) to approximate the “inverse” of the unknown nonlinear distortion, but the NN used there has positive weights for regularity purposes. This may cause performance losses—although NNs are universal function approximators, the function-approximation capacity of positive NNs is unknown. Optimization involving positivity-constrained NNs is also challenging.

In this work, we put forth a new solution for the nonlinear model identification problem in [20] and its extensions. Our contribution is twofold: First, we offer a new identifiability result that does not rely on the restrictive assumptions used in [20]. This may substantially enlarge the spectrum of applicable cases for this nonlinear component analysis model. Second, we propose a general-purpose neural network (other than a special positive neural network) based implementation for the formulated problem. This way, the risk of not being able to approximate certain nonlinear functions is removed. The associated optimization problem is also much easier to handle, leveraging existing optimizers for NNs, e.g., Adam [23]—which makes our implementation easily scalable. Numerical experiments are employed to showcase the performance of the proposed approach.

II. BACKGROUND

A. Linear and Nonlinear Independent Component Analysis

Many classic latent component analysis models start with the following linear mixture model (LMM):

$$\mathbf{x}_\ell = \mathbf{A}\mathbf{s}_\ell, \quad \ell = 1, 2, \dots \quad (1)$$

where $\mathbf{x}_\ell \in \mathbb{R}^M$ denotes the ℓ th observed data sample, $\mathbf{A} \in \mathbb{R}^{M \times K}$ the “mixing system” (or the basis of the subspace where \mathbf{x}_ℓ ’s reside), and $\mathbf{s}_\ell \in \mathbb{R}^K$ the vector that holds the K latent components. When $M \geq K$, the system is considered “over-determined” and can be reduced to an $M = K$ case via dimensionality reduction approaches, e.g., PCA [24]. In this work, we let $M = K$ for simplicity.

Many applications are concerned with identifying \mathbf{A} and/or $\{\mathbf{s}_\ell\}$ from $\{\mathbf{x}_\ell\}$. Model identifiability has been established under various conditions. For example, ICA assumes that $s_{k,\ell}$ ’s are statistically independent [25], and NMF assumes that \mathbf{A} and \mathbf{s}_ℓ are element-wise nonnegative [4], [26], [27]. Beyond the classic LMMs, nonlinear mixtures have also been considered, mostly under the umbrella of nonlinear ICA. The notable line of work in [12]–[14], [19] considers the model

$$\mathbf{x}_\ell = \mathbf{g}(\mathbf{s}_\ell), \quad \ell = 1, 2, \dots$$

where $\mathbf{g}(\cdot) : \mathbb{R}^K \rightarrow \mathbb{R}^M$ is an invertible continuous nonlinear distortion applied onto the latent components. This model is

in general not identifiable, even if one assumes that $s_{k,\ell}$ ’s are statistically independent [12]. A series of additional structural information on \mathbf{s}_ℓ (e.g., temporal correlations) has been exploited to establish identifiability. Another way for establishing model identifiability is to exploit structural information of the nonlinear distortions, other than that of the latent components. The post-nonlinear mixture (PNM) model is often considered [15]–[18], [20], where we have

$$\mathbf{x}_\ell = \mathbf{g}(\mathbf{As}_\ell), \quad \ell = 1, 2, \dots \quad (2)$$

in which $\mathbf{g}(\cdot) = [g_1(\cdot), \dots, g_M(\cdot)]^\top$ and $g_m(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is a scalar-to-scalar nonlinear continuous invertible function. The PNM model has a variety of applications in sensing-related problems, e.g., radar, wireless communications, and biomedical sensing; see discussions in [17]. Under the PNM model, the identifiability of \mathbf{s}_ℓ has also been established—using the statistical independence of the latent components [17], [18].

B. Nonlinear Dependent Component Analysis: Prior Art

Very recently, Yang *et al.* [20] considered a problem under the PNM model. Instead of having $s_{k,\ell}$ ’s to be statistically independent, the model assumption is that

$$\mathbf{s}_\ell \in \Delta, \quad \Delta = \{\mathbf{s} \in \mathbb{R}^K \mid \mathbf{1}^\top \mathbf{s} = 1, \mathbf{s} \geq \mathbf{0}\}. \quad (3)$$

Under this model, the observations (i.e., \mathbf{x}_ℓ ’s) are generated as weighted combinations (or, more precisely, convex combinations) of $\mathbf{a}_1, \dots, \mathbf{a}_K$ and then distorted by $g_1(\cdot), \dots, g_M(\cdot)$. Note that the convex combination model is a particularly important one that finds applications in topic modeling [8], soft clustering [21], and hyperspectral unmixing [22]. The \mathbf{g} -part is also critical in modeling nonlinear effects happening in practice, e.g., those hyperspectral imaging [28]. Note that since $\mathbf{1}^\top \mathbf{s}_\ell = 1$, the latent components are *dependent*.

The work in [20] utilizes the sum-to-one structure to construct a functional equation, and shows that under some conditions a carefully constructed model identification criterion can “remove” $\mathbf{g}(\cdot)$ through learning a nonlinear function $\mathbf{f}(\cdot)$. Then, the problem for identifying \mathbf{s}_ℓ becomes a classic NMF problem. There are a number of caveats. First, the assumptions for \mathbf{g} -removal might be restrictive. The assumptions include that \mathbf{A} being nonnegative and incoherent, and that the elements of $\mathbf{f} \circ \mathbf{g}$ are convex or concave functions. The last condition is particularly hard to enforce since one has no control for it—or even a way for checking it. Second, when implementing the learning criterion, the authors in [20] use a neural network to represent \mathbf{f} . However, the NN there is with positive network weights for enforcing function invertibility. This construction may have hindered the function-approximation capability of NNs. The associated constrained optimization problem is also challenging to handle. The work in [20] uses a trust-region based nonconvex quadratic programming method, which may not be scalable.

III. PROPOSED APPROACH

In this work, we offer a new solution under the PNM model and (3) that effectively circumvent the challenges in [20].

A. A Functional Equation-based Formulation

Ideally, we expect to learn element-wise invertible nonlinear function \mathbf{f} such that the following holds:

$$\mathbf{1}^\top \mathbf{f}(\mathbf{x}_\ell) = \mathbf{1}^\top \mathbf{f} \circ \mathbf{g}(\mathbf{A}\mathbf{s}_\ell) = \mathbf{1}^\top \mathbf{h}(\mathbf{A}\mathbf{s}_\ell) = 1, \quad (4)$$

for all ℓ . Here \mathbf{h} is also an element-wise function with $h_i = f_i \circ g_i$. In other words, our learning objective is to find an invertible function to reverse the distortions introduced by \mathbf{g} so that the sum-to-one condition can be satisfied.

Formally, we wish to have the following criterion satisfied in terms of f -searching:

$$\text{find } \mathbf{f} \quad (5a)$$

$$\text{s.t. } \mathbf{1}^\top \mathbf{f}(\mathbf{x}) = 1, \forall \mathbf{x} \in \mathcal{X} \quad (5b)$$

$$f_i \text{ is invertible over } \mathcal{X}, \forall i, \quad (5c)$$

where $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^M | \mathbf{x} = \mathbf{g}(\mathbf{As}), \forall \mathbf{s} \in \text{int}\Delta, \mathbf{s} \in \mathbb{R}^K\}$, in which $\text{int}\Delta$ means the interior of Δ . Note that the criterion is identical to what was proposed in [20], where the functional equations in (5b) play key role in removing the nonlinear distortions. The difference lies in model assumptions. In particular, Yang *et al.* assumed that \mathbf{A} is generic, nonnegative and incoherent in [20]. In this work, we only require that \mathbf{A} is generic (i.e., the entries are drawn from any jointly continuous distribution). Yang *et al.* showed the following theorem:

Theorem 1 [20] Consider the post-nonlinear mixture model $\mathbf{x}_\ell = \mathbf{g}(\mathbf{As}_\ell)$ where $\mathbf{A} \in \mathbb{R}^{K \times K}$ with the constraint $\mathbf{s}_\ell \in \text{int}\Delta$, where $g_i(\cdot)$ for all K are continuous and invertible. Assume that $N \rightarrow \infty$ and all points in \mathcal{X} are available. Also assume that \mathbf{A} is drawn from any joint continuous distribution.

In addition, assume that

- 1) \mathbf{A} is nonnegative and is incoherent (see the definition in [20]); and that
- 2) by solving problem (5), the resulting $h_i = f_i \circ g_i$'s are all convex or concave for $i = 1, \dots, K$.

Then, h_i has to be an affine functions almost surely; i.e., any f_i that is a solution of (5) satisfies

$$h_i(x) = f_i \circ g_i(x) = c_i x + d_i, \quad i = 1, \dots, K,$$

where c_i, d_i are constants. In addition, if $\sum_{i=1}^K d_i \neq 0$, we have $h_i(x) = f_i \circ g_i(x) = \alpha_i x, \quad i = 1, \dots, K, \quad \alpha_i \neq 0, \forall i$.

The theorem showed that the PNM model is identifiable even under dependent latent components. The challenge is that both conditions 1) and 2) may be restrictive—especially condition 2). In this work, we show that these conditions are in fact not needed. To proceed, we show the following:

Lemma 1 Consider $\mathbf{s} = [s_1, \dots, s_K]^\top \in \text{int}\Delta$. Then, $\frac{\partial s_i}{\partial s_j} = 0$ for $i \neq j$ where $i, j = 1, \dots, K - 1$.

Proof: Lemma 1 can be shown as follows. First, for $\mathbf{s}_\ell \in \text{int } \Delta_K$, we only have $K - 1$ free variables, i.e., without loss of generality, s_i for $i = 1, \dots, K - 1$. For any fixed \bar{s}_i , s_j can be any possible values in a nonempty continuous domain

(e.g., if $s_i = 0.5$ then the domain of s_j is $(0, 0.5)$ regardless of other components). Hence, if one treats s_i as a function of s_j , then the sensitivity of s_i w.r.t. s_j is defined as

$$\frac{\partial s_i}{\partial s_j} = \lim_{\Delta s_j \rightarrow 0} \frac{s_i(s_j + \Delta s_j) - s_i(s_j)}{\Delta s_j} = \lim_{\Delta s_j \rightarrow 0} \frac{\bar{s}_i - \bar{s}_i}{\Delta s_j} = 0.$$

This completes the proof. \blacksquare

This lemma is important for deriving our main theorem:

Theorem 2 (Nonlinearity Removal) Consider the post-nonlinear mixture model $\mathbf{x}_\ell = \mathbf{g}(\mathbf{As}_\ell)$ with the constraint $\mathbf{s}_\ell \in \text{int}\Delta$, where $g_i(\cdot)$ for all i are continuous and invertible. Assume that $N \rightarrow \infty$ and all points in \mathcal{X} are available. In addition, assume that $K \geq 3$, that $\mathbf{A} \in \mathbb{R}^{K \times K}$ is drawn from any joint continuous distribution, and that the learned $h_i = f_i \circ g_i$ is twice differentiable for all $i = 1, \dots, K$. Then, by solving problem (5), the resulting h_i 's are affine functions almost surely; i.e., any f_i that is a solution of (5) satisfies $h_i(x) = f_i \circ g_i(x) = c_i x + d_i, \quad i = 1, \dots, K$, where c_i, d_i are constants. In addition, if $\sum_{i=1}^K d_i \neq 0$, we have

$$h_i(x) = f_i \circ g_i(x) = \alpha_i x, \quad i = 1, \dots, K, \quad \alpha_i \neq 0, \forall i. \quad (6)$$

The proof sketch is as follows. According to Lemma 1, we have $\frac{\partial s_i}{\partial s_j} = 0$ for $\mathbf{s} \in \text{int } \Delta$. Then, by taking second order derivatives of the equality constraint in (4) w.r.t. s_i and s_j , it ends up with a system of linear equations that involves the vector $\mathbf{h}'' = [h_1'', \dots, h_M'']^\top$, i.e., $\mathbf{H}\mathbf{h}'' = \mathbf{0}$. By utilizing the assumptions, one can show that \mathbf{H} has full column rank—which immediately implies $\mathbf{h}'' = \mathbf{0}$. This further leads to that all h_i 's are affine. We defer the detailed proof to a pertinent journal version.

We would like to remark that Theorem 2 offers a set of conditions that are substantially more relaxed relative to Theorem 1. Notably, the conditions in 1) and 2) of Theorem 1 are not used in our theorem. Relaxing the nonnegativity of \mathbf{A} makes the method applicable to a lot more problems where the mixing system can have negative entries (e.g., speech separation). Removing 2) is also quite desirable, since this condition cannot be guaranteed or checked when implementing the criterion.

B. Latent Component Identification

Note that under (6), the following holds:

$$\mathbf{f}(\mathbf{x}_\ell) = \mathbf{C}\mathbf{A}\mathbf{s}_\ell = \mathbf{B}\mathbf{s}_\ell, \quad \ell = 1, 2, \dots, N$$

where $\mathbf{C} = \text{Diag}(\alpha_1, \dots, \alpha_M)$. This model is identical to the structural matrix factorization model in [24], [26], [29], which is identifiable if $\{\mathbf{s}_\ell\}$ satisfies certain conditions, e.g., the separability condition or the sufficiently scattered condition; see details in [4], [26], [29], [30]. Hence, a simple strategy is to first implement the criterion in (4) for nonlinearity removal. Then, any structural matrix factorization algorithm proposed in the literature, e.g., those in [29], [30], can be employed for identifying \mathbf{s}_ℓ from $\mathbf{f}(\mathbf{x}_\ell)$. In this work, we utilize the minimum-volume enclosing simplex (MVES) algorithm from [30] for \mathbf{s}_ℓ -identification after nonlinearity removal.

C. Neural Network-based Implementation

We have shown that solving Problem (5) removes the nonlinear distortions. However, Problem (5) is not really “workable” since it involves continuous functional searching. To approach this formulation, we parameterize the function f with neural networks due to their universal approximation ability. Each f_i is approximated by an individual neural network. Hence, the practical formulation is as follows:

$$\begin{aligned} \min_{\theta_f, \theta_g} & \sum_{\ell=1}^N (1 - \mathbf{1}^\top f_{\text{NN}}(\mathbf{x}_\ell))^2 \\ & + \lambda \sum_{\ell=1}^N \|\mathbf{x}_\ell - g_{\text{NN}}(f_{\text{NN}}(\mathbf{x}_\ell))\|_2^2 \end{aligned} \quad (7)$$

where two neural networks $f_{\text{NN}} = [f_{\text{NN}}^{(1)}, \dots, f_{\text{NN}}^{(M)}]^\top$ and $g_{\text{NN}} = [g_{\text{NN}}^{(1)}, \dots, g_{\text{NN}}^{(M)}]^\top$ are parameterized by θ_f and θ_g , respectively. Note that g_{NN} can be regarded as estimation for the ground-truth g .

To explain the above formulation, note that the first fitting term is for approximating the equality constraint in (5b). The second term articulates the difference between our implementation and that in [20]. The latter does not have the second term in (7). Instead, a constraint $\theta_f > \mathbf{0}$ is employed. The reason is that under this positivity constraint, the function f_{NN} is always invertible, which satisfies the problem specification in (5c). However, this may be problematic since positive NNs may not retain the universal approximation property for nonlinear functions—while the universal function approximation ability is the reason why one uses NNs in the first place.

In our implementation, we use the second term in (7) to promote invertibility of the learned f_{NN} . It is not hard to see the following:

Lemma 2 Assume that there exists a function g_{NN} such that for all $\mathbf{x} \in \mathcal{X}$, and that $\mathbf{x} = g_{\text{NN}}(f_{\text{NN}}(\mathbf{x}))$ holds. Then, f_{NN} is invertible over \mathcal{X} .

Hence, when N is large, the regularization approximately enforces invertibility over \mathcal{X} .

Another benefit of employing our formulation other than the positivity constraint formulation as in [20] is that unconstrained optimization for NNs is much easier. A number of off-the-shelf optimizers developed for large-scale NN-related optimization, e.g., Adam [23] based stochastic gradient, can be utilized to handle the formulated problem.

IV. NUMERICAL RESULTS

We use two baselines in our simulations, i.e., nonlinear matrix factor recovery (NMFR) [20] that was developed under the same model and the linear model MVES [30].

Our formulation is tackled by PyTorch-based Adam algorithm [23] with the initial step size being 10^{-3} . Adam is a stochastic gradient algorithm that works under mini-batch settings. The batch size is 5,000 in our simulations. The algorithm stops after running 5,000 epochs. The parameter λ

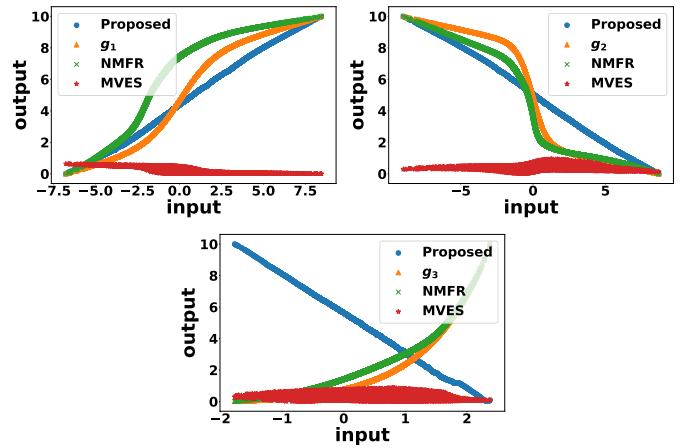


Fig. 1. Nonlinearity removal effects compared with baselines.

is set to be 10^{-5} . For f_{NN} and g_{NN} , each channel is modeled with a single hidden layer neural network with 256 neurons.

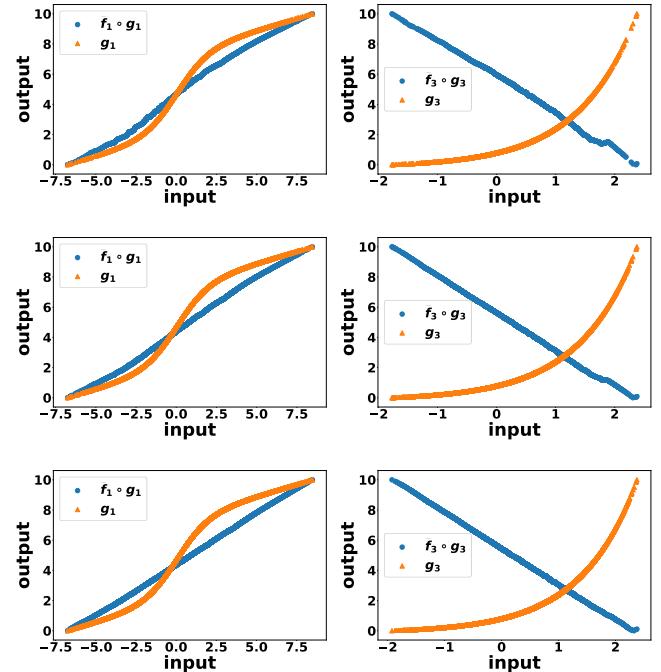


Fig. 2. Impact of N . From top to bottom, N is 5000, 10000 and 20000, respectively.

For the first simulation, we set $K = 3$ and the three nonlinear functions are: $g_1(x) = 5\text{sigmoid}(x) + 0.3x$, $g_2(x) = -3\tanh(x) - 0.2x$ and $g_3 = 0.4 \exp(x)$. The number of samples is 10,000. The matrix A is drawn from standard Gaussian distribution. The learned $f_{\text{NN}}^{(i)} \circ g_i$'s are shown in Fig. 1. One can see that the proposed method works remarkably better than the NMFR from [20]. Although the two methods start with the same conceptual formulation in (5), the performance difference may come from implementation strategies: since Yang *et al.*'s implementation uses positive NNs, it may not be able to approximate the true solution f ;

TABLE I
MSE BETWEEN \mathbf{S} AND ESTIMATED $\hat{\mathbf{S}}$

N	Proposed	NMFR	MVES
5000	$3.67e^{-3} \pm 1.53e^{-3}$	$1.30e^{-1} \pm 4.48e^{-2}$	$5.49e^{-2} \pm 4.50e^{-3}$
10000	$1.11e^{-3} \pm 5.54e^{-4}$	$9.68e^{-2} \pm 2.84e^{-2}$	$5.84e^{-2} \pm 8.60e^{-3}$
20000	$1.58e^{-4} \pm 1.65e^{-4}$	$8.01e^{-2} \pm 1.78e^{-2}$	$4.96e^{-2} \pm 3.11e^{-3}$

in addition, constrained optimization may be much harder than dealing with our unconstrained formulation.

For the next simulation, we qualitatively show the influence of the sample size N . With the same setting as in Fig. 1, we randomly select two channels of the observations and show the learned composite functions in Fig. 2. The figure clearly illustrates that as more samples are available, the nonlinearity removal performance improves substantially. This also echoes our main theorem—which was developed under $N = \infty$ (more precisely, uncontably infinite “samples”).

In the last simulation, we combine nonlinearity removal and s_ℓ -identification, where the second phase is conducted by applying MVES to $f_{\text{NN}}(\mathbf{x}_\ell)$ for $\ell = 1, \dots, N$. The performance measure here is the mean squared error (MSE) of the estimated $\mathbf{S} = [s_1, \dots, s_N]$ [29]. The results are shown in Table. I, which are averaged over 10 random trials. It can be seen that the performance of the proposed approach exhibits the best performance—and shows a notable margin over the baselines. In particular, the MSE performance is one or two orders of magnitude lower compared to NMFR and MVES.

V. CONCLUSION

In this work, we address the nonlinearity removal and latent component identification problems under the post-nonlinear model in the presence of nonnegative and sum-to-one dependent components. Our contribution is two fold: first, we have tightened the sufficient conditions under which the nonlinearity is removable—which offers substantially more relaxed conditions relative to a recently derived result; second, we offer a new NN-based formulation that has better function approximation ability and is easier to optimize. As a result, the numerical performance is largely improved compared to prior work.

REFERENCES

- [1] S. Wold, K. Esbensen, and P. Geladi, “Principal component analysis,” *Chemometrics and intelligent laboratory systems*, vol. 2, no. 1-3, pp. 37–52, 1987.
- [2] P. Comon and C. Jutten, *Handbook of Blind Source Separation*. Elsevier, 2010.
- [3] A. Hyvärinen and E. Oja, “Independent component analysis: algorithms and applications,” *Neural networks*, vol. 13, no. 4-5, pp. 411–430, 2000.
- [4] X. Fu, K. Huang, N. D. Sidiropoulos, and W.-K. Ma, “Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications,” *IEEE Signal Process. Mag.*, vol. 36, no. 2, pp. 59–80, March 2019.
- [5] D. Lee and H. Seung, “Learning the parts of objects by non-negative matrix factorization,” *Nature*, vol. 401, no. 6755, pp. 788–791, 1999.
- [6] I. Tasic and P. Frossard, “Dictionary learning,” *IEEE Signal Process. Mag.*, vol. 28, no. 2, pp. 27–38, 2011.
- [7] N. D. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, and C. Faloutsos, “Tensor decomposition for signal processing and machine learning,” *IEEE Trans. Signal Process.*, vol. 65, no. 13, pp. 3551–3582, 2017.
- [8] K. Huang, X. Fu, and N. D. Sidiropoulos, “Anchor-free correlated topic modeling: Identifiability and algorithm,” in *Proceedings of NIPS 2016*, 2016, pp. 1786–1794.
- [9] K. Huang and X. Fu, “Detecting overlapping and correlated communities without pure nodes: Identifiability and algorithm,” in *Proceedings of ICML 2019*, vol. 97, 09–15 Jun 2019, pp. 2859–2868.
- [10] N. Gillis, “The why and how of nonnegative matrix factorization,” *Regularization, Optimization, Kernels, and Support Vector Machines*, vol. 12, p. 257, 2014.
- [11] M. A. Kramer, “Nonlinear principal component analysis using autoassociative neural networks,” *AIChE journal*, vol. 37, no. 2, pp. 233–243, 1991.
- [12] A. Hyvärinen and P. Pajunen, “Nonlinear independent component analysis: Existence and uniqueness results,” *Neural Networks*, vol. 12, no. 3, pp. 429–439, 1999.
- [13] A. Hyvärinen and H. Morioka, “Unsupervised feature extraction by time-contrastive learning and nonlinear ICA,” in *Proceedings of NIPS 2016*, 2016, pp. 3765–3773.
- [14] ———, “Nonlinear ICA of temporally dependent stationary sources,” in *Proceedings of AISTATS 2017*, vol. 54, 20–22 Apr 2017, pp. 460–469.
- [15] E. Oja, “The nonlinear PCA learning rule in independent component analysis,” *Neurocomputing*, vol. 17, no. 1, pp. 25–45, 1997.
- [16] A. Ziehe, M. Kawanabe, S. Harmeling, and K.-R. Müller, “Blind separation of post-nonlinear mixtures using linearizing transformations and temporal decorrelation,” *Journal of Machine Learning Research*, vol. 4, no. Dec, pp. 1319–1338, 2003.
- [17] A. Taleb and C. Jutten, “Source separation in post-nonlinear mixtures,” *IEEE Trans. Signal Process.*, vol. 47, no. 10, pp. 2807–2820, 1999.
- [18] S. Achard and C. Jutten, “Identifiability of post-nonlinear mixtures,” *IEEE Signal Process. Lett.*, vol. 12, no. 5, pp. 423–426, 2005.
- [19] A. Hyvärinen, H. Sasaki, and R. E. Turner, “Nonlinear ica using auxiliary variables and generalized contrastive learning,” in *The 22nd International Conference on Artificial Intelligence and Statistics. Journal of Machine Learning Research*, 2019, pp. 859–868.
- [20] B. Yang, X. Fu, N. D. Sidiropoulos, and K. Huang, “Learning nonlinear mixtures: Identifiability and algorithm,” *IEEE Transactions on Signal Processing*, 2020.
- [21] G. Hamerly and C. Elkan, “Alternatives to the k-means algorithm that find better clusterings,” in *Proceedings of the eleventh international conference on Information and knowledge management*, 2002, pp. 600–607.
- [22] W.-K. Ma, J. Bioucas-Dias, T.-H. Chan, N. Gillis, P. Gader, A. Plaza, A. Ambikapathi, and C.-Y. Chi, “A signal processing perspective on hyperspectral unmixing,” *IEEE Signal Process. Mag.*, vol. 31, no. 1, pp. 67–81, Jan 2014.
- [23] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [24] X. Fu, W.-K. Ma, K. Huang, and N. D. Sidiropoulos, “Blind separation of quasi-stationary sources: Exploiting convex geometry in covariance domain,” *IEEE Trans. Signal Process.*, vol. 63, no. 9, pp. 2306–2320, May 2015.
- [25] P. Comon, “Independent component analysis, a new concept?” *Signal Processing*, vol. 36, no. 3, pp. 287 – 314, 1994.
- [26] X. Fu, K. Huang, and N. D. Sidiropoulos, “On identifiability of nonnegative matrix factorization,” *IEEE Signal Process. Lett.*, vol. 25, no. 3, pp. 328–332, 2018.
- [27] K. Huang, N. Sidiropoulos, and A. Swami, “Non-negative matrix factorization revisited: Uniqueness and algorithm for symmetric decomposition,” *IEEE Trans. Signal Process.*, vol. 62, no. 1, pp. 211–224, 2014.
- [28] N. Dobigeon, J.-Y. Tourneret, C. Richard, J. C. M. Bermudez, S. McLaughlin, and A. O. Hero, “Nonlinear unmixing of hyperspectral images: Models and algorithms,” *IEEE Signal Process. Mag.*, vol. 31, no. 1, pp. 82–94, 2014.
- [29] X. Fu, K. Huang, B. Yang, W. Ma, and N. D. Sidiropoulos, “Robust volume minimization-based matrix factorization for remote sensing and document clustering,” *IEEE Trans. Signal Process.*, vol. 64, no. 23, pp. 6254–6268, Dec 2016.
- [30] T.-H. Chan, C.-Y. Chi, Y.-M. Huang, and W.-K. Ma, “A convex analysis-based minimum-volume enclosing simplex algorithm for hyperspectral unmixing,” *IEEE Trans. Signal Process.*, vol. 57, no. 11, pp. 4418–4432, Nov. 2009.