



HAL
open science

Multilinear compressive sensing and an application to convolutional linear networks

François Malgouyres, Joseph Landsberg

► **To cite this version:**

François Malgouyres, Joseph Landsberg. Multilinear compressive sensing and an application to convolutional linear networks. *SIAM Journal on Mathematics of Data Science*, 2019, 1 (3), pp.446-475. hal-01494267v4

HAL Id: hal-01494267

<https://hal.science/hal-01494267v4>

Submitted on 1 Feb 2023

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

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

MULTILINEAR COMPRESSIVE SENSING AND AN APPLICATION TO CONVOLUTIONAL LINEAR NETWORKS*

FRANÇOIS MALGOUYRES [†] AND JOSEPH LANDSBERG [‡]

Abstract. We study a deep linear network endowed with the following structure: a matrix X is obtained by multiplying K matrices (called factors and corresponding to the action of the layers). The action of each layer (i.e. factor) is obtained by applying a fixed linear operator to a vector of parameters satisfying a constraint. The number of layers is not limited. Assuming that X is given and factors have been estimated, the error between the product of the estimated factors and X (i.e. the reconstruction error) is either the statistical or the empirical risk.

We provide necessary and sufficient conditions on the network topology under which a stability property holds. The stability property requires that the error on the parameters defining the near-optimal factors scales linearly with the reconstruction error (i.e. the risk). Therefore, under these conditions on the network topology, any successful learning task leads to stably defined features that can be interpreted.

In order to do so, we first evaluate how the Segre embedding and its inverse distort distances. Then, we show that any deep structured linear network can be cast as a generic multilinear problem that uses the Segre embedding. This is the *tensorial lifting*. Using the tensorial lifting, we provide a necessary and sufficient conditions for the identifiability of the factors up to a scale rearrangement. We finally provide necessary and sufficient condition called deep-Null Space Property (because of the analogy with the usual Null Space Property in the compressed sensing framework) which guarantees that the stability property holds.

We illustrate the theory with a practical example where the deep structured linear network is a convolutional linear network. We obtain a condition on the scattering of the supports which is strong but not empty. A simple test on the network topology can be implemented to test if the condition holds.

Key words. Interpretable learning, stable recovery, matrix factorization, deep linear networks, convolutional networks.

AMS subject classifications. 68T05; 90C99; 15-02

1. Introduction.

1.1. The aim of the paper. Deep learning has led to many practical breakthroughs and has led to significant improvements and state of the art performances in many fields such as computer vision, natural language processing, signal processing, robotics *etc.* The range of applications grows at a strong pace. Despite these empirical successes, the theory supporting deep learning is still far from satisfactory. For instance, sharp and accurate answers to the most natural questions on: the efficiency of optimization algorithms when applied to the objective function minimized in deep learning ([55, 37, 22, 23, 68, 10, 11, 44, 79]); and the expressiveness of the networks ([12, 3, 31, 26, 45, 27, 63, 76]); and guarantees on the statistical risk ([42, 34, 80, 69]) for learned neural networks are still missing. This makes it difficult to optimize and configure neural networks. Moreover, the absence of answers to these questions prevents certification that systems built with deep learning algorithms are robust.

*Submitted to the editors 7/3/2018.

[†]Institut de Mathématiques de Toulouse ; UMR5219 Université de Toulouse ; CNRS UPS IMT; F-31062 Toulouse Cedex 9, France; (Francois.Malgouyres@math.univ-toulouse.fr) — and — Institut de Recherche Technologique Saint Exupéry

Funding: This work has been supported by the DEEL program on Dependable and Explainable Learning (www.deel.ai)

[‡]Department of Mathematics; Mailstop 3368; Texas A& M University; College Station, TX 77843-3368; (jml@math.tamu.edu)

Funding: Joseph Landsberg is supported by NSF DMS-1405348 and AF1814254

The reasons explaining the outcome of a neural network are often difficult to highlight ([8, 62, 41]). Even worse, despite the settings described in [6, 4, 14, 53, 71, 83], the instability of the parameters optimizing the deep learning objective does not allow the interpretation of the features defined by these parameters. This last problem is the one we investigate in this work.

Our goal in this paper is to evaluate how far the architectures used in applications are from architectures for which we can guarantee that the parameters returned by the algorithm, and therefore the features defined using these parameters, are stably defined. To do so, we consider two families of networks and establish necessary and sufficient conditions on their topology guaranteeing that the features learned by the algorithm are stably defined.

More precisely, we establish statements of the following form for two families of deep networks. Below, the action of the network parameterized by \mathbf{h} is denoted $f_{\mathbf{h}}$.

INFORMAL THEOREM 1.1. Stability guarantee

We assume a known parameterized family of functions $f_{\mathbf{h}}$ and a metric¹ d between parameter pairs. We establish a necessary and sufficient condition on the family $f_{\mathbf{h}}$ guaranteeing that:

There exists a constant $C > 0$ such that for any input/output pairs I, X and any pair of parameters $\mathbf{h}^, \bar{\mathbf{h}}$ for which*

$$\delta = \|X - f_{\mathbf{h}^*}(I)\|,$$

and

$$\eta = \|X - f_{\bar{\mathbf{h}}}(I)\|,$$

are sufficiently small, we have

$$(1.1) \quad d(\bar{\mathbf{h}}, \mathbf{h}^*) \leq C(\delta + \eta).$$

Considering a regression problem, the values δ and η can be interpreted as the statistical or the empirical risk for the parameters $\bar{\mathbf{h}}$ and \mathbf{h}^* . The inequality (1.1) therefore guarantees that the set made of the parameters leading to a small risk has a small diameter. The features defined using such parameters are therefore stably defined. This seems to be the minimal condition allowing the interpretation of the features. The condition on the family of functions $f_{\mathbf{h}}$ is typically a condition on the topology of the network.

In [Informal theorem 1.1](#), $\bar{\mathbf{h}}$ and \mathbf{h}^* might have different roles. For instance, if we know that the input/output pairs have been generated using a particular $\bar{\mathbf{h}}$, possibly up to some error as modeled by δ , then (1.1) guarantees that \mathbf{h}^* is close to $\bar{\mathbf{h}}$ and provides a way to control the statistical risk.

The existing stability guarantees [6, 4, 14, 53, 71, 83, 60] consider this setting and describe both a network topology and an algorithm whose output \mathbf{h}^* is guaranteed to be close to $\bar{\mathbf{h}}$. In this study, we do not make any assumption on the construction of $\bar{\mathbf{h}}$ and \mathbf{h}^* and our objective is more modest. With regard to their objective, giving a necessary and sufficient condition of stability plays the same role as a complexity theory statement saying that a particular configuration is NP-hard. It rules out some network topologies.

Notice that, when I and X are such that it is possible to have $\delta = \eta = 0$, the above stability guarantee implies that the minimizer of the network objective function is unique. [Theorem 6.5](#) and [Theorem 6.4](#) will show that this uniqueness

¹The metric takes into account inter-layer rescaling.

condition is strongly related to the level of over-specification of the network. The simplified and intuitive statement is that optimal solutions of overspecified networks are not unique and are unstable. This explains the instability observed in applications. The theorems analyze this property in detail. This might be viewed as a negative result since overspecification is currently the main hypothesis of statements guaranteeing the success of the neural network optimization.

1.2. The considered deep networks.

1.2.1. Overview. We consider two kinds of deep networks: A general family of deep structured linear networks² in Sections [section 6](#) and [section 7](#); and a family of convolutional linear networks in [section 8](#). The formal statements for the deep structured linear networks are in Theorems [Theorem 7.2](#) and [Theorem 7.3](#). The statements for convolutional linear networks are in [Theorem 8.4](#). Below, we describe the deep structured linear networks.

1.2.2. Deep structured linear networks. The term deep linear network usually corresponds to fully-connected feed-forward networks, without bias, in which the activation function is the identity. In the general results described in this paper we consider deep linear networks and provide two means to enforce some structure to the network. As we describe below, the structures can be used to include: feed-forward linear networks; convolutional linear networks (as is done in [section 8](#)); the action of a ReLU activation function; sparse networks; non-negative networks; and combinations of the above. The family also includes most matrix factorization problems.

We model a deep structured linear network as a product of matrices called factors. The factors depend linearly on parameters in \mathbb{R}^S , for $S \in \mathbb{N}$.

More precisely, consider an arbitrary depth parameter $K \geq 1$. The number of layers is $K + 1$ and the layers are enumerated in such a way that the layer receiving the input is $K + 1$ and the layer returning the output is 1. We consider sizes $m_1 \dots m_{K+1} \in \mathbb{N}$, write $m_1 = m$, $m_{K+1} = n$. We consider, for $k = 1 \dots K$, the linear map

$$(1.2) \quad \begin{aligned} M_k : \mathbb{R}^S &\longrightarrow \mathbb{R}^{m_k \times m_{k+1}} \\ h &\longmapsto M_k(h) \end{aligned}$$

Given some parameters, $h_1, \dots, h_K \in \mathbb{R}^S$, the action of the deep structured linear network is the product

$$M_1(h_1) \cdots M_K(h_K)$$

The factor $M_K(h_K)$ might involve the inputs of the samples by considering $M_K(h_K) = M'_K(h_K)I$ for: a linear map M'_K ; and for a matrix I whose columns contain the inputs. Given outputs $X \in \mathbb{R}^{m \times n}$, the optimization of the parameters h_1, \dots, h_K defining the network aims at getting

$$M_1(h_1) \cdots M_K(h_K) \simeq X.$$

To model feed-forward linear networks, the mappings M_k , $k = 1 \dots K - 1$ (and M'_K) construct the matrix by placing the entry of h_k corresponding to an edge in the network in the corresponding entry in $M_k(h_k)$.

²We call this family *deep structured linear networks* because the family is endowed with tools to impose structures. We analyze the impact of the structure on the stability property. However, these tools might be used to define the usual deep linear network.

For convolutional layers, M_k and M'_K concatenate convolution matrices³ defined by a portion of the entries in \mathbf{h}_k . Each convolution matrix is at the location corresponding to a prescribed edge.

The main argument for studying deep structured linear networks is due to their strong connection to non-linear networks that uses the rectified linear unit (ReLU)⁴ activation function. We explain it in detail. The action of the ReLU activation function at the layer k treats every entry independently of the other entries and multiplies it by either 1 (the entry is kept) or 0 (the entry is canceled). More precisely, denoting $\mathbf{h} = (h_k)_{k=1..K}$, the action of the ReLU activation function on the layer k is to apply the map $A_k : \mathbb{R}^{m_k \times n} \mapsto \mathbb{R}^{m_k \times n}$ (where $m_k \times n$ is the size data in the layer k) such that:

$$(A_k M)_{i,j} = a_k(\mathbf{h})_{i,j} M_{i,j}, \quad \text{for } (i, j) \in \{1, \dots, m_k\} \times \{1, \dots, n\}$$

where $a_k(\mathbf{h}) \in \{0, 1\}^{m_k \times n}$ is defined by

$$a_k(\mathbf{h})_{i,j} = \begin{cases} 1 & \text{if } \left(M_{k+1}(h_{k+1}) A_{k+1} M_{k+2}(h_{k+2}) \cdots A_{K-1} M_K(h_K) \right)_{i,j} \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

The function

$$\begin{aligned} a_k : \mathbb{R}^{S \times K} &\longrightarrow \{0, 1\}^{m_k \times n} \\ \mathbf{h} &\longmapsto a_k(\mathbf{h}) \end{aligned}$$

is piecewise constant because $\{0, 1\}^{m_k \times n}$ is finite. (This has already been used in [68].) As a consequence, the parameter space $\mathbb{R}^{S \times K}$ is partitioned into subsets such that on every subset a_k is constant, for all $k = 1..K$. Therefore, on every subset the action of the non-linear network coincides with the action of a deep structured linear network that groups at every layer A_k and M_{k+1} . Further, the landscape of the objective function of the non-linear neural network that uses ReLU coincides, on every part of the partition, with the landscape of a deep structured linear network. This is a strong argument in favor of the study of deep structured linear networks.

Notice that deep structured linear networks have also been obtained in [22, 23, 44] by modeling the action of the activation function as random, independent of the input and when considering the expectation of the network action. However, these assumptions are not satisfied by the deep networks used in applications (see [23]) and it is not clear that this link can be exploited to obtain theoretical guarantees for realistic deep networks.

In addition to the structure induced by the operators M_k , we also consider structure imposed on the vectors h . We assume that we know a collection of models $\mathcal{M} = (\mathcal{M}^L)_{L \in \mathbb{N}}$ with the property that for every L , $\mathcal{M}^L \subset \mathbb{R}^{S \times K}$ is a given subset. We will assume that the parameters $\mathbf{h} \in \mathbb{R}^{S \times K}$ defining the factors are such that there exists $L \in \mathbb{N}$ such that $\mathbf{h} \in \mathcal{M}^L$. For instance, the constraint $\mathbf{h} \in \mathcal{M}^L$ might be used to impose sparsity, grouped sparsity or co-sparsity. One might also use the constraint $\mathbf{h} \in \mathcal{M}^L$ to impose non-negativity, orthogonality, equality, compactness, *etc.* Generally speaking \mathcal{M} is used to impose some prior or some form of regularity or to compress the parameter space and obtain better bounds [7]. The models might also

³Depending on the situation: Toeplitz, block-Toeplitz, circulant or block-circulant matrices. The matrices often involve downsampling.

⁴ReLU is the most common activation function.

be used to alleviate ambiguities. For instance, if the operators M_k and M_{k+1} allow permutations (i.e. there exists $(h_k, h_{k+1}) \neq (g_k, g_{k+1})$ and a permutation matrix C such that $M_k(h_k) = M_k(g_k)C$ and $M_{k+1}(h_{k+1}) = C^{-1}M_{k+1}(g_{k+1})$), we can use a complete ordering of the parameter space $\mathbb{R}^{K \times S}$ and impose, using \mathcal{M} , the largest of all the equivalent versions of a parameters to be considered.

1.3. Bibliography.

1.3.1. Other matrix factorization and compressed sensing. The content of this paper is strongly related to and can be considered as an extension of the research field usually named *compressed sensing*. Because of the importance of this field of research and to simplify the reading for readers whose main interest is in deep learning, we have separated this part of the bibliography and placed it in [section 2](#). Notice that the statement of [Informal theorem 1.1](#) can be interpreted in the context of signal recovery. In particular, the results on deep structured linear networks can probably be specialized to be applicable to matrix factorization problems for which stability properties have not been established [[77](#), [20](#), [21](#), [59](#), [58](#), [46](#), [56](#)]. We have not investigated this potential.

1.3.2. Tensors and deep networks. The analysis conducted in this paper is based on a connection, named *tensorial lifting*, between deep structured linear networks and a tensor problem (see [section 5](#)). The tensorial lifting has already been described in [[60](#)] but other connections between tensor and network problems have been described by other authors. In particular, in [[26](#), [27](#), [45](#)], the authors define a score function using a tensor. They highlight a network topology that computes the score function defined by a tensor decomposable using a CP-decomposition, a Hierarchical Tucker [[26](#), [27](#)] or a tensor train decomposition [[45](#)]. They then deduce the expressive power of the network topology from the connections between the tensor decompositions. These results highlight and analyze why deep networks are more expressive than shallow ones. Tensors and tensor decomposition have also been used to represent the cross-moment and construct a solver [[42](#)], encode the convolution layers with a tensor of order 4 and manipulate this tensor to improve the network [[49](#), [65](#), [82](#)], to represent a tensor layer [[73](#), [81](#)].

1.3.3. Stability property for neural networks. To the best of our knowledge, the articles establishing stability properties are [[4](#), [14](#), [53](#), [71](#), [83](#), [60](#)].

Among them, [[14](#), [53](#), [83](#)] consider a family of networks of depth 1 or 2 (depending on the article, the definition of the depth may vary). The article [[71](#)] contains a study on deep networks (the depth can be large), but the study only focuses on the recovery of one layer. The articles [[4](#), [60](#)] consider networks without depth limitation.

In [[14](#)], the authors consider the minimization of the statistical risk (not the empirical risk). The input is assumed Gaussian and the output is generated by a network involving one linear layer followed by ReLU and a mean. The number of intermediate nodes is smaller than the input size. They provide conditions guaranteeing that, with high probability, a randomly initialized gradient descent algorithm converges to the true parameters. The authors of [[53](#)] consider a feed forward network made of one unknown linear layer, followed by ReLU and a sum. The size of the intermediate layer equals the size of the data, the size of the output is 1. Again, they assume Gaussian input data and consider the minimization of the risk (not the empirical risk). They show that the stochastic gradient descent converges to the true solution. In [[83](#)], the authors consider a non-linear layer followed by a linear layer. The size of the intermediate layer is smaller than the size of the input and the size of the output is 1. They

describe an initialization algorithm based on a tensor decomposition such that with high probability, the gradient algorithm minimizing the empirical risk converges to the true parameters that generated the data.

The authors of [71] consider a feed-forward neural network and show that, if the input is Gaussian or its distribution is known, a method based on moments and sparse dictionary learning can retrieve the parameters defining the first layer. Nothing is said about the stability or the estimation of the other layers.

The authors of [4] consider deep feed-forward networks which are very sparse and randomly generated. They show that they can be learned with high probability one layer after another. However, very sparse and randomly generated networks are not used in practice and one might want to study more versatile structures. The article [60] studies deep structured linear networks (without the models \mathcal{M}) and uses the same tensorial lifting we use here. However, in [60] the function d measuring the error between parameters is only defined using the ℓ^∞ norm and is not a metric. The transversality condition of [60] is sufficient to guarantee the stability but is not necessary. All these weaknesses are corrected in this extended version. The general result is also specialized to deep convolutional linear networks.

1.4. Organization of the paper. Because it is strongly related, we give an extensive bibliography on compressive sensing and stable recovery properties for matrix factorization problems in section 2. We describe the framework of the paper and our notations in section 3.

The main contributions of this paper are:

- In section 4, we investigate and recall several results on tensors, tensor rank and the Segre embedding. In particular, we investigate how the Segre embedding distorts distances.
- In section 5, we describe the *tensorial lifting*. It expresses any deep structured linear networks in a generic multilinear format. The latter composes a linear lifting operator and the Segre embedding.
- When $\delta = \eta = 0$ (see section 6):
 - We establish a simple geometric condition on the intersection of two sets which is necessary and sufficient to guarantee the identifiability of the parameters up to scale ambiguity (Proposition 6.3).
 - We provide simpler conditions which involve the rank of the lifting operator (defined in section 5) such that:
 - * Under-specified case: If the lifting operator rank is large (e.g. larger than $2K(S-1) + 2$, when $\mathcal{M} = \mathbb{R}^{S \times K}$) and the lifting operator is random, for almost every lifting operator, the solution of

$$M_1(h_1) \cdots M_K(h_K) = X$$

is identifiable (Theorem 6.4).

- * Over-specified case: If the lifting operator rank is small (e.g. smaller than $2S - 1$, when $\mathcal{M} = \mathbb{R}^{S \times K}$), the solution of

$$M_1(h_1) \cdots M_K(h_K) = X$$

is not identifiable (Theorem 6.5);

- We also provide a simple algorithm to compute the rank of the lifting operator (Proposition 5.3).
- Stability guarantee statements for deep structured linear networks are in section 7:

- We define the deep-Null Space Property (Definition [Definition 7.1](#)): a generalization of the usual Null Space Property [[25](#)] that also applies to the deep problems.
- We establish that the deep-Null Space Property is a necessary and sufficient condition to guarantee stability (see the informal statement above or [Theorem 7.2](#) and [Theorem 7.3](#)).
- We specialize the results to convolutional linear networks in [section 8](#) and establish a simple condition that can be computed (see Algorithm [Algorithm 8.1](#)). The is such that (see [Theorem 8.4](#))
 - If the condition is satisfied the convolutional linear networks can be stably recovered;
 - If the condition is not satisfied, the convolutional linear network is not identifiable.

In simple words, the condition holds when the supports of the convolution kernels are sufficiently scattered. This is not satisfied by the convolutional kernel used in applications and explains their instability.

2. Bibliography on matrix factorization and compressed sensing. Before describing the bibliography on compressed sensing, we interpret this stability statement of [Informal theorem 1.1](#) in the context of signal processing. In signal processing, we usually know that \bar{h} exists and δ represents the sum of a modeling error and noise. The inequality [\(1.1\)](#) guarantees that, when the condition is satisfied, even an approximative minimizer of

$$(2.1) \quad \operatorname{argmin}_{L \in \mathbb{N}, (h_k)_{k=1..K} \in \mathcal{M}^L} \|M_1(h_1) \cdots M_K(h_K) - X\|^2.$$

leads to a solution h^* close \bar{h} . This property is often named: *stable recovery guarantee*.

When $\delta = 0$ (i.e., the data exactly fits the model and is not noisy) and $\eta = 0$ (i.e., [\(2.1\)](#) is perfectly solved) this is an *identifiability guarantee*. This is a necessary condition of stable recovery.

In this section, we distinguish the cases $K = 1$, $K = 2$ and $K \geq 3$.

2.1. $K = 1$: Linear inverse problems. The simplest version consists of a model with one layer (i.e., $K = 1$) and $\mathcal{M} = \mathbb{R}^{S \times K}$. Recovering h_1 from X is a linear inverse problem. The data X can be vectorized to form a column vector and the operator M_1 simply multiplies the column vector h_1 by a fixed (rectangular) matrix. Typically, when the linear inverse problem is over-determined, the latter matrix has more rows than columns, the uniqueness of a solution to [\(2.1\)](#) depends on the column rank of the matrix and the stable recovery constant depends on the smallest singular value of M_1 .

When the matrix is not full column rank, the identifiability and stable recovery for this problem has been intensively studied for many constraints \mathcal{M} . In particular, for sparsity constraints this is the compressed/compressive sensing problem (see the seminal articles [[15](#), [30](#)]). Some compressed sensing statements (especially the ones guaranteeing that any minimizer of the ℓ^0 problem stably recovers the unknown problem) are special cases ($K = 1$) of the statements provided in this paper. We will not give a complete review on compressed sensing but would like to highlight the Null Space Property described in [[25](#)]. The fundamental limits of compressed sensing (for a solution of the ℓ^0 problem) have been analyzed in detail in [[13](#)].

Although the main novelty of the paper is to investigate stable recovery properties for any $K \geq 1$, we specialize the statements made for $K \geq 1$ to the case $K = 1$ in

order to illustrate the new statements and to provide a way of comparison with well known results.

2.2. $K = 2$: Bilinear inverse problems and bilinear parameterizations.

When $k \geq 2$, the problem becomes non-linear because of the product in (2.1). This significantly complicates the analysis. What follows are the main instances studied in the literature when $K = 2$.

Non-negative Matrix factorization (NMF) and low rank prior. In non-negative matrix factorization [50], M_1 and M_2 map the entries in h_1 and h_2 at prescribed locations in the factors (say, one column after another). The constraints \mathcal{M} imposes that all the entries in h_1 and h_2 are non-negative. The NMF has been widely used for many applications.

Conditions guaranteeing that the factors provided by the NMF identify⁵ the correct factors (up to rescaling and permutation) were first established in the pioneering work [29]. To the best of our knowledge, this is the first paper addressing recovery guarantees for a problem of depth $K = 2$. It emphasizes a separability condition that guarantees identifiability. The proof is purely geometric and relies on the analysis of inclusions of simplicial cones. This result is significantly extended in [48]. In this paper, the continuity of the NMF estimator is established. Concerning computational aspects, NMF is NP-complete [78]. However, under the separability hypothesis of [29], the solution of the NMF problem can be computed in polynomial time [5].

We can slightly generalize⁶ the problem and introduce a linear degradation operator

$$H : \mathbb{R}^{m \times n} \longrightarrow \mathbb{R}^{m \times n}.$$

Use the same mapping M_1 and M_2 as for the NMF, with $\mathcal{M} = \mathbb{R}^{S \times K}$, but with a small number of lines (resp columns) in $M_2(h_2)$ (resp. $M_1(h_1)$). Any solution of the problem

$$(h_1^*, h_2^*) \in \operatorname{argmin}_{(h_1, h_2) \in \mathbb{R}^{S \times K}} \|H(M_1(h_1)M_2(h_2)) - X\|^2$$

leads to a low rank approximation $M_1(h_1^*)M_2(h_2^*)$ of an inverse of H , at X . Again, a large corpus of literature exists on the low rank prior [66, 17, 32, 19].

Phase Retrieval. Phase retrieval fits the framework described in the present paper when we take

$$M_1(h_1) = \operatorname{diag}(\mathcal{F}h_1) \quad M_2(h_2) = (\mathcal{F}h_2)^*$$

and

$$\mathcal{M} = \{(h, h) \in \mathbb{R}^{S \times K} \mid h \in \mathbb{R}^S\}$$

where S is the size of the signal, \mathcal{F} computes N linear measurements of any element in \mathbb{R}^S (typically Fourier measurements), $\operatorname{diag}(\cdot)$ creates an $N \times N$ diagonal matrix whose diagonal contains the input and $*$ is the (entry-wise) complex conjugate.

The tensorial lifting at the core of the present paper generalizes the lifting used in the inspiring work on PhaseLift [52, 18, 16]. As is often the case when $K = 2$, PhaseLift is a semidefinite program that can be efficiently solved when the unknown is of moderate size. These papers also provide conditions on the measurements guaranteeing that the phases are stably recovered by PhaseLift.

The benefit of the generalization introduced with the tensorial lifting is that it applies to any multilinear inverse problem.

⁵Stable recovery is not established.

⁶The interested readers can check that this generalization only leads to a small change of the Lifting operator introduced in section 5. It is therefore done at no cost.

Self-calibration and de-mixing. Measuring operators often depend linearly on parameters that are not perfectly known. The estimation of these parameters is crucial to restore the data measured by the device. This is the self-calibration problem. This naturally fits the setting of this article: - we let h_1 be the parameters defining the sensing matrix and $M_1(h_1)$ be the sensing matrix. Then h_2 defines the signal (or signals) contained in the column(s) of $M_2(h_2)$.

Many instances of this problem have been studied and much progress has been made to obtain algorithms that can be applied to problems of larger and larger size. This leads to a very interesting line of research.

To the best of our knowledge, the first stable recovery statements concern the blind-deconvolution problem. In [2], the authors use a lifting to transform the blind-deconvolution problem into a semidefinite program with an unknown whose size is the product of the sizes⁷ of h_1 and h_2 . Such problems can be solved for unknowns of moderate size. The authors of [2] provide explicit conditions guaranteeing the stable recovery with high probability. This idea has been generalized and applied to other similar problems in [24, 9]. The authors of [54] consider a significantly more general calibration model. In this model, $M_1(h_1)$ is diagonal and its diagonal contains the entries of h_1 . $M_2(h_2)$ simply multiplies h_2 by a fixed known matrix (the theorems consider a random matrix). The constraint imposes h_2 to be sparse. For this problem, they prove that with high probability the numerical method called SparseLift is stable with a controlled accuracy. SparseLift returns the left and right singular vectors of the solutions of an ℓ^1 optimization problem whose unknown is the same as in [2]. However, solving an ℓ^1 minimization problem is much simpler than a semi-definite problem. This is a very significant practical improvement.

As emphasized in [51], in order to motivate its non-convex approach, the only drawback of the numerical methods described in [2, 54] is their complexity. The extra complexity is due to the fact that they optimize a variable in the product space $\mathbb{R}^{S \times S}$ and then deduce an approximate solution of the un-lifted problem. This is what motivates the authors of [51] to propose a non-convex approach. The constructed algorithm provably stably recovers the sensing parameters and the signals with a geometric convergence rate.

Sparse coding and dictionary learning: Sparse coding and dictionary learning is another kind of bilinear problem (see [67] for an overview). In that framework, the columns of X contain the data. Most often, people consider two layers: $K = 2$. The layer $M_1(h_1)$ is an optimized dictionary of atoms defined by the parameters h_1 and each column of $M_2(h_2)$ contains the code (or coordinates) of the corresponding column in X . Most often, h_2 is assumed sparse.

The identifiability and stable recovery of the factors has been studied in many dictionary learning contexts and provides guarantees on the approximate recovery of both an incoherent dictionary and sparse coefficients when the number of samples is sufficiently large (i.e., in our setting when n is large). In [36], the authors developed local optimality conditions in the noiseless case, as well as sample complexity bounds for local recovery when $M_1(h_1)$ is square and $M_2(h_2)$ are iid Bernoulli-Gaussian. This was extended to overcomplete dictionaries in [33] (see also [70] for tight frames) and to the noisy case in [43]. The authors of [74] provide exact recovery results for dictionary learning, when the coefficient matrix has Bernoulli-Gaussian entries and the dictionary matrix has full column rank. This was extended to overcomplete dictionaries in [1] and in [6] but only for approximate recovery. Finally, [35] provides such guarantees

⁷With our notations this is simply $S \times S$ but this can be much more favorable.

under general conditions which cover many practical settings.

Contributions in these frameworks. The present article considers the identifiability and stability of the recovery for any $K \geq 1$ in a general and unifying framework. As was already mentioned, we do not investigate computational issues. As will appear, later the paper, the analogue of the lifting at the core of the algorithms described in the above papers (in particular the papers on phase retrieval and self-calibration) is a *tensorial lifting* (see [section 5](#)) and involves tensors that cannot be manipulated in practice. Even when we are able to manipulate the tensors, the computation of the best rank 1 approximation of such tensors is an open non-convex problem. Therefore, there is no numerically efficient and reliable way to extract the un-lifted parameters from an optimized tensor. Because of that, we have not yet pursued the construction of a numerical scheme based on the tensorial lifting when $K \geq 3$. As was already mentioned, at this writing, the success of algorithms for $K \geq 3$ is mostly supported by empirical evidence. Proving their efficiency is a wide open problem (see [\[55, 37, 22, 23, 68, 10, 11, 44, 79\]](#)). The purpose of the paper is to provide guarantees on the stability of the solution when such an empirical success occurs.

The specialization of the presented results to problems with $K = 2$ leads to necessary and sufficient conditions for the stable recovery. This is slightly different from the usual approach. Usually, authors provide sufficient conditions and argue their sharpness by comparing the number of samples required by their method and the information theoretic limit (typically, the number of independent variables of the problem).

It would of course be interesting to see how far it is possible to unify the different problems with $K = 2$ using the framework of this paper. We have however not pursued this route and instead focus on the situation $K \geq 3$.

2.3. $K \geq 3$. The difficulties, when $K \geq 3$, come from the fact that tools used for problems with $K = 2$ are not applicable. In particular, we cannot use the usual lifting, the singular value decomposition or the $\sin-\theta$ theorem in [\[28\]](#). Often, these tools are replaced by analogous objects involving tensors. This complicates the analysis and prohibits the use of numerical schemes that manipulate lifted variables.

To the best of our knowledge, little is known concerning the identifiability and the stability of matrix factorization when $K \geq 3$. The uniqueness of the factorization corresponding to the Fast Fourier Transform was proved in [\[57\]](#). Other results consider the identifiability of the factors which are sparse and random [\[64\]](#). The authors of the present paper have announced preliminary versions of the results described here in [\[60\]](#). They are significantly extended here.

3. Notation and summary of the hypotheses. We continue to use the notation introduced in the introduction. For an integer $k \in \mathbb{N}$, set $[k] = \{1, \dots, k\}$.

We consider $K \geq 1$ and $S \geq 2$ and real-valued tensors of order K whose axes are of size S , denoted $T \in \mathbb{R}^{S \times \dots \times S}$. The space of tensors is abbreviated \mathbb{R}^{S^K} . The entries of T are denoted T_{i_1, \dots, i_K} , where $(i_1, \dots, i_K) \in [S]^K$. For $\mathbf{i} \in [S]^K$, the entries of \mathbf{i} are $i = (i_1, \dots, i_K)$ (for $\mathbf{j} \in [S]^K$ we let $\mathbf{j} = (j_1, \dots, j_K)$, etc). We either write $T_{\mathbf{i}}$ or T_{i_1, \dots, i_K} .

To simplify notations, from now on, the parameters defining the factors are gathered in a single matrix and are denoted with bold fonts $\mathbf{h} \in \mathbb{R}^{S \times K}$. The k^{th} vector containing the parameters for the layer k is denoted $\mathbf{h}_k \in \mathbb{R}^S$. The i^{th} entry of the k^{th} vector is denoted $\mathbf{h}_{k,i} \in \mathbb{R}$. A vector not related to an element in $\mathbb{R}^{S \times K}$ is denoted

$h \in \mathbb{R}^S$ (i.e., using a light font). Throughout the paper we assume

$$\mathcal{M} = (\mathcal{M}^L)_{L \in \mathbb{N}}, \text{ with } \mathcal{M}^L \subset \mathbb{R}^{S \times K}.$$

We also assume that, for all $L \in \mathbb{N}$, $\mathcal{M}^L \neq \emptyset$. They can however be equal or constant after a given L' .

All the vector spaces \mathbb{R}^{S^K} , $\mathbb{R}^{S \times K}$, \mathbb{R}^S etc. are equipped with the usual Euclidean norm. This norm is denoted $\|\cdot\|$ and the scalar product $\langle \cdot, \cdot \rangle$. In the particular case of matrices, $\|\cdot\|$ corresponds to the Frobenius norm. We also use the usual p norm, for $p \in [1, \infty]$, and denote it by $\|\cdot\|_p$. In particular, for $\mathbf{h} \in \mathbb{R}^{S \times K}$ and $T \in \mathbb{R}^{S^K}$, we have for $p < +\infty$

$$\|\mathbf{h}\|_p = \left(\sum_{k=1}^K \sum_{i=1}^S |\mathbf{h}_{k,i}|^p \right)^{1/p}, \quad \|T\|_p = \left(\sum_{\mathbf{i} \in [S]^K} |T_{\mathbf{i}}|^p \right)^{1/p}$$

and

$$\|\mathbf{h}\|_{\infty} = \max_{\substack{k \in [K] \\ i \in [S]}} |\mathbf{h}_{k,i}|, \quad \|T\|_{\infty} = \max_{\mathbf{i} \in [S]^K} |T_{\mathbf{i}}|.$$

Set

$$(3.1) \quad \mathbb{R}_*^{S \times K} = \{\mathbf{h} \in \mathbb{R}^{S \times K} \mid \forall k \in [K], \|\mathbf{h}_k\| \neq 0\}.$$

Define an equivalence relation on $\mathbb{R}_*^{S \times K}$: for any $\mathbf{h}, \mathbf{g} \in \mathbb{R}_*^{S \times K}$, $\mathbf{h} \sim \mathbf{g}$ if and only if there exist $(\lambda_k)_{k \in [K]} \in \mathbb{R}^K$ such that

$$(3.2) \quad \prod_{k=1}^K \lambda_k = 1 \quad \text{and} \quad \mathbf{h}_k = \lambda_k \mathbf{g}_k, \forall k \in [K].$$

Denote the equivalence class of $\mathbf{h} \in \mathbb{R}_*^{S \times K}$ by $\langle \mathbf{h} \rangle$.

The zero tensor is of rank 0. A non-zero tensor $T \in \mathbb{R}^{S^K}$ is of *rank* 1 (or decomposable) if and only if there exists $\mathbf{h} \in \mathbb{R}_*^{S \times K}$ such that T is the outer product of the vectors \mathbf{h}_k , for $k \in [K]$. That is, for any $\mathbf{i} \in [S]^K$,

$$T_{\mathbf{i}} = \mathbf{h}_{1,i_1} \cdots \mathbf{h}_{K,i_K}.$$

Let $\Sigma_1 \subset \mathbb{R}^{S^K}$ denote the set of tensors of rank 0 or 1.

The *rank* of a tensor $T \in \mathbb{R}^{S^K}$ is

$$\text{rk}(T) = \min\{r \in \mathbb{N} \mid \text{there exists } T_1, \dots, T_r \in \Sigma_1 \text{ such that } T = T_1 + \dots + T_r\}.$$

For $r \in \mathbb{N}$, let

$$\Sigma_r = \{T \in \mathbb{R}^{S^K} \mid \text{rk}(T) \leq r\}.$$

The $*$ superscript refers to optimal solutions. A set with a $*$ subscript means that 0 is ruled out of the set. In particular, $\Sigma_{1,*}$ denotes the non-zero tensors of rank 1. Attention should be paid to $\mathbb{R}_*^{S \times K}$ (see (3.1)).

4. Facts on the Segre embedding and tensors of rank 1 and 2. Parametrize $\Sigma_1 \subset \mathbb{R}^{S \times K}$ by the map

$$(4.1) \quad \begin{aligned} P : \mathbb{R}^{S \times K} &\longrightarrow \Sigma_1 \subset \mathbb{R}^{S \times K} \\ \mathbf{h} &\longmapsto (\mathbf{h}_{1,i_1} \mathbf{h}_{2,i_2} \cdots \mathbf{h}_{K,i_K})_{i \in [S]^K}. \end{aligned}$$

The map P is called the Segre embedding and is often denoted \widehat{Seg} in the algebraic geometry literature.

Standard Facts:

1. **Identifiability of $\langle \mathbf{h} \rangle$ from $P(\mathbf{h})$:** For \mathbf{h} and $\mathbf{g} \in \mathbb{R}_*^{S \times K}$, $P(\mathbf{h}) = P(\mathbf{g})$ if and only if $\langle \mathbf{h} \rangle = \langle \mathbf{g} \rangle$.
2. **Geometrical description of $\Sigma_{1,*}$:** $\Sigma_{1,*}$ is a smooth (i.e., C^∞) manifold of dimension $K(S-1) + 1$ (see, e.g., [47], chapter 4, pp. 103).
3. **Geometrical description of Σ_2 :** We recall that the singular locus $(\overline{\Sigma}_2)_{sing}$ of the closure $\overline{\Sigma}_2$ of Σ_2 has dimension strictly less than that of $\overline{\Sigma}_2$ and that $\overline{\Sigma}_2 \setminus (\overline{\Sigma}_2)_{sing}$ is a smooth manifold. The dimension of $\overline{\Sigma}_2 \setminus (\overline{\Sigma}_2)_{sing}$ is $2K(S-1) + 2$ when $K > 2$, and is $4(S-1)$ when $K = 2$ (see, e.g., [47], chapter 5).

We can improve Standard Fact [Item 1](#) and obtain a stability result guaranteeing that if we know a rank 1 tensor sufficiently close to $P(\mathbf{h})$, we approximately know $\langle \mathbf{h} \rangle$. In order to state this, we need to define a metric on $\mathbb{R}_*^{S \times K} / \sim$ (where \sim is defined by [\(3.2\)](#)). This has to be considered with care since, whatever $\mathbf{h} \in \mathbb{R}_*^{S \times K}$, the subset $\{h \mid h \in \langle \mathbf{h} \rangle\}$ is not compact. In particular, considering

$$\mathbf{h}'_k = \begin{cases} \lambda \mathbf{h}_k & \text{if } k = 1 \\ \lambda^{-\frac{1}{K-1}} \mathbf{h}_k & \text{otherwise} \end{cases}$$

when λ goes to infinity, we easily construct examples that make the standard metric on equivalence classes useless⁸.

This leads us to consider

$$\mathbb{R}_{diag}^{S \times K} = \{\mathbf{h} \in \mathbb{R}_*^{S \times K} \mid \forall k \in [K], \|\mathbf{h}_k\|_\infty = \|\mathbf{h}_1\|_\infty\}.$$

The interest in this set comes from the fact that, whatever $\mathbf{h} \in \mathbb{R}_*^{S \times K}$, the set $\langle \mathbf{h} \rangle \cap \mathbb{R}_{diag}^{S \times K}$ is finite. Indeed, if $\mathbf{g} \in \langle \mathbf{h} \rangle \cap \mathbb{R}_{diag}^{S \times K}$ the $(\lambda_k)_{k \in [K]} \in \mathbb{R}^K$ such that, for all $k \in [K]$, $\mathbf{h}_k = \lambda_k \mathbf{g}_k$ must all satisfy $|\lambda_k| = 1$, i.e., $\lambda_k = \pm 1$.

DEFINITION 4.1. For any $p \in [1, \infty]$, we define the mapping $d_p : (\mathbb{R}_*^{S \times K} / \sim \times \mathbb{R}_*^{S \times K} / \sim) \rightarrow \mathbb{R}$ by

$$d_p(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) = \inf_{\substack{\mathbf{h}' \in \langle \mathbf{h} \rangle \cap \mathbb{R}_{diag}^{S \times K} \\ \mathbf{g}' \in \langle \mathbf{g} \rangle \cap \mathbb{R}_{diag}^{S \times K}}} \|\mathbf{h}' - \mathbf{g}'\|_p \quad \forall \mathbf{h}, \mathbf{g} \in \mathbb{R}_*^{S \times K}.$$

⁸For instance, if \mathbf{h} and $\mathbf{g} \in \mathbb{R}_*^{S \times K}$ are such that $\mathbf{h}_1 = \mathbf{g}_1$, we have

$$\inf_{\mathbf{h}' \in \langle \mathbf{h} \rangle, \mathbf{g}' \in \langle \mathbf{g} \rangle} \|\mathbf{h}' - \mathbf{g}'\|_p = 0$$

even though we might have $\mathbf{h}_2 \neq \mathbf{g}_2$ (and therefore $\langle \mathbf{h} \rangle \neq \langle \mathbf{g} \rangle$). This does not define a metric.

Also, when \mathbf{h} and \mathbf{g} are such that $\mathbf{h}_k \neq \mathbf{g}_k$, whatever $k \in [K]$, we have

$$\sup_{\mathbf{h}' \in \langle \mathbf{h} \rangle} \inf_{\mathbf{g}' \in \langle \mathbf{g} \rangle} \|\mathbf{h}' - \mathbf{g}'\|_p = +\infty.$$

Therefore, the Hausdorff distance between $\langle \mathbf{h} \rangle$ and $\langle \mathbf{g} \rangle$ is infinite for almost every pair (\mathbf{h}, \mathbf{g}) . This metric is therefore not very useful in the present context.

PROPOSITION 4.2. *For any $p \in [1, \infty]$, d_p is a metric on $\mathbb{R}_*^{S \times K} / \sim$.*

The proof is in [subsection 10.1](#).

Notice that the equivalence relationship and metric defined above are not adapted to operators M_k allowing invariance such as permutations. More precisely, for some operators M_k , there exists \mathbf{h}, \mathbf{g} and a permutation matrix C such that $(\mathbf{h}_k, \mathbf{h}_{k+1}) \neq (\mathbf{g}_k, \mathbf{g}_{k+1})$ and $M_k(\mathbf{h}_k) = M_k(\mathbf{g}_k)C$ and $M_{k+1}(\mathbf{h}_{k+1}) = C^{-1}M_{k+1}(\mathbf{g}_{k+1})$. In such a case, we have $d_p(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) \neq 0$. However, the features defined at the layer k are just permuted and can still be interpreted. As already said, in such a case, it is possible to use the models \mathcal{M} to select one the equally interpretable \mathbf{h} .

Using the above metric, we can state that not only $\langle \mathbf{h} \rangle$ is uniquely determined by $P(\mathbf{h})$, but this operation is stable.

THEOREM 4.3. **Stability of $\langle \mathbf{h} \rangle$ from $P(\mathbf{h})$**

Let \mathbf{h} and $\mathbf{g} \in \mathbb{R}_*^{S \times K}$ be such that $\|P(\mathbf{g}) - P(\mathbf{h})\|_\infty \leq \frac{1}{2} \max(\|P(\mathbf{h})\|_\infty, \|P(\mathbf{g})\|_\infty)$.
For all $p, q \in [1, \infty]$,

$$(4.2) \quad d_p(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) \leq 7(KS)^{\frac{1}{p}} \min\left(\|P(\mathbf{h})\|_\infty^{\frac{1}{K}-1}, \|P(\mathbf{g})\|_\infty^{\frac{1}{K}-1}\right) \|P(\mathbf{h}) - P(\mathbf{g})\|_q.$$

The proof of the theorem is in [subsection 10.2](#).

In the final result, the bound established in [Theorem 4.3](#) plays a role similar to the *sin* - θ Theorem of [28] in [54, 18, 2].

The following proposition shows that the upper bound in (4.2) cannot be improved by a significant factor, in particular when q is large.

PROPOSITION 4.4. *There exist \mathbf{h} and $\mathbf{g} \in \mathbb{R}_*^{S \times K}$ such that $\|P(\mathbf{g})\|_\infty \leq \|P(\mathbf{h})\|_\infty$, $\|P(\mathbf{g}) - P(\mathbf{h})\|_\infty \leq \frac{1}{2} \|P(\mathbf{h})\|_\infty$ and*

$$7(KS)^{\frac{1}{p}} \|P(\mathbf{h})\|_\infty^{\frac{1}{K}-1} \|P(\mathbf{h}) - P(\mathbf{g})\|_q \leq C_q d_p(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle),$$

where

$$C_q = \begin{cases} 28(KS)^{\frac{1}{q}} & \text{if } q < +\infty, \\ 28 & \text{if } q = +\infty. \end{cases}$$

The proof of the theorem is in [subsection 10.3](#).

As stated in the following theorem, we have a more valuable upper bound in the general case.

THEOREM 4.5. **“Lipschitz continuity” of P**

For any $q \in [1, \infty]$ and any \mathbf{h} and $\mathbf{g} \in \mathbb{R}_*^{S \times K}$,

$$(4.3) \quad \|P(\mathbf{h}) - P(\mathbf{g})\|_q \leq S^{\frac{K-1}{q}} K^{1-\frac{1}{q}} \max\left(\|P(\mathbf{h})\|_\infty^{1-\frac{1}{K}}, \|P(\mathbf{g})\|_\infty^{1-\frac{1}{K}}\right) d_q(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle).$$

The theorem is proved in [subsection 10.4](#).

Notice that, considering \mathbf{h} and $\mathbf{g} \in \mathbb{R}_*^{S \times K}$ such that $\mathbf{h}_{k,i} = 1$ and $\mathbf{g}_{k,i} = \varepsilon$, for all $k \in [K]$ and $i \in [S]$ and for a $0 < \varepsilon \ll 1$, we easily calculate

$$S^{\frac{K-1}{q}} K^{1-\frac{1}{q}} \max\left(\|P(\mathbf{h})\|_\infty^{1-\frac{1}{K}}, \|P(\mathbf{g})\|_\infty^{1-\frac{1}{K}}\right) d_q(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) \leq K \|P(\mathbf{h}) - P(\mathbf{g})\|_q.$$

As a consequence, the upper bound in [Theorem 4.5](#) is tight up to at most a factor K .

5. The tensorial lifting . The following proposition is clear (it can be shown by induction on K):

PROPOSITION 5.1. *Let M_k , $k \in [K]$ be as in (1.2). The entries of the matrix*

$$M_1(\mathbf{h}_1)M_2(\mathbf{h}_2) \cdots M_K(\mathbf{h}_K)$$

are multivariate polynomials whose variables are the entries of $\mathbf{h} \in \mathbb{R}^{S \times K}$. Moreover, every entry is the sum of monomials of degree K . Each monomial is a constant times $\mathbf{h}_{1,i_1} \cdots \mathbf{h}_{K,i_K}$, for some $\mathbf{i} \in [S]^K$.

Notice that any monomial $\mathbf{h}_{1,i_1} \cdots \mathbf{h}_{K,i_K}$ is the entry $P(\mathbf{h})_{\mathbf{i}}$ in the tensor $P(\mathbf{h})$. Therefore every polynomial in the previous proposition takes the form $\sum_{\mathbf{i} \in [S]^K} c_{\mathbf{i}} P(\mathbf{h})_{\mathbf{i}}$ for some constants $(c_{\mathbf{i}})_{\mathbf{i} \in [S]^K}$ independent of \mathbf{h} . In words, every entry of the matrix $M_1(\mathbf{h}_1)M_2(\mathbf{h}_2) \cdots M_K(\mathbf{h}_K)$ is obtained by applying a linear form to $P(\mathbf{h})$. Moreover, the polynomial coefficients defining the linear form are uniquely determined by the linear maps M_1, \dots, M_K . This leads to the following statement.

COROLLARY 5.2. Tensorial Lifting

Let M_k , $k \in [K]$ be as in (1.2). The map

$$(\mathbf{h}_1, \dots, \mathbf{h}_K) \mapsto M_1(\mathbf{h}_1)M_2(\mathbf{h}_2) \cdots M_K(\mathbf{h}_K),$$

uniquely determines a linear map

$$\mathcal{A} : \mathbb{R}^{S^K} \longrightarrow \mathbb{R}^{m \times n},$$

such that for all $\mathbf{h} \in \mathbb{R}^{S \times K}$

$$(5.1) \quad M_1(\mathbf{h}_1)M_2(\mathbf{h}_2) \cdots M_K(\mathbf{h}_K) = \mathcal{A}P(\mathbf{h}).$$

We call (5.1) and its use the *tensorial lifting*. When $K = 1$, we simply have $\mathcal{A} = M_1$. When $K = 2$ it corresponds to the usual lifting already exploited to establish stability results for phase recovery, blind-deconvolution, self-calibration, sparse coding, *etc.* Notice that, when $K \geq 2$, it may be difficult to provide a closed form expression for the operator \mathcal{A} . We can however determine simple properties of \mathcal{A} . In most reasonable cases, \mathcal{A} is sparse. If the operators M_k simply embed the values of h in a matrix, the matrix representing \mathcal{A} only contains zeros and ones. Since the operators M_k are known, we can compute $\mathcal{A}P(\mathbf{h})$, for any $\mathbf{h} \in \mathbb{R}^{S \times K}$, using (5.1). Said differently, we can compute \mathcal{A} for any rank 1 tensor. Therefore, since \mathcal{A} is linear, we can compute $\mathcal{A}T$ for any low rank tensor T . If the dimensions of the problem permit, one can manipulate \mathcal{A} in a basis of \mathbb{R}^{S^K} .

Since $\text{rk}(\mathcal{A})$ is an important quantity, we emphasize that $\text{rk}(\mathcal{A}) \leq mn$. It is also possible to compute $\text{rk}(\mathcal{A})$, when mn is not too large, using the following proposition.

PROPOSITION 5.3. *For R independent random \mathbf{h}^r , with $r = 1..R$, according to the normal distribution in $\mathbb{R}^{S \times K}$, we have with probability 1*

$$(5.2) \quad \dim(\text{Span}((\mathcal{A}P(\mathbf{h}^r))_{r=1..R})) = \begin{cases} R & \text{if } R \leq \text{rk}(\mathcal{A}) \\ \text{rk}(\mathcal{A}) & \text{otherwise.} \end{cases}$$

The proof is in [subsection 10.5](#).

Using [Corollary 5.2](#), when (2.1) has a minimizer, we rewrite in the form

$$(5.3) \quad \mathbf{h}^* \in \underset{L \in \mathbb{N}, \mathbf{h} \in \mathcal{M}^L}{\text{argmin}} \|\mathcal{A}P(\mathbf{h}) - X\|^2.$$

We now decompose this problem into two sub-problems: A least-squares problem

$$(5.4) \quad T^* \in \operatorname{argmin}_{T \in \mathbb{R}^{S^K}} \|\mathcal{A}T - X\|^2$$

and a non-convex problem

$$(5.5) \quad \mathbf{h}^* \in \operatorname{argmin}_{L \in \mathbb{N}, \mathbf{h} \in \mathcal{M}^L} \|\mathcal{A}(P(\mathbf{h}) - T^*)\|^2.$$

PROPOSITION 5.4. *Let X and \mathcal{A} be such that (2.1) has a minimizer:*

1. *Let \mathbf{h}^* be a solution of (5.3). Then, for any solution T^* of (5.4), \mathbf{h}^* also minimizes (5.5).*
2. *Let T^* be a solution of (5.4) and \mathbf{h}^* a solution of (5.5). Then, \mathbf{h}^* also minimizes (5.3).*

The proof is in [subsection 10.6](#).

From now on, because of the equivalence between solutions of (5.5) and (5.3), we stop using the notation \mathbf{h}'^* and write $\mathbf{h}^* \in \operatorname{argmin}_{L \in \mathbb{N}, \mathbf{h} \in \mathcal{M}^L} \|\mathcal{A}(P(\mathbf{h}) - T^*)\|^2$.

6. Identifiability (error free case). Throughout this section, we assume that X is such that there exists \bar{L} and $\bar{\mathbf{h}} \in \mathcal{M}^{\bar{L}}$ such that

$$(6.1) \quad X = M_1(\bar{\mathbf{h}}_1) \cdots M_K(\bar{\mathbf{h}}_K).$$

Under this assumption, $X = \mathcal{A}P(\bar{\mathbf{h}})$, so

$$P(\bar{\mathbf{h}}) \in \operatorname{argmin}_{T \in \mathbb{R}^{S^K}} \|\mathcal{A}T - X\|^2.$$

Moreover, we trivially have $P(\bar{\mathbf{h}}) \in \Sigma_1$ and therefore $\bar{\mathbf{h}}$ minimizes (5.5), (2.1) and (5.3). As a consequence, (2.1) has a minimizer.

We ask whether there exist guarantees that the resolution of (2.1) allows one to recover $\bar{\mathbf{h}}$ up to the usual uncertainties.

In this regard, for any $\mathbf{h} \in \langle \bar{\mathbf{h}} \rangle$, we have $P(\mathbf{h}) = P(\bar{\mathbf{h}})$ and therefore $\mathcal{A}P(\mathbf{h}) = \mathcal{A}P(\bar{\mathbf{h}}) = X$. Thus unless we make further assumptions on $\bar{\mathbf{h}}$, we cannot expect to distinguish any particular element of $\langle \bar{\mathbf{h}} \rangle$ using only X . In other words, recovering $\langle \bar{\mathbf{h}} \rangle$ is the best we can hope for.

DEFINITION 6.1. **Identifiability**

We say that $\langle \bar{\mathbf{h}} \rangle$ is identifiable if the elements of $\langle \bar{\mathbf{h}} \rangle$ are the only solutions of (2.1).

We say that \mathcal{M} is identifiable if for every $L \in \mathbb{N}$ and every $\bar{\mathbf{h}} \in \mathcal{M}^L$, $\langle \bar{\mathbf{h}} \rangle$ is identifiable.

PROPOSITION 6.2. **Characterization of the global minimizers**

For any $L^ \in \mathbb{N}$ and any $\mathbf{h}^* \in \mathcal{M}^{L^*}$, $(L^*, \mathbf{h}^*) \in \operatorname{argmin}_{L \in \mathbb{N}, \mathbf{h} \in \mathcal{M}} \|\mathcal{A}P(\mathbf{h}) - X\|^2$ if and only if*

$$P(\mathbf{h}^*) \in P(\bar{\mathbf{h}}) + \operatorname{Ker}(\mathcal{A}).$$

The Proposition is proved in [subsection 10.7](#).

In order to state the following proposition, we define for any L and $L' \in \mathbb{N}$

$$P(\mathcal{M}^L) - P(\mathcal{M}^{L'}) := \left\{ P(\mathbf{h}) - P(\mathbf{g}) \mid \mathbf{h} \in \mathcal{M}^L \text{ and } \mathbf{g} \in \mathcal{M}^{L'} \right\} \subset \mathbb{R}^{S^K}.$$

PROPOSITION 6.3. **Necessary and sufficient conditions of identifiability**

1. For any \bar{L} and $\bar{\mathbf{h}} \in \mathcal{M}^{\bar{L}}$: $\langle \bar{\mathbf{h}} \rangle$ is identifiable if and only if for any $L \in \mathbb{N}$

$$(P(\bar{\mathbf{h}}) + \text{Ker}(\mathcal{A})) \cap P(\mathcal{M}^L) \subset \{P(\bar{\mathbf{h}})\}.$$

2. \mathcal{M} is identifiable if and only if for any L and $L' \in \mathbb{N}$

$$(6.2) \quad \text{Ker}(\mathcal{A}) \cap (P(\mathcal{M}^L) - P(\mathcal{M}^{L'})) \subset \{0\}.$$

The proposition is proved in [subsection 10.8](#).

In the context of the usual compressed sensing (i.e., when $K = 1$, \mathcal{M} contains L -sparse signals, \mathcal{A} is a rectangular matrix with full row rank and X is a vector), the proposition is already stated in Lemma 3.1 of [\[25\]](#).

In reasonably small cases and when $P(\mathcal{M})$ is algebraic, one can use tools from numerical algebraic geometry such as those described in [\[39, 40\]](#) to check whether the condition [\(6.2\)](#) holds or not. The drawback of [Proposition 6.3](#) is that, given a deep structured linear network as described by \mathcal{A} , the condition [\(6.2\)](#) might be difficult to verify.

We therefore establish simpler conditions related to the identifiability of \mathcal{M} . First we establish a condition such that for almost every \mathcal{A} satisfying it, \mathcal{M} is identifiable. The main benefit of this condition is that its constituents can be computed in many practical situations.

Before that, we recall a few facts of algebraic geometry, for $X, Y \subset \mathbb{R}^N$, the *join* of X and Y (see, e.g., [\[38, Ex. 8.1\]](#)) is

$$J(X, Y) := \overline{\{sx + ty \mid x \in X, y \in Y, s, t \in \mathbb{R}\}}.$$

If for all $L \in \mathbb{N}$, \mathcal{M}^L is Zariski closed and invariant under rescaling (e.g., if they are all linear spaces), then $P(\mathcal{M}^L) - P(\mathcal{M}^{L'})$ is a Zariski open subset of $J(P(\mathcal{M}^L), P(\mathcal{M}^{L'}))$. In general, it is contained in this join.

Recall the following fact (*): for complex algebraic varieties $X, Y \subset \mathbb{C}^N$, any component Z of $X \cap Y$ has $\dim(Z) \geq \dim(X) + \dim(Y) - N$, and equality holds generically (we make “generically” precise in our context below). Moreover, if X, Y are invariant under rescaling, since $0 \in X \cap Y$, we have $X \cap Y \neq \emptyset$. (See, e.g., [\[72, §I.6.2\]](#).)

This intersection result indicates that if there exists L, L' such that

$$\text{rk}(\mathcal{A}) < \dim(P(\mathcal{M}^L) - P(\mathcal{M}^{L'}))$$

we expect to have non-identifiability; and if the rank is larger, for all pair L, L' , we expect identifiability. More precisely:

THEOREM 6.4. Almost surely sufficient condition for Identifiability

For almost every \mathcal{A} such that

$$\text{rk}(\mathcal{A}) \geq \dim(J(P(\mathcal{M}^L), P(\mathcal{M}^{L'}))), \quad \text{for all } L, L',$$

\mathcal{M} is identifiable.

The theorem is proved in [subsection 10.9](#).

Since $\dim(J(P(\mathcal{M}^L), P(\mathcal{M}^{L'}))) \leq \dim(P(\mathcal{M}^L)) + \dim(P(\mathcal{M}^{L'}))$, if D_{max} is the maximum dimension of $P(\mathcal{M}^L)$ over all L , one has the same conclusion if $\text{rk}(\mathcal{A}) \geq 2D_{max}$.

When $K = 1$, we illustrate this result by interpreting it in the context of compressive sensing, where \mathbf{h} is a vector, X is a vector, \mathcal{A} is a rectangular sampling matrix of full row rank and $\text{Ker}(\mathcal{A})$ is large. The statement analogous to [Theorem 6.4](#) in the compressive sensing framework takes the form: “For almost every sampling matrix, any L sparse signal \mathbf{h} can be recovered from $\mathcal{A}\mathbf{h}$ as soon as $2L \leq \text{rk}(\mathcal{A})$.” Moreover, the constituent of the ℓ^0 minimization model used to recover the signal are also the constituents of [\(5.3\)](#). Again, the main novelty is to extend this result to the identifiability of the factors of deep matrix products.

In order to establish a necessary condition for identifiability, first note that if we extend $P(\mathcal{M}^L) - P(\mathcal{M}^{L'})$ to be scale invariant, this will not affect whether or not it intersects $\text{ker}(\mathcal{A})$ outside of the origin. We immediately conclude that in the complex setting where $\mathcal{M}^L, \mathcal{M}^{L'}$ are both Zariski closed, that \mathcal{M} is non-identifiable whenever $\text{rk}(\mathcal{A}) < \dim(P(\mathcal{M}^L) - P(\mathcal{M}^{L'}))$. This indicates that we should always expect non-identifiability whenever $\text{rk}(\mathcal{A}) < \dim(P(\mathcal{M}^L) - P(\mathcal{M}^{L'}))$ but is not adequate to prove it because real algebraic varieties need not satisfy (*). However it is true for real linear spaces, so we immediately conclude the following weak result:

THEOREM 6.5. Necessary condition for Identifiability

Let $C(P(\mathcal{M}^L) - P(\mathcal{M}^{L'}))$ be the set of all points on all lines through the origin intersecting $P(\mathcal{M}^L) - P(\mathcal{M}^{L'})$, and let q be the maximal dimension of a linear space on $C(P(\mathcal{M}^L) - P(\mathcal{M}^{L'}))$. Then if $q > \text{rk}(\mathcal{A})$, \mathcal{M} is not identifiable. In particular when the $P(\mathcal{M}^L)$'s contain linear space and if we let S' be the largest dimension of these vector space, if $2S' > \text{rk}(\mathcal{A})$, then \mathcal{M} is not identifiable.

Let us illustrate the theorems by considering a deep feed-forward ReLU network and consider the structured linear network obtained by fixing the action of ReLU (as is done in [subsection 1.2.2](#)). The matrix X contains the outputs, the operator M_K multiplies the matrix containing the inputs by the weights between the first and the second layer. For every input/output pair, the action of ReLU is different and removes paths from the input entries to the output entries. We assume however that every entry of every output is reached by at least one path in the network starting at a non-zero entry of the input. In that case, it is not difficult to see that \mathcal{A} is a surjection and therefore $\text{rk}(\mathcal{A}) = mn$, where m is the size of the output and n is the number of learning samples.

The condition in [Theorem 6.4](#) becomes

$$mn \geq \dim(\Sigma_2) = 2K(S - 1) + 2,$$

and KS is typically the number of parameters of the network. The intuition behind [Theorem 6.4](#) is that, if the action of ReLU is sufficiently random and if the above inequality holds, we can expect the network to be identifiable with high probability⁹. This situation corresponds to an under-parameterized case (favorable for identifiability).

The condition in [Theorem 6.5](#) is

$$2S > mn.$$

When this inequality holds the network is not identifiable. It corresponds to an over-parameterized configuration. In the intermediate situation, when $2S \leq mn <$

⁹This statement gives the intuition behind [Theorem 6.4](#) but it should be made precise, as emphasized in the perspectives of this paper.

$2K(S - 1) + 2$, and when the action of the activation function does not introduce sufficiently randomness, the theorems are inconclusive.

Notice that such networks can also be analyzed using [Proposition 6.3](#). It is indeed not difficult to see that if there exists two paths that: 1/ start from the same entry of the input layer; 2/ end at the same entry of the output layer; 3/ if both paths are present (despite the action of ReLU) for every input/output pair; then [\(6.2\)](#) does not hold¹⁰ and $\mathbb{R}^{S \times K}$ is not identifiable. It is not clear at this point that the conditions 1, 2, 3 are met by all non-identifiable structured linear feed-forward networks. However, removing paths from the network (as is done by ReLU and Dropout) is a way to avoid conditions 1, 2, 3 to be met.

7. Stability guarantee. In this section, we consider errors of different natures. We assume that there exists \bar{L} and $L^* \in \mathbb{N}$, $\bar{\mathbf{h}} \in \mathcal{M}^{\bar{L}}$ and $\mathbf{h}^* \in \mathcal{M}^{L^*}$, such that

$$(7.1) \quad \|M_1(\bar{\mathbf{h}}_1) \cdots M_K(\bar{\mathbf{h}}_K) - X\| \leq \delta,$$

and

$$(7.2) \quad \|M_1(\mathbf{h}_1^*) \cdots M_K(\mathbf{h}_K^*) - X\| \leq \eta,$$

for δ and η typically small.

Again, this corresponds to existing unknown parameters $\bar{\mathbf{h}}$ that we estimate from a noisy observation X , using an inaccurate solution \mathbf{h}^* of [\(2.1\)](#) (as in [\[13\]](#) where the case $K = 1$ is studied). Otherwise, $\bar{\mathbf{h}}$ and \mathbf{h}^* shall be interpreted as different learned parameters; δ and η are the corresponding risks.

Notice that the above hypothesis does not even require [\(2.1\)](#) to have a solution. Algorithms which do not come with a guarantee sometimes manage to reach small δ and η values. In those cases, the analysis we conduct in this section permits to get the stability guarantee, despite the lack of a guarantee of the algorithm. Finally, the hypotheses [\(7.1\)](#) and [\(7.2\)](#) enable one to obtain guarantees for algorithms that, instead of minimizing [\(2.1\)](#), minimize an objective function which approximates the one in [\(2.1\)](#). This is particularly relevant for machine learning applications when [\(2.1\)](#) can be an empirical risk that needs to be regularized or is not truly minimized (for instance, when using *dropout* [\[75\]](#)).

A necessary and sufficient condition for the identifiability of \mathcal{M} is stated in [Proposition 6.3](#). The condition is on the way $\text{Ker}(\mathcal{A})$ and $P(\mathcal{M}^L) - P(\mathcal{M}^{L'})$ intersect. In order to get a stability guarantee, we need a stronger condition on the geometry of this intersection to hold for every L and $L' \in \mathbb{N}$. This condition is provided in the next definition.

DEFINITION 7.1. Deep-Null Space Property

Let $\gamma > 0$ and $\rho > 0$. We say that $\text{Ker}(\mathcal{A})$ satisfies the deep-Null Space Property (deep-NSP) with respect to the collection of models \mathcal{M} with constants (γ, ρ) if for any L and $L' \in \mathbb{N}$, any $T \in P(\mathcal{M}^L) - P(\mathcal{M}^{L'})$ satisfying $\|\mathcal{A}T\| \leq \rho$ and any $T' \in \text{Ker}(\mathcal{A})$, we have

$$(7.3) \quad \|T\| \leq \gamma \|T - T'\|.$$

The deep-NSP implies that, for $T \in P(\mathcal{M}^L) - P(\mathcal{M}^{L'})$ close to $\text{Ker}(\mathcal{A})$ in the sense that $\|\mathcal{A}T\| \leq \rho$ we must have, by decomposing $T = T' + T''$, with $T' \in \text{Ker}(\mathcal{A})$ and

¹⁰Simply consider two rank one tensors, each tensor being a Dirac at the position corresponding to one of the two paths.

T'' in its orthogonal complement

$$\|T\| \leq \gamma \|T - T'\| = \gamma \|T''\| \leq \frac{\gamma}{\sigma_{min}} \|\mathcal{A}T''\| \leq \frac{\gamma}{\sigma_{min}} \rho,$$

where σ_{min} is the smallest non-zero singular value of \mathcal{A} . In words, $\|T\|$ must be small. We can conclude that under the deep-NSP, $P(\mathcal{M}^L) - P(\mathcal{M}^{L'})$ and $\{T \in \mathbb{R}^{S^K} \mid \|\mathcal{A}T\| \leq \rho\}$ intersect at most in the neighborhood of 0.

Additionally, (7.3) implies that in the neighborhood of 0, $\text{Ker}(\mathcal{A})$ and $P(\mathcal{M}^L) - P(\mathcal{M}^{L'})$ are not tangential, i.e., their intersection is transverse.

If $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the collection of models \mathcal{M} with constants (γ, ρ) , then for all $T' \in \text{Ker}(\mathcal{A})$ and all $T \in P(\mathcal{M}^L) - P(\mathcal{M}^{L'})$ satisfying $\|\mathcal{A}T\| \leq \rho$,

$$\|T'\| \leq \|T\| + \|T' - T\| \leq (\gamma + 1)\|T' - T\|.$$

Therefore,

$$(7.4) \quad \forall T' \in \text{Ker}(\mathcal{A}), \quad \|T'\| \leq (\gamma + 1)d_{loc}(T', P(\mathcal{M}^L) - P(\mathcal{M}^{L'})),$$

where we have set for any $C \subset \mathbb{R}^{S^K}$

$$d_{loc}(T', C) = \inf_{T \in C, \|\mathcal{A}T\| \leq \rho} \|T' - T\|.$$

The converse is also true: if $\text{Ker}(\mathcal{A})$ satisfies (7.4), it satisfies the deep-NSP with respect to the collection of models \mathcal{M} with appropriate constants. In the context of the usual compressed sensing (i.e., when $K = 1$, \mathcal{M}^L contains L -sparse signals, \mathcal{A} is a rectangular matrix with full row rank and X is a vector), the localization appearing in d_{loc} can be discarded since the inequality must hold when T' is small and since in this case this localization has no effect. Therefore, in the compressed sensing context, (7.4) (and therefore deep-NSP) is the usual Null Space Property with respect to L -sparse vectors, as defined in [25]. However, deep-NSP is generalized to take into account deep structured linear network. This motivates the name.

In the general case, the deep-NSP can be understood as a local version of the generalized-NSP for \mathcal{A} relative to $P(\cup_{L \in \mathbb{N}} \mathcal{M}^L) - P(\cup_{L \in \mathbb{N}} \mathcal{M}^L)$, as defined in [13]. Our interest in locality (as imposed by the constraint $\|\mathcal{A}T\| \leq \rho$) is motivated by the fact that we want to use the deep-NSP when the signal to noise ratio is controlled (i.e., the hypotheses of Theorem 4.3 are satisfied). The condition for the stability property therefore includes such hypotheses.

We have not adapted the robust-NSP defined in [13]. The benefit in not using this definition is to obtain a simpler definition for deep-NSP. In particular (7.3) does not involve the geometry of \mathcal{A} in the orthogonal complement of $\text{Ker}(\mathcal{A})$. Looking in detail at the benefit of this adaptation is of great interest.

Finally, we trivially have the following facts:

- If $\text{Ker}(\mathcal{A}) = \{0\}$, then $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the model $\mathbb{R}^{S \times K}$ with constant $(1, +\infty)$.
- For any $\gamma' \geq \gamma$: If $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the collection of models \mathcal{M} with constants (γ, ρ) , then $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the collection of models \mathcal{M} with constants (γ', ρ) .
- For any $\mathcal{M}' \subset \mathcal{M}$: If $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the collection of models \mathcal{M} with constant (γ, ρ) , then $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the collection of models \mathcal{M}' with constants (γ, ρ) . In

particular, if $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the model $\mathbb{R}^{S \times K}$ with constant (γ, ρ) , it satisfies the deep-NSP with respect to any collection of models, with constants (γ, ρ) .

THEOREM 7.2. Sufficient condition for the stability property

Assume $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the collection of models \mathcal{M} and with the constants (γ, ρ) . For any \mathbf{h}^* as in (7.2) with η and δ (see (7.2) and (7.1)) such that $\delta + \eta \leq \rho$, we have

$$\|P(\mathbf{h}^*) - P(\bar{\mathbf{h}})\| \leq \frac{\gamma}{\sigma_{\min}} (\delta + \eta),$$

where σ_{\min} is the smallest non-zero singular value of \mathcal{A} . Moreover, if $\bar{\mathbf{h}} \in \mathbb{R}_*^{S \times K}$ and $\frac{\gamma}{\sigma_{\min}} (\delta + \eta) \leq \frac{1}{2} \max(\|P(\bar{\mathbf{h}})\|_{\infty}, \|P(\mathbf{h}^*)\|_{\infty})$ then

$$(7.5) \quad d_p(\langle \mathbf{h}^* \rangle, \langle \bar{\mathbf{h}} \rangle) \leq \frac{7(KS)^{\frac{1}{p}} \gamma}{\sigma_{\min}} \min\left(\|P(\bar{\mathbf{h}})\|_{\infty}^{\frac{1}{K}-1}, \|P(\mathbf{h}^*)\|_{\infty}^{\frac{1}{K}-1}\right) (\delta + \eta).$$

The first part of the proof is very similar to standard proofs in the Compressed Sensing and stable recovery literature. The second part simply uses [Theorem 4.3](#). The theorem is proved in [subsection 10.10](#).

[Theorem 7.2](#) provides a sufficient condition to obtain stability. The only significant hypothesis made on the deep structured linear network is that $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the collection of models \mathcal{M} . One might ask whether this hypothesis is sharp or not. The next theorem shows that the answer to this question is positive.

THEOREM 7.3. Necessary condition for the stability property

Assume the stability property holds: There exists C and $\delta > 0$ such that for any $\bar{L} \in \mathbb{N}$, $\bar{\mathbf{h}} \in \mathcal{M}^{\bar{L}}$, any $X = \mathcal{A}P(\bar{\mathbf{h}}) + e$, with $\|e\| \leq \delta$, any $L^* \in \mathbb{N}$ and any $\mathbf{h}^* \in \mathcal{M}^{L^*}$ such that

$$\|\mathcal{A}P(\mathbf{h}^*) - X\|^2 \leq \|e\|$$

we have

$$d_2(\langle \mathbf{h}^* \rangle, \langle \bar{\mathbf{h}} \rangle) \leq C \min\left(\|P(\bar{\mathbf{h}})\|_{\infty}^{\frac{1}{K}-1}, \|P(\mathbf{h}^*)\|_{\infty}^{\frac{1}{K}-1}\right) \|e\|.$$

Then, $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to the collection of models \mathcal{M} with constants

$$(\gamma, \rho) = (CS^{\frac{K-1}{2}} \sqrt{K} \sigma_{\max}, \delta)$$

where σ_{\max} is the spectral radius of \mathcal{A} .

The first part of the proof is inspired by and similar to the proof of the analogous converse statement in [25]. The second part simply uses [Theorem 4.5](#). The theorem is proved in [subsection 10.11](#).

The sharpness of the known results when $K = 2$ is usually argued by comparing the number of samples necessary for the recovery and the information theoretic limit of the problem. As far as the authors know, the above theorem is therefore new even when $K = 2$.

As is usually the case with Null Space Property or Restricted Isometry Property, it will often be difficult or impossible to establish that a particular operator \mathcal{A} satisfies the deep-NSP with respect to the collection of models \mathcal{M} . To find favorable cases, we need to consider random operators \mathcal{A} such that the distribution of \mathcal{A} enables one to establish that the deep-NSP holds with high probability, when in the

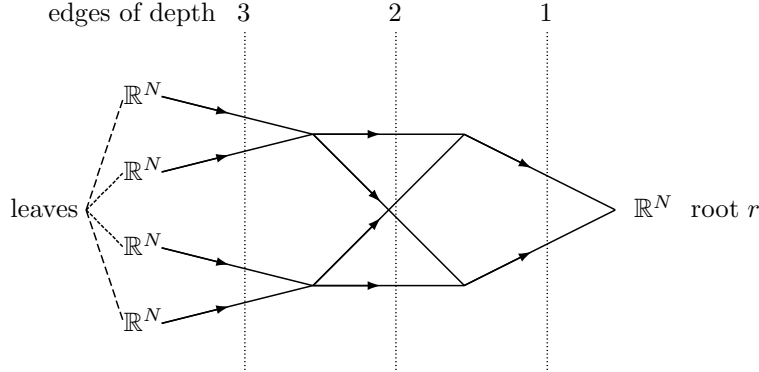


FIG. 1. Example of the considered convolutional linear network. To every edge is attached a convolution kernel. The network does not involve non-linearities or sampling.

right configurations (see the bibliography in section 2 whose references contain many examples of such arguments). The most common distribution for the analog of \mathcal{A} include operators/matrices whose coefficients are Gaussian or Bernoulli. Collections of models inducing sparsity, non-negativity, low-rank constraints *etc.* are the most studied. In this regard, the fact that there exists a low complexity test guaranteeing that the networks considered in section 8 can be stably recovered is an exception.

8. Application to convolutional linear network. We consider a convolutional linear network as depicted in Figure 1. The network typically aims at performing a linear analysis or synthesis of a signal living in \mathbb{R}^N . The considered convolutional linear network is defined from a rooted directed acyclic graph $\mathcal{G}(\mathcal{E}, \mathcal{N})$ composed of nodes \mathcal{N} and edges \mathcal{E} . Each edge connects two nodes. The root of the graph is denoted by r and the set containing all its leaves is denoted by \mathcal{F} . We denote by \mathcal{P} the set of all paths connecting the leaves and the root. We assume, without loss of generality, that the length of any path between any leaf and the root is independent of the considered leaf and equal to some constant $K \geq 0$. We also assume that, for any edge $e \in \mathcal{E}$, the number of edges separating e and the root is the same for all paths between e and r . It is called the depth of e . We also say that e belongs to the layer k . For any $k \in [K]$, we denote the set containing all the edges of depth k , by $\mathcal{E}(k)$.

Moreover, to any edge e is attached a convolution kernel of support $\mathcal{S}_e \subset [N]$. We assume (without loss of generality) that $\sum_{e \in \mathcal{E}(k)} |\mathcal{S}_e|$ is independent of k ($|\mathcal{S}_e|$ denotes the cardinality of \mathcal{S}_e). We take

$$S = \sum_{e \in \mathcal{E}(1)} |\mathcal{S}_e|.$$

For any edge e , we consider the mapping $\mathcal{T}_e : \mathbb{R}^S \rightarrow \mathbb{R}^N$ that maps any $h \in \mathbb{R}^S$ into the convolution kernel h_e , attached to the edge e , whose support is \mathcal{S}_e . It simply writes at the right location (i.e., those in \mathcal{S}_e) the entries of h defining the kernel on the edge e .

At each layer k , the convolutional linear network computes, for all $e \in \mathcal{E}(k)$, the convolution between the signal at the origin of e ; then, it attaches to any ending node the sum of all the convolutions arriving at that node. Examples of such convolutional linear networks includes wavelets, wavelet packets [61] or the fast transforms optimized

in [20, 21]. It is clear that the operation performed at any layer depends linearly on the parameters $h \in \mathbb{R}^S$ and that its results serve as inputs for the next layer. The convolutional linear network therefore depends on parameters $\mathbf{h} \in \mathbb{R}^{S \times K}$ and takes the form

$$X = M_1(\mathbf{h}_1) \cdots M_K(\mathbf{h}_K),$$

where the operators M_k satisfy (1.2).

This section aims at identifying conditions such that any unknown parameters $\bar{\mathbf{h}} \in \mathbb{R}^{S \times K}$ can be identified or stably recovered from $X = M_1(\bar{\mathbf{h}}_1) \cdots M_K(\bar{\mathbf{h}}_K)$ (possibly corrupted by an error).

In order to do so, we introduce some notation. We apply the convolutional linear network to an input $x \in \mathbb{R}^{N|\mathcal{F}|}$, where x is the concatenation of the signals $x^f \in \mathbb{R}^N$ for $f \in \mathcal{F}$. Therefore, X is the (horizontal) concatenation of $|\mathcal{F}|$ matrices $X^f \in \mathbb{R}^{N \times N}$ such that

$$(8.1) \quad Xx = \sum_{f \in \mathcal{F}} X^f x^f \quad , \text{ for all } x \in \mathbb{R}^{N|\mathcal{F}|}.$$

Consider the convolutional linear network defined by $\mathbf{h} \in \mathbb{R}^{S \times K}$ as well as $f \in \mathcal{F}$ and $n \in [N]$. The column of X corresponding to the entry n in the leaf f is the translation by n of

$$(8.2) \quad \sum_{p \in \mathcal{P}(f)} \mathcal{T}^p(\mathbf{h})$$

where $\mathcal{P}(f)$ contains all the paths of \mathcal{P} starting from the leaf f and

$$\mathcal{T}^p(\mathbf{h}) = \mathcal{T}_{e^1}(\mathbf{h}_1) * \cdots * \mathcal{T}_{e^K}(\mathbf{h}_K) \quad , \text{ with } p = (e^1, \dots, e^K),$$

is the composition of convolutions along the path p .

For any $k \in [K]$, define the mapping $\mathbf{e}_k : [S] \rightarrow \mathcal{E}(k)$ which provides for any $i \in [S]$ the unique edge of $\mathcal{E}(k)$ such that the i^{th} entry of $h \in \mathbb{R}^S$ contributes to $\mathcal{T}_{\mathbf{e}_k(i)}(h)$. For any $\mathbf{i} \in [S]^K$, let $\mathbf{p}_i = (\mathbf{e}_1(\mathbf{i}_1), \dots, \mathbf{e}_K(\mathbf{i}_K))$ and

$$\mathbf{I} = \{\mathbf{i} \in [S]^K \mid \mathbf{p}_i \in \mathcal{P}\}.$$

The latter contains all the indices corresponding to a valid path in the network. For any set of parameters $\mathbf{h} \in \mathbb{R}^{S \times K}$ and any path $\mathbf{p} \in \mathcal{P}$, we also let $\mathbf{h}^{\mathbf{p}}$ denote the restriction of \mathbf{h} to its indices contributing to the kernels on the path \mathbf{p} . We let $\mathbf{1} \in \mathbb{R}^S$ denote a vector of size S with all its entries equal to 1. For any edge e , $1^e \in \mathbb{R}^S$ consists of zeros except for the entries corresponding to the edge e which are equal to 1. For any $\mathbf{p} = (e^1, \dots, e^K) \in \mathcal{P}$, the support of $M_1(1^{e^1}) \cdots M_K(1^{e^K})$ is denoted by $\text{Supp}(\mathbf{p})$.

Finally, by Corollary 5.2 there exists a unique mapping

$$\mathcal{A} : \mathbb{R}^{S^K} \rightarrow \mathbb{R}^{N \times N|\mathcal{F}|}$$

such that

$$\mathcal{A}P(\mathbf{h}) = M_1(\mathbf{h}_1) \cdots M_K(\mathbf{h}_K) \quad , \text{ for all } \mathbf{h} \in \mathbb{R}^{S \times K},$$

where P is the Segre embedding defined in (4.1).

DEFINITION 8.1. *We say the topology of the network is sufficiently scattered if and only if all the entries of $M_1(1) \cdots M_K(1)$ belong to $\{0, 1\}$.*

The following statements will show that having a sufficiently scattered topology is a necessary and sufficient condition for the stability of the optimal parameters. Before going into this, we illustrate the scattering property with a simple example.

Consider a simple composition of two convolutions ($K = 2$ and $|\mathcal{P}| = 1$). At first, we make an assumption on the supports \mathcal{S}_e , imposing that the supports of both kernels are in $\{1, 2, 3\}$. The topology is obviously not sufficiently scattered. Indeed, some of the entries of the convolution kernel corresponding to the matrix $M_1(1)M_2(1)$ are equal to 2.

Now consider an assumption on the network topology imposing $\{1, 2, 3\}$ for the support of the first kernel and $\{1, 10\}$ for the second. When we observe the convolution of two kernels having such supports, we see two replicas of the first kernel; the amplitudes of the replicas depend on the second kernel; and both kernels are identifiable. In this last example, the network topology is sufficiently scattered.

The scattering condition can easily be computed using [Algorithm 8.1](#). Indeed, when applying the network to a Dirac in the leaf f , using (8.1), we obtain the convolution kernel of X_f . We can then easily test if X_f only contains 0's and 1's. The numerical complexity of [Algorithm 8.1](#) is essentially the cost for applying $|\mathcal{F}|$ times the network. It is usually low. Notice that a network is sufficiently scattered if and only if, for all leaves $f \in \mathcal{F}$, the sub-networks originating at f are sufficiently scattered. The scattering of these subnetworks are independent. The fact the the convolution kernels, for the different leaves, overlap does not affect the scattering property.

Algorithm 8.1 Algorithm testing if the topology of the convolutional network leads to the stability guarantee.

Input: The network topology

Output: Boolean output = ‘true’, if the topology is sufficiently scattered; ‘false’, otherwise.

output = true

For each $f \in \mathcal{F}$ **do**

 Build x : a Dirac positioned at the leaf f

 Apply the network to x in order to compute $y = M_1(1) \cdots M_K(1)x$

 If some of the entries of y are outside $\{0, 1\}$ then set output = false

end For each

Finally, beside the known examples in blind-deconvolution (i.e., when $K = 2$ and $|\mathcal{P}| = 1$) [2, 9], there are (truly deep) convolutional linear networks satisfying the condition of the first statement of [Proposition 8.2](#). For instance, the convolutional linear network corresponding to the un-decimated Haar (wavelet)¹¹ transform is a tree and for any of its leaves $f \in \mathcal{F}$, $|\mathcal{P}(f)| = 1$. Moreover, the support of the kernel living on the edge e , of depth k , on this path is $\{0, 2^k\}$. It is therefore not difficult to check that the first condition of [Proposition 8.2](#) holds.

PROPOSITION 8.2. Necessary condition of identifiability of convolutional linear network

- *Either the topology of the network is sufficiently scattered and then*
 1. *for any distinct \mathbf{p} and $\mathbf{p}' \in \mathcal{P}$, $\text{Supp}(\mathbf{p}) \cap \text{Supp}(\mathbf{p}') = \emptyset$.*
 2. *$\text{Ker}(\mathcal{A}) = \{T \in \mathbb{R}^{S^K} \mid \forall \mathbf{i} \in \mathbf{I}, T_{\mathbf{i}} = 0\}$.*

¹¹Un-decimated means computed with the “Algorithme à trous”, [61], Section 5.5.2 and 6.3.2. The Haar wavelet is described in [61], Section 7.2.2, p. 247 and Example 7.7, p. 235.

- or the topology of the network is not sufficiently scattered and then $\mathbb{R}^{S \times K}$ is not identifiable.

The Proposition is proved in [subsection 10.12](#).

PROPOSITION 8.3. *If $|\mathcal{P}| = 1$ and the topology of the network is sufficiently scattered, then $\text{Ker}(\mathcal{A}) = \{0\}$ and $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with respect to any model collection \mathcal{M} with constant $(\gamma, \rho) = (1, +\infty)$. Moreover, we have $\sigma_{\min} = \sqrt{N}$.*

The Proposition is proved in [subsection 10.13](#).

In what follows, we establish stability results for a convolutional linear network estimator. In order to do so, we consider a convolutional linear network of known structure $\mathcal{G}(\mathcal{E}, \mathcal{N})$ and $(\mathcal{S}_e)_{e \in \mathcal{E}}$. We consider parameters $\bar{\mathbf{h}} \in \mathbb{R}^{S \times K}$ and $\mathbf{h}^* \in \mathbb{R}^{S \times K}$ such that

$$(8.3) \quad \|M_1(\bar{\mathbf{h}}_1) \cdots M_K(\bar{\mathbf{h}}_K) - X\| \leq \delta,$$

and

$$(8.4) \quad \|M_1(\mathbf{h}_1^*) \cdots M_K(\mathbf{h}_K^*) - X\| \leq \eta.$$

We say that two networks sharing the same structure and defined by \mathbf{h} and $\mathbf{g} \in \mathbb{R}^{S \times K}$ are equivalent if and only if

$$\forall \mathbf{p} \in \mathcal{P}, \exists (\lambda_e)_{e \in \mathbf{p}} \in \mathbb{R}^{\mathbf{p}}, \text{ such that } \prod_{e \in \mathbf{p}} \lambda_e = 1 \text{ and } \forall e \in \mathbf{p}, \mathcal{T}_e(\mathbf{g}) = \lambda_e \mathcal{T}_e(\mathbf{h}).$$

The equivalence class of $\mathbf{h} \in \mathbb{R}^{S \times K}$ is denoted by $\{\mathbf{h}\}$. For any $p \in [1, +\infty]$, we define

$$\delta_p(\{\mathbf{h}\}, \{\mathbf{g}\}) = \left(\sum_{\mathbf{p} \in \mathcal{P}} d_p(\langle \mathbf{h}^{\mathbf{p}} \rangle, \langle \mathbf{g}^{\mathbf{p}} \rangle)^p \right)^{\frac{1}{p}},$$

where recall that $\mathbf{h}^{\mathbf{p}}$ (resp $\mathbf{g}^{\mathbf{p}}$) denotes the restriction of \mathbf{h} (resp. \mathbf{g}) to the path \mathbf{p} and d_p is defined in [Definition 4.1](#). Since d_p is a metric, it follows that δ_p is a metric between network classes.

We summarize the results concerning convolutional networks in the following theorem.

THEOREM 8.4. Necessary and sufficient condition of stable recovery of convolutional linear network

If [Algorithm 8.1](#) returns ‘false’, the network topology is not sufficiently scattered and the network is not identifiable.

If [Algorithm 8.1](#) returns ‘true’, if $\bar{\mathbf{h}}$ and \mathbf{h}^ satisfy (8.3) and (8.4) and*

- *if all the edges support a significant convolution kernel: there exists $\varepsilon > 0$ such that for all $e \in \mathcal{E}$, $\|\mathcal{T}_e(\bar{\mathbf{h}})\|_{\infty} \geq \varepsilon$,*
- *if the “signal to noise ratio” is sufficient: $\delta + \eta \leq \frac{\sqrt{N}\varepsilon^K}{2}$,*

then the network defined by \mathbf{h}^ and $\bar{\mathbf{h}}$ are close to each other*

$$\delta_p(\{\mathbf{h}^*\}, \{\bar{\mathbf{h}}\}) \leq 7(KS')^{\frac{1}{p}} \varepsilon^{1-K} \frac{\delta + \eta}{\sqrt{N}},$$

where $S' = \max_{e \in \mathcal{E}} |\mathcal{S}_e|$ is the size the largest convolution kernel.

The theorem is proved in [subsection 10.14](#).

9. Conclusion and perspectives. In this paper, we have established necessary and sufficient conditions for the identifiability and stable recovery of deep structured linear networks. They rely on the lifting of the problem in a tensor space. The technique is called *tensorial lifting*. The main results are proved using compressed sensing technics and properties of the Segre embedding (the embedding that maps the parameters in the tensor space). The general results are then specialized to establish necessary and sufficient conditions for the stable recovery of a convolutional linear network of any depth $K \geq 1$.

Among the most salient perspectives, we mention the possibility to study deep feed-forward ReLU networks. For such a network, the action of ReLU is different for every sample; this leads to a different operator \mathcal{A} for every sample; and all the different \mathcal{A} 's sense (linearly) the same rank one tensor. We can concatenate these operators to form a unique sensing operator. For instance, when modeling the action of ReLU as a Bernouilli variable applied to every path of the network, we expect to obtain sample complexity bounds (for instance) under the favorable hypothesis that an oracle has given us the action of ReLU.

A natural perspective of this work is also to study compressed networks (see [7]), when the compression preserves the expressivity of the network.

Finally, the model considered in this paper approximately solves polynomial equations: $\mathcal{A}P(\mathbf{h}) \sim X$. The structure of the polynomials is induced by the operators M_k (i.e. the network topology) is very particular and restrictive. For instance, we only consider homogeneous polynomials in $P(\mathbf{h})$. Extending this work to larger families of polynomials as well as limits of polynomials seem natural.

10. Appendices.

10.1. Proof of Proposition 4.2. Notice that, the sets $\langle \mathbf{h} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}$ and $\langle \mathbf{g} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}$ are finite and therefore the infimum in the definition of d is reached. We also have whatever $\mathbf{h}, \mathbf{g} \in \mathbb{R}_*^{S \times K}$

$$(10.1) \quad d_p(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) = \inf_{\mathbf{h}' \in \langle \mathbf{h} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}} \left(\inf_{\mathbf{g}' \in \langle \mathbf{g} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}} \|\mathbf{h}' - \mathbf{g}'\|_p \right).$$

Moreover, whatever $\mathbf{h} \in \mathbb{R}_*^{S \times K}$ and \mathbf{h}' and $\mathbf{h}'' \in \langle \mathbf{h} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}$ there exist $(s_k)_{k \in [K]} \in \{-1, 1\}^K$ such that $\prod_{k \in [K]} s_k = 1$ and

$$\mathbf{h}'_k = s_k \mathbf{h}''_k, \quad \forall k \in [K].$$

Using the above two properties, we can check that

$$\inf_{\mathbf{g}' \in \langle \mathbf{g} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}} \|\mathbf{h}' - \mathbf{g}'\|_p = \inf_{\mathbf{g}' \in \langle \mathbf{g} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}} \|\mathbf{h}'' - \mathbf{g}'\|_p$$

As a consequence, the outer infimum in (10.1) is irrelevant and we have

$$d_p(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) = \inf_{\mathbf{g}' \in \langle \mathbf{g} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}} \|\mathbf{h}' - \mathbf{g}'\|_p, \quad \forall \mathbf{h}, \mathbf{g} \in \mathbb{R}_*^{S \times K} \text{ and } \mathbf{h}' \in \langle \mathbf{h} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}.$$

Using this last property, we easily check that d_p is a metric on $\mathbb{R}_*^{S \times K} / \sim$. \square

10.2. Proof of Theorem 4.3. Notice first that when $K = 1$ the inequality is a straightforward consequence of the usual inequalities between l^p norms. We therefore assume from now on that $K \geq 2$.

All along the proof, we consider \mathbf{h} and $\mathbf{g} \in \mathbb{R}_*^{S \times K}$ and assume that $\|P(\mathbf{h})\|_\infty \geq \|P(\mathbf{g})\|_\infty$. We also assume that $\|P(\mathbf{g}) - P(\mathbf{h})\|_\infty \leq \frac{1}{2}\|P(\mathbf{h})\|_\infty$. We first prove the inequality when $p = q = +\infty$.

In order to do so, we consider

$$\mathbf{i} \in \operatorname{argmax}_{\mathbf{j} \in [S]^K} |P(\mathbf{h})_{\mathbf{j}}|$$

and assume, without loss of generality (otherwise, we can multiply one vector of \mathbf{h} and \mathbf{g} by -1 to get this property and multiply back once the inequality has been established), that $P(\mathbf{h})_{\mathbf{i}} \geq 0$. We therefore have $P(\mathbf{h})_{\mathbf{i}} = \|P(\mathbf{h})\|_\infty$. Notice also that we have, under the above hypotheses,

$$(10.2) \quad \|P(\mathbf{g})\|_\infty \geq P(\mathbf{g})_{\mathbf{i}} \geq P(\mathbf{h})_{\mathbf{i}} - \|P(\mathbf{g}) - P(\mathbf{h})\|_\infty \geq \frac{1}{2}\|P(\mathbf{h})\|_\infty > 0.$$

Moreover, we consider the operator $E_{\mathbf{i}}$ that extracts the K signals of size S that are obtained when freezing, at the index \mathbf{i} in a tensor T , all coordinates but one. Formally, we denote

$$\begin{aligned} E_{\mathbf{i}} : \mathbb{R}^{S^K} &\longrightarrow \mathbb{R}^{S \times K} \\ T &\longmapsto E_{\mathbf{i}}(T) \end{aligned}$$

where for all $k \in [K]$ and all $j \in [S]$

$$E_{\mathbf{i}}(T)_{k,j} = T_{\mathbf{i}_1, \dots, \mathbf{i}_{k-1}, j, \mathbf{i}_{k+1}, \dots, \mathbf{i}_K}.$$

We consider

$$\mathbf{h}' = (P(\mathbf{h})_{\mathbf{i}})^{-1 + \frac{1}{K}} E_{\mathbf{i}}(P(\mathbf{h})) \quad \text{and} \quad \mathbf{g}' = (P(\mathbf{g})_{\mathbf{i}})^{-1 + \frac{1}{K}} E_{\mathbf{i}}(P(\mathbf{g})).$$

We have for all $\mathbf{j} \in [S]^K$

$$\begin{aligned} P(\mathbf{h}')_{\mathbf{j}} &= (P(\mathbf{h})_{\mathbf{i}})^{-K+1} P(E_{\mathbf{i}}(P(\mathbf{h})))_{\mathbf{j}}, \\ &= (P(\mathbf{h})_{\mathbf{i}})^{-K+1} \prod_{k=1}^K P(\mathbf{h})_{\mathbf{i}_1, \dots, \mathbf{i}_{k-1}, \mathbf{j}_k, \mathbf{i}_{k+1}, \dots, \mathbf{i}_K} \\ &= (P(\mathbf{h})_{\mathbf{i}})^{-K+1} \prod_{k=1}^K \mathbf{h}_{1, \mathbf{i}_1} \dots \mathbf{h}_{k-1, \mathbf{i}_{k-1}} \mathbf{h}_{k, \mathbf{j}_k} \mathbf{h}_{k+1, \mathbf{i}_{k+1}} \dots \mathbf{h}_{K, \mathbf{i}_K} \\ &= \prod_{k=1}^K \mathbf{h}_{k, \mathbf{j}_k} = P(\mathbf{h})_{\mathbf{j}}. \end{aligned}$$

We therefore have $P(\mathbf{h}') = P(\mathbf{h})$. This can be written $\mathbf{h}' \in \langle \mathbf{h} \rangle$. Similarly, we have $\mathbf{g}' \in \langle \mathbf{g} \rangle$.

Also, because of the definition of \mathbf{i} and \mathbf{h}' , we are guaranteed that, whatever $k \in [K]$,

$$\begin{aligned} \|\mathbf{h}'_k\|_\infty &= (P(\mathbf{h})_{\mathbf{i}})^{-1 + \frac{1}{K}} \|E_{\mathbf{i}}(P(\mathbf{h}))\|_\infty \\ &= \|P(\mathbf{h})\|_\infty^{-1 + \frac{1}{K}} \|P(\mathbf{h})\|_\infty = \|P(\mathbf{h})\|_\infty^{\frac{1}{K}} \end{aligned}$$

The latter being independent of k , we have $\mathbf{h}' \in \mathbb{R}_{\text{diag}}^{S \times K}$. Unfortunately, unless for instance $\mathbf{i} \in \operatorname{argmax}_{\mathbf{j} \in [S]^K} |P(\mathbf{g})_{\mathbf{j}}|$, it might occur that $\mathbf{g}' \notin \mathbb{R}_{\text{diag}}^{S \times K}$. However, if we consider

$$\mathbf{g}'' \in \operatorname{argmin}_{\mathbf{f} \in \langle \mathbf{g} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}} \|\mathbf{f} - \mathbf{g}'\|_{\infty},$$

we have since $\mathbf{h}' \in \langle \mathbf{h} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}$ and $\mathbf{g}'' \in \langle \mathbf{g} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}$

$$d_{\infty}(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) \leq \|\mathbf{h}' - \mathbf{g}''\|_{\infty}$$

and therefore

$$(10.3) \quad d_{\infty}(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) \leq \|\mathbf{h}' - \mathbf{g}'\|_{\infty} + \|\mathbf{g}' - \mathbf{g}''\|_{\infty}.$$

In the sequel we will successively calculate upper bounds of $\|\mathbf{h}' - \mathbf{g}'\|_{\infty}$ and $\|\mathbf{g}' - \mathbf{g}''\|_{\infty}$ in order to find an upper bound of $d_{\infty}(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle)$.

Upper bound of $\|\mathbf{h}' - \mathbf{g}'\|_{\infty}$:

We have

$$\begin{aligned} \|\mathbf{h}' - \mathbf{g}'\|_{\infty} &= \|(P(\mathbf{h})_{\mathbf{i}})^{-1+\frac{1}{K}} E_{\mathbf{i}}(P(\mathbf{h})) - (P(\mathbf{g})_{\mathbf{i}})^{-1+\frac{1}{K}} E_{\mathbf{i}}(P(\mathbf{g}))\|_{\infty} \\ &\leq \|(P(\mathbf{h})_{\mathbf{i}})^{-1+\frac{1}{K}} (E_{\mathbf{i}}(P(\mathbf{h})) - E_{\mathbf{i}}(P(\mathbf{g})))\|_{\infty} \\ &\quad + \left\| \left((P(\mathbf{h})_{\mathbf{i}})^{-1+\frac{1}{K}} - (P(\mathbf{g})_{\mathbf{i}})^{-1+\frac{1}{K}} \right) E_{\mathbf{i}}(P(\mathbf{g})) \right\|_{\infty} \\ &\leq \|P(\mathbf{h})\|_{\infty}^{-1+\frac{1}{K}} \|E_{\mathbf{i}}(P(\mathbf{h})) - E_{\mathbf{i}}(P(\mathbf{g}))\|_{\infty} + \|P(\mathbf{g})\|_{\infty} |(P(\mathbf{h})_{\mathbf{i}})^{-1+\frac{1}{K}} - (P(\mathbf{g})_{\mathbf{i}})^{-1+\frac{1}{K}}| \\ &\leq \|P(\mathbf{h})\|_{\infty}^{-1+\frac{1}{K}} \|P(\mathbf{h}) - P(\mathbf{g})\|_{\infty} + \|P(\mathbf{h})\|_{\infty} |(P(\mathbf{h})_{\mathbf{i}})^{-1+\frac{1}{K}} - (P(\mathbf{g})_{\mathbf{i}})^{-1+\frac{1}{K}}| \end{aligned}$$

But we also have using the mean value theorem and (10.2)

$$\begin{aligned} |(P(\mathbf{h})_{\mathbf{i}})^{-1+\frac{1}{K}} - (P(\mathbf{g})_{\mathbf{i}})^{-1+\frac{1}{K}}| &\leq \left(1 - \frac{1}{K}\right) P(\mathbf{g})_{\mathbf{i}}^{-2+\frac{1}{K}} |P(\mathbf{h})_{\mathbf{i}} - P(\mathbf{g})_{\mathbf{i}}| \\ &\leq \left(1 - \frac{1}{K}\right) \left(\frac{1}{2} \|P(\mathbf{h})\|_{\infty}\right)^{-2+\frac{1}{K}} \|P(\mathbf{h}) - P(\mathbf{g})\|_{\infty} \\ &\leq 4 \|P(\mathbf{h})\|_{\infty}^{-2+\frac{1}{K}} \|P(\mathbf{h}) - P(\mathbf{g})\|_{\infty} \end{aligned}$$

We therefore finally obtain that

$$(10.4) \quad \|\mathbf{h}' - \mathbf{g}'\|_{\infty} \leq 5 \|P(\mathbf{h})\|_{\infty}^{-1+\frac{1}{K}} \|P(\mathbf{h}) - P(\mathbf{g})\|_{\infty}.$$

Upper bound of $\|\mathbf{g}' - \mathbf{g}''\|_{\infty}$:

First, since $\mathbf{g}'' \in \langle \mathbf{g} \rangle = \langle \mathbf{g}' \rangle$, we know that there exists $(\lambda_k)_{k \in [K]} \in \mathbb{R}^K$ such that

$$(10.5) \quad \prod_{k=1}^K \lambda_k = 1$$

and

$$\mathbf{g}''_k = \lambda_k \mathbf{g}'_k, \quad \text{for all } k \in [K].$$

Furthermore, we have for all $k \in [K]$

$$(10.6) \quad \|\mathbf{g}'_k - \mathbf{g}''_k\|_{\infty} = |1 - \lambda_k| \|\mathbf{g}'_k\|_{\infty}.$$

Also, if there is k' such that $\lambda_{k'} < 0$, since (10.5) holds, there necessarily exist another k'' such that $\lambda_{k''} < 0$. If we replace $\mathbf{g}_{k'}$ by $-\mathbf{g}_{k'}$ and replace $\mathbf{g}_{k''}$ by $-\mathbf{g}_{k''}$ we remain in $\langle \mathbf{g} \rangle \cap \mathbb{R}_{\text{diag}}^{S \times K}$ and can only make $\|\mathbf{g}' - \mathbf{g}''\|_\infty$ decrease. Repeating this process until all the λ_k 's are non-negative, we can assume without loss of generality that

$$\lambda_k \geq 0 \quad , \text{ whatever } k \in [K].$$

This being said, we establish two other simple facts that motivate the structure of the proof. First, in order to find an upper bound for (10.6), we easily establish (using (10.2)) that

$$\begin{aligned} \|\mathbf{g}'_k\|_\infty &= (P(\mathbf{g})_i)^{-1+\frac{1}{K}} \|E_i(P(\mathbf{g}))\|_\infty \\ &\leq \left(\frac{1}{2}\|P(\mathbf{h})\|_\infty\right)^{-1+\frac{1}{K}} \|P(\mathbf{h})\|_\infty \end{aligned}$$

and therefore

$$(10.7) \quad \|\mathbf{g}'_k\|_\infty \leq 2\|P(\mathbf{h})\|_\infty^{\frac{1}{K}}.$$

Second, the value λ_k appearing in (10.6), can be bounded by using bounds on $\|\mathbf{g}'_k\|_\infty$ and the identity

$$(10.8) \quad \|\mathbf{g}''_k\|_\infty = \|P(\mathbf{g})\|_\infty^{\frac{1}{K}} = \lambda_k \|\mathbf{g}'_k\|_\infty.$$

Qualitatively, the latter identity indeed guarantees that, as $\|P(\mathbf{g}) - P(\mathbf{h})\|_\infty$ goes to 0, λ_k goes to 1. Let us now establish this quantitatively.

Recalling that

$$\mathbf{g}' = (P(\mathbf{g})_i)^{-1+\frac{1}{K}} E_i(P(\mathbf{g})),$$

and using (10.2) again, we obtain

$$\|\mathbf{g}'_k\|_\infty \leq \left(\|P(\mathbf{h})\|_\infty - \frac{1}{2}\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty\right)^{-1+\frac{1}{K}} \|P(\mathbf{g})\|_\infty.$$

We also have (again, using (10.2))

$$\begin{aligned} \|\mathbf{g}'_k\|_\infty &\geq (P(\mathbf{g})_i)^{-1+\frac{1}{K}} |P(\mathbf{g})_i| \\ &= (P(\mathbf{g})_i)^{\frac{1}{K}} \\ &\geq \left(\|P(\mathbf{h})\|_\infty - \frac{1}{2}\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty\right)^{\frac{1}{K}}. \end{aligned}$$

Plugging the upper bound of $\|\mathbf{g}'_k\|_\infty$ in (10.8), using successively (10.2), the mean

value theorem and the hypothesis on the size of $P(\mathbf{h}) - P(\mathbf{g})$ gives:

$$\begin{aligned}
\lambda_k - 1 &= \frac{\|P(\mathbf{g})\|_\infty^{\frac{1}{K}}}{\|\mathbf{g}'_k\|_\infty} - 1 \\
&\geq \|P(\mathbf{g})\|_\infty^{-1+\frac{1}{K}} \left(\|P(\mathbf{h})\|_\infty - \frac{1}{2}\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty \right)^{1-\frac{1}{K}} - 1 \\
&\geq \left(1 - \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{2\|P(\mathbf{h})\|_\infty} \right)^{1-\frac{1}{K}} - 1 \\
&\geq -\left(1 - \frac{1}{K}\right) \left(1 - \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{2\|P(\mathbf{h})\|_\infty} \right)^{-\frac{1}{K}} \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{2\|P(\mathbf{h})\|_\infty} \\
&\geq -\left(1 - \frac{1}{4}\right)^{-\frac{1}{K}} \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{2\|P(\mathbf{h})\|_\infty} \\
&\geq -\frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{\|P(\mathbf{h})\|_\infty}.
\end{aligned}$$

Similarly, plugging the lower bound of $\|\mathbf{g}'_k\|_\infty$ in (10.8), we obtain using successively (10.2), the mean value theorem and the hypothesis on the size of $P(\mathbf{h}) - P(\mathbf{g})$:

$$\begin{aligned}
\lambda_k - 1 &\leq \|P(\mathbf{g})\|_\infty^{\frac{1}{K}} \left(\|P(\mathbf{h})\|_\infty - \frac{1}{2}\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty \right)^{-\frac{1}{K}} - 1 \\
&\leq \left(1 - \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{2\|P(\mathbf{h})\|_\infty} \right)^{-\frac{1}{K}} - 1 \\
&\leq \frac{1}{K} \left(1 - \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{2\|P(\mathbf{h})\|_\infty} \right)^{-1-\frac{1}{K}} \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{2\|P(\mathbf{h})\|_\infty} \\
&\leq \frac{1}{K} \left(1 - \frac{1}{4} \right)^{-1-\frac{1}{K}} \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{2\|P(\mathbf{h})\|_\infty} \\
&\leq \frac{4^2}{2K3^2} \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{\|P(\mathbf{h})\|_\infty} \\
&\leq \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{\|P(\mathbf{h})\|_\infty}.
\end{aligned}$$

Finally, we get

$$(10.9) \quad |\lambda_k - 1| \leq \frac{\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty}{\|P(\mathbf{h})\|_\infty}.$$

By combining (10.6), (10.7) and (10.9), we obtain

$$\|\mathbf{g}'_k - \mathbf{g}''_k\|_\infty \leq 2 \|P(\mathbf{h})\|_\infty^{-1+\frac{1}{K}} \|P(\mathbf{h}) - P(\mathbf{g})\|_\infty.$$

Combining the latter inequality with (10.3) and (10.4) provides

$$d_\infty(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) \leq 7 \|P(\mathbf{h})\|_\infty^{-1+\frac{1}{K}} \|P(\mathbf{h}) - P(\mathbf{g})\|_\infty,$$

and concludes the proof when $p = q = +\infty$.

In order to establish the property when $1 \leq p \leq +\infty$ and $1 \leq q \leq +\infty$, we simply use the fact that

$$d_p(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle) \leq (KS)^{\frac{1}{p}} d_\infty(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle)$$

and

$$\|P(\mathbf{h}) - P(\mathbf{g})\|_\infty \leq \|P(\mathbf{h}) - P(\mathbf{g})\|_q.$$

□

10.3. Proof of Proposition 4.4. In the example, we consider \mathbf{h} and \mathbf{g} such that for all $k \in [K]$ and all $i \in [S]$

$$\mathbf{h}_{k,i} = \begin{cases} 1 & \text{if } i = 0, \\ 0 & \text{otherwise,} \end{cases} \quad \text{and} \quad \mathbf{g}_{k,i} = \begin{cases} \left(\frac{1}{2}\right)^{\frac{1}{K}} & \text{if } i = 0, \\ \epsilon_q & \text{otherwise,} \end{cases}$$

where $\epsilon_{+\infty} = \left(\frac{1}{2}\right)^{\frac{1}{K}}$ and $\epsilon_q = \min\left(\left(\frac{1 - \left(\frac{1}{2}\right)^{\frac{q}{K}}}{S-1}\right)^{\frac{1}{q}}, \left(\frac{1}{2}\right)^{\frac{1}{K}}\right)$, if $q < +\infty$. We immediately obtain

$$\|P(\mathbf{h})\|_\infty = 1, \quad \|P(\mathbf{g})\|_\infty = \frac{1}{2} \quad \text{and} \quad \|P(\mathbf{h}) - P(\mathbf{g})\|_\infty = \frac{1}{2}.$$

We also have,

$$d_p(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle)^p = \|\mathbf{h} - \mathbf{g}\|_p^p \geq K(S-1) \epsilon_q^p \geq \frac{KS}{2} \epsilon_q^p.$$

Decomposing the sum necessary to the calculation of the l^q norm of a tensor according to number of index different from 0 (which corresponds to l in the sum below), we obtain

$$\begin{aligned} \|P(\mathbf{h}) - P(\mathbf{g})\|_q^q &= \sum_{l=0}^K \binom{l}{K} (S-1)^l \epsilon_q^{lq} \left(\frac{1}{2}\right)^{\frac{(K-l)q}{K}}, \\ &= \left(\left(\frac{1}{2}\right)^{\frac{q}{K}} + (S-1)\epsilon_q^q \right)^K \leq 1. \end{aligned}$$

We then easily obtain that

$$7 \|P(\mathbf{h})\|_\infty^{-1 + \frac{1}{K}} (KS)^{\frac{1}{p}} \|P(\mathbf{h}) - P(\mathbf{g})\|_q \leq 7(KS)^{\frac{1}{p}},$$

and therefore

$$(10.10) \quad 7 \|P(\mathbf{h})\|_\infty^{-1 + \frac{1}{K}} (KS)^{\frac{1}{p}} \|P(\mathbf{h}) - P(\mathbf{g})\|_q \leq 7 \frac{d_p(\langle \mathbf{h} \rangle, \langle \mathbf{g} \rangle)}{\epsilon_q} 2^{\frac{1}{p}}.$$

We first calculate a lower bound of ϵ_q when $\epsilon_q = \left(\frac{1 - \left(\frac{1}{2}\right)^{\frac{q}{K}}}{S-1}\right)^{\frac{1}{q}}$ (which, in particular, rules out $q = +\infty$). Using the mean value theorem, we obtain

$$1 - \left(\frac{1}{2}\right)^{\frac{q}{K}} \geq \min_{t \in [\frac{1}{2}, 1]} \left(\frac{q}{K} t^{\frac{q}{K}-1}\right) \left(1 - \frac{1}{2}\right).$$

Distinguishing, whether $q \leq K$ or not, we find after a short calculation that, since $q \geq 1$,

$$1 - \left(\frac{1}{2}\right)^{\frac{q}{K}} \geq \min\left(\frac{1}{2K}, \frac{1}{K} \left(\frac{1}{2}\right)^{\frac{q}{K}}\right) = \frac{1}{K} \min\left(\left(\frac{1}{2}\right)^{\frac{1}{q}}, \left(\frac{1}{2}\right)^{\frac{1}{K}}\right)^q \geq \frac{1}{K2^q}.$$

We then deduce

$$\epsilon_q \geq \frac{1}{2(KS)^{\frac{1}{q}}}.$$

Of course, when $\epsilon_q = \left(\frac{1}{2}\right)^{\frac{1}{K}}$ (which includes $q = +\infty$), we immediately obtain

$$\epsilon_q \geq \frac{1}{2}.$$

Using this lower bound in (10.10) leads to the bounds stated in the proposition. \square

10.4. Proof of Theorem 4.5. Before starting the proof, we define for any $k \in \{0, \dots, K\}$

$$P_k(\mathbf{h}, \mathbf{g})_{\mathbf{i}} = \mathbf{g}_{1, \mathbf{i}_1} \dots \mathbf{g}_{k, \mathbf{i}_k} \mathbf{h}_{k+1, \mathbf{i}_{k+1}} \dots \mathbf{h}_{K, \mathbf{i}_K} \quad , \text{ for all } \mathbf{h}, \mathbf{g} \in \mathbb{R}^{S \times K} \text{ and all } \mathbf{i} \in [S]^K. \blacksquare$$

We consider \mathbf{g} and $\mathbf{h} \in \mathbb{R}^{S \times K}$. Let us first assume that $\|\mathbf{g}\|_{\infty} \leq \|\mathbf{h}\|_{\infty} = 1$. We have for any $\mathbf{i} \in [S]^K$, using this hypothesis and standard inequalities between l^p norms, when $q < +\infty$

$$\begin{aligned} |P(\mathbf{g})_{\mathbf{i}} - P(\mathbf{h})_{\mathbf{i}}|^q &= \left| \sum_{k=0}^{K-1} (P_{k+1}(\mathbf{h}, \mathbf{g})_{\mathbf{i}} - P_k(\mathbf{h}, \mathbf{g})_{\mathbf{i}}) \right|^q \\ &\leq K^{q-1} \sum_{k=0}^{K-1} |P_{k+1}(\mathbf{h}, \mathbf{g})_{\mathbf{i}} - P_k(\mathbf{h}, \mathbf{g})_{\mathbf{i}}|^q \\ &\leq K^{q-1} \sum_{k=0}^{K-1} |\mathbf{g}_{k+1, \mathbf{i}_{k+1}} - \mathbf{h}_{k+1, \mathbf{i}_{k+1}}|^q \end{aligned}$$

The same calculation when $q = +\infty$ leads to

$$|P(\mathbf{g})_{\mathbf{i}} - P(\mathbf{h})_{\mathbf{i}}| \leq K \max_{k=1..K} |\mathbf{g}_{k, \mathbf{i}_k} - \mathbf{h}_{k, \mathbf{i}_k}|.$$

Therefore, we have when $q < +\infty$

$$\begin{aligned} \|P(\mathbf{h}) - P(\mathbf{g})\|_q^q &= \sum_{\mathbf{i} \in [S]^K} |P(\mathbf{h})_{\mathbf{i}} - P(\mathbf{g})_{\mathbf{i}}|^q \\ &\leq K^{q-1} \sum_{k=1}^K \sum_{\mathbf{i} \in [S]^K} |\mathbf{g}_{k, \mathbf{i}_k} - \mathbf{h}_{k, \mathbf{i}_k}|^q \\ &= K^{q-1} \sum_{k=1}^K S^{K-1} \|\mathbf{g}_k - \mathbf{h}_k\|_q^q \\ &= K^{q-1} S^{K-1} \|\mathbf{g} - \mathbf{h}\|_q^q \end{aligned}$$

and therefore

$$\|P(\mathbf{h}) - P(\mathbf{g})\|_q \leq K^{1-\frac{1}{q}} S^{\frac{K-1}{q}} \|\mathbf{g} - \mathbf{h}\|_q.$$

Again, a similar calculus for $q = +\infty$ leads to

$$\|P(\mathbf{h}) - P(\mathbf{g})\|_{+\infty} \leq K \|\mathbf{g} - \mathbf{h}\|_{+\infty}.$$

Remember that the two last inequalities hold for \mathbf{g} and $\mathbf{h} \in \mathbb{R}^{S \times K}$ such that $\|\mathbf{g}\|_\infty \leq \|\mathbf{h}\|_\infty = 1$.

Let us now consider any \mathbf{g}' and $\mathbf{h}' \in \mathbb{R}^{S \times K}$ and any $\mathbf{g} \in \mathbb{R}_{\text{diag}}^{S \times K} \cap \langle \mathbf{g}' \rangle$ and $\mathbf{h} \in \mathbb{R}_{\text{diag}}^{S \times K} \cap \langle \mathbf{h}' \rangle$. We denote $\delta = \max(\|\mathbf{g}\|_{+\infty}, \|\mathbf{h}\|_{+\infty})$. Notice first that $\|\mathbf{g}\|_{+\infty} = \|P(\mathbf{g}')\|_{+\infty}^{\frac{1}{K}}$ and $\|\mathbf{h}\|_{+\infty} = \|P(\mathbf{h}')\|_{+\infty}^{\frac{1}{K}}$. Therefore

$$(10.11) \quad \delta = \max(\|P(\mathbf{g}')\|_{+\infty}, \|P(\mathbf{h}')\|_{+\infty})^{\frac{1}{K}}.$$

We can apply the above inequality to $\frac{\mathbf{h}}{\delta}$ and $\frac{\mathbf{g}}{\delta}$ (we might need to switch \mathbf{h} and \mathbf{g} but it does not change the final inequality) and obtain when $q < +\infty$

$$\|P\left(\frac{\mathbf{h}}{\delta}\right) - P\left(\frac{\mathbf{g}}{\delta}\right)\|_q \leq K^{1-\frac{1}{q}} S^{\frac{K-1}{q}} \left\| \frac{\mathbf{g}}{\delta} - \frac{\mathbf{h}}{\delta} \right\|_q.$$

This leads to

$$\|P(\mathbf{h}) - P(\mathbf{g})\|_q \leq K^{1-\frac{1}{q}} S^{\frac{K-1}{q}} \delta^{K-1} \|\mathbf{g} - \mathbf{h}\|_q.$$

Similarly, when $q = +\infty$, we obtain

$$\|P(\mathbf{h}) - P(\mathbf{g})\|_{+\infty} \leq K \delta^{K-1} \|\mathbf{g} - \mathbf{h}\|_{+\infty}.$$

The fact that these two last inequalities hold for any $\mathbf{g} \in \mathbb{R}_{\text{diag}}^{S \times K} \cap \langle \mathbf{g}' \rangle$ and any $\mathbf{h} \in \mathbb{R}_{\text{diag}}^{S \times K} \cap \langle \mathbf{h}' \rangle$, together with (10.11), leads to the statement provided in [Theorem 4.5](#). \square

10.5. Proof of Proposition 5.3. The span of the Segre variety $P(\mathbb{R}^{S \times K})$ is the full ambient space \mathbb{R}^{S^K} , so there exists sets of $R \leq S^K$ points on it that are linearly independent. The set of R -tuples of points on $P(\mathbb{R}^{S \times K})$ that fail to be linearly independent is a proper subvariety of the variety of sets of R -tuples of points on $P(\mathbb{R}^{S \times K})$ because being a linearly independent set of points is an open condition and there exists sets of points that are linearly independent. Therefore $R \leq S^K$ independent and randomly chosen points according to a continuous distribution on $P(\mathbb{R}^{S \times K})$ will be linearly independent.

The intersection $P(\mathbb{R}^{S \times K}) \cap \text{Ker}(\mathcal{A})$ is a proper subvariety of $P(\mathbb{R}^{S \times K})$, so with probability one, $R \leq S^K$ independent randomly chosen points according to a continuous distribution will not intersect it and be linearly independent. This is indeed the intersection of two non-empty open conditions. Therefore, all spans of subsets of the points will intersect $\text{Ker}(\mathcal{A})$ transversely (in particular, the span of fewer than $\text{rk}(\mathcal{A})$ points will not intersect it). Thus their image under \mathcal{A} will have dimension as large as possible. The same argument works if $R > S^K$. \square

10.6. Proof of Proposition 5.4. The proof relies on the fact that for any $T^* \in \operatorname{argmin}_{T \in \mathbb{R}^{S^K}} \|\mathcal{A}T - X\|^2$, we have

$$\mathcal{A}^t(\mathcal{A}T^* - X) = 0,$$

where $\mathcal{A}^t : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{S^K}$ is the adjoint linear map. This implies that for any

$$T^* \in \operatorname{argmin}_{T \in \mathbb{R}^{S^K}} \|\mathcal{A}T - X\|^2,$$

any $L \in \mathbb{N}$ and any $\mathbf{h} \in \mathcal{M}^L$

$$\begin{aligned} \|\mathcal{A}P(\mathbf{h}) - X\|^2 &= \|\mathcal{A}(P(\mathbf{h}) - T^*) + (\mathcal{A}T^* - X)\|^2, \\ &= \|\mathcal{A}(P(\mathbf{h}) - T^*)\|^2 + \|\mathcal{A}T^* - X\|^2 + 2\langle \mathcal{A}(P(\mathbf{h}) - T^*), \mathcal{A}T^* - X \rangle, \\ &= \|\mathcal{A}(P(\mathbf{h}) - T^*)\|^2 + \|\mathcal{A}T^* - X\|^2. \end{aligned}$$

In words, $\|\mathcal{A}P(\mathbf{h}) - X\|^2$ and $\|\mathcal{A}(P(\mathbf{h}) - T^*)\|^2$ only differ by an additive constant. Moreover, since the value of the objective function $\|\mathcal{A}T^* - X\|^2$ is independent of the particular minimizer T^* we are considering, this additive constant is independent of T^* . As a consequence, a minimizer of $\|\mathcal{A}P(\mathbf{h}) - X\|^2$ also minimizes $\|\mathcal{A}(P(\mathbf{h}) - T^*)\|^2$ and vice versa. \square

10.7. Proof of Proposition 6.2. Write $\bar{T} = P(\bar{\mathbf{h}})$ and let L^* and \mathbf{h}^* be a minimizer of (2.1). Proposition 5.4 and the fact that \bar{T} minimizes (5.4) implies that $(L^*, \mathbf{h}^*) \in \operatorname{argmin}_{L \in \mathbb{N}, \mathbf{h} \in \mathcal{M}^L} \|\mathcal{A}(P(\mathbf{h}) - \bar{T})\|^2$. As a consequence,

$$\|\mathcal{A}(P(\mathbf{h}^*) - \bar{T})\|^2 = 0$$

and

$$P(\mathbf{h}^*) \in \bar{T} + \operatorname{Ker}(\mathcal{A}),$$

proving the first implication.

Conversely, let $L^* \in \mathbb{N}$ and $\mathbf{h}^* \in \mathcal{M}^{L^*}$ be such that $P(\mathbf{h}^*) \in \bar{T} + \operatorname{Ker}(\mathcal{A})$, then

$$\|\mathcal{A}(P(\mathbf{h}^*) - \bar{T})\|^2 = 0 = \min_{L \in \mathbb{N}, \mathbf{h} \in \mathcal{M}^L} \|\mathcal{A}(P(\mathbf{h}) - \bar{T})\|^2.$$

As a consequence, $(L^*, \mathbf{h}^*) \in \operatorname{argmin}_{L \in \mathbb{N}, \mathbf{h} \in \mathcal{M}^L} \|\mathcal{A}(P(\mathbf{h}) - \bar{T})\|^2$ and, using Proposition 5.4, \mathbf{h}^* is a minimizer of (2.1). \square

10.8. Proof of Proposition 6.3.

- Proof of the first statement of Proposition 6.3:

We first assume that $\langle \bar{\mathbf{h}} \rangle$ is identifiable. We consider L^* and \mathbf{h}^* such that there is L^* such that $P(\mathbf{h}^*) \in (P(\bar{\mathbf{h}}) + \operatorname{Ker}(\mathcal{A})) \cap P(\mathcal{M}^{L^*})$. We know from Proposition 6.2 that $\mathbf{h}^* \in \operatorname{argmin}_{L \in \mathbb{N}, \mathbf{h} \in \mathcal{M}^L} \|\mathcal{A}P(\mathbf{h}) - X\|^2$. Using that $\langle \bar{\mathbf{h}} \rangle$ is identifiable, $\langle \mathbf{h}^* \rangle = \langle \bar{\mathbf{h}} \rangle$ and, from Standard Fact Item 1 (at the beginning of section 4), we get $P(\mathbf{h}^*) = P(\bar{\mathbf{h}})$. Finally, we can conclude, that if $\langle \bar{\mathbf{h}} \rangle$ is identifiable we have $(P(\bar{\mathbf{h}}) + \operatorname{Ker}(\mathcal{A})) \cap P(\mathcal{M}) \subset \{P(\bar{\mathbf{h}})\}$.

Let us assume now that for all $L \in \mathbb{N}$, $(P(\bar{\mathbf{h}}) + \operatorname{Ker}(\mathcal{A})) \cap P(\mathcal{M}^L) \subset \{P(\bar{\mathbf{h}})\}$ and consider

$$(L^*, \mathbf{h}^*) \in \operatorname{argmin}_{L \in \mathbb{N}, \mathbf{h} \in \mathcal{M}} \|\mathcal{A}P(\mathbf{h}) - X\|^2.$$

Using [Proposition 6.2](#), we know that $P(\mathbf{h}^*) \in (P(\bar{\mathbf{h}}) + \text{Ker}(\mathcal{A})) \cap P(\mathcal{M}^{L^*})$. Using the hypothesis, we have $P(\mathbf{h}^*) = P(\bar{\mathbf{h}})$ and using [Standard Fact Item 1](#), we finally conclude that $\langle \mathbf{h}^* \rangle = \langle \bar{\mathbf{h}} \rangle$. This completes the proof of the first statement.

- Proof of the second statement of [Proposition 6.3](#):

Assume that there is L and $L' \in \mathbb{N}$ such that $\text{Ker}(\mathcal{A}) \cap (P(\mathcal{M}^L) - P(\mathcal{M}^{L'})) \not\subseteq \{0\}$ then there exist $\mathbf{h} \in \mathcal{M}^L$ and $\bar{\mathbf{h}} \in \mathcal{M}^{L'}$ such that $P(\mathbf{h}) \neq P(\bar{\mathbf{h}})$ and $P(\mathbf{h}) - P(\bar{\mathbf{h}}) \in \text{Ker}(\mathcal{A})$. Using the first statement of the proposition, we obtain that $\bar{\mathbf{h}}$ is not identifiable. As a conclusion, \mathcal{M} is not identifiable.

Conversely, assume that there exists L' and some non-identifiable $\bar{\mathbf{h}} \in \mathcal{M}^{L'}$. Using the first statement of the proposition, we know that there exists $L \in \mathbb{N}$ and $\mathbf{h} \in \mathcal{M}^L$ such that $P(\mathbf{h}) \neq P(\bar{\mathbf{h}})$ and $P(\mathbf{h}) - P(\bar{\mathbf{h}}) \in \text{Ker}(\mathcal{A})$. Therefore $\text{Ker}(\mathcal{A}) \cap (P(\mathcal{M}) - P(\mathcal{M})) \not\subseteq \{0\}$. □

10.9. Proof of [Theorem 6.4](#). We first make the “equality holds generically” statement precise in our context. Fix any variety X and assume Y is a linear space, say of dimension y . Let $G(y, \mathbb{C}^N)$ denote the Grassmannian of y -planes through the origin in \mathbb{C}^N . The Grassmannian is both a smooth manifold and an algebraic variety. We can interpret “equality holds generically” in this context as saying for a Zariski open subset of $G(y, \mathbb{C}^N)$, equality will hold. In our situation, if we fix $\text{rk}(\mathcal{A})$ and allow $\text{ker}(\mathcal{A})$ to vary as a point in the Grassmannian, with probability one, it will intersect $J(P(\mathcal{M}^L), P(\mathcal{M}^{L'}))$ only in the origin, and this assertion is also true over \mathbb{R} because complex numbers are only needed to assure existence of intersections, not non-existence. □

10.10. Proof of [Theorem 7.2](#). We have

$$\begin{aligned} \|\mathcal{A}(P(\mathbf{h}^*) - P(\bar{\mathbf{h}}))\| &\leq \|\mathcal{A}P(\mathbf{h}^*) - X\| + \|\mathcal{A}P(\bar{\mathbf{h}}) - X\| \\ &\leq \delta + \eta \end{aligned}$$

Geometrically, this means that $P(\mathbf{h}^*)$ belongs to a cylinder centered at $P(\bar{\mathbf{h}})$ whose direction is $\text{Ker}(\mathcal{A})$ and whose section is defined using the operator \mathcal{A} . If we further decompose (the decomposition is unique)

$$P(\mathbf{h}^*) - P(\bar{\mathbf{h}}) = T + T',$$

where $T' \in \text{Ker}(\mathcal{A})$ and T is orthogonal to $\text{Ker}(\mathcal{A})$, we have

$$(10.12) \quad \|\mathcal{A}(P(\mathbf{h}^*) - P(\bar{\mathbf{h}}))\| = \|\mathcal{A}T\| \geq \sigma_{\min} \|T\|,$$

where σ_{\min} is the smallest non-zero singular value of \mathcal{A} . We finally obtain

$$\|P(\mathbf{h}^*) - P(\bar{\mathbf{h}}) - T'\| = \|T\| \leq \frac{\delta + \eta}{\sigma_{\min}}.$$

The term on the left-hand side corresponds to the distance between a point in $P(\mathcal{M}^{L^*}) - P(\mathcal{M}^L)$ (namely $P(\mathbf{h}^*) - P(\bar{\mathbf{h}})$) and a point in $\text{Ker}(\mathcal{A})$ (namely T'). ■

Since $\text{Ker}(\mathcal{A})$ satisfies the deep-NSP with constants (γ, ρ) , when $\delta + \eta \leq \rho$, we obtain the first inequality of the theorem

$$\|P(\mathbf{h}^*) - P(\bar{\mathbf{h}})\| \leq \gamma \frac{\delta + \eta}{\sigma_{\min}}.$$

When $\bar{\mathbf{h}} \in \mathbb{R}_*^{S \times K}$, for $\frac{\gamma}{\sigma_{\min}} (\delta + \eta) \leq \frac{1}{2} \max(\|P(\bar{\mathbf{h}})\|_\infty, \|P(\mathbf{h}^*)\|_\infty)$, we can apply [Theorem 4.3](#) and obtain [\(7.5\)](#). \square

10.11. Proof of [Theorem 7.3](#). Let \bar{L} and $\bar{L}' \in \mathbb{N}$ and $\bar{\mathbf{h}} \in \mathcal{M}^{\bar{L}}$ and $\bar{\mathbf{h}}' \in \mathcal{M}^{\bar{L}'}$ be such that $\|\mathcal{A}(P(\bar{\mathbf{h}}) - P(\bar{\mathbf{h}}'))\| \leq \delta$. We also consider throughout the proof $T' \in \text{Ker}(\mathcal{A})$. We assume that $\|P(\bar{\mathbf{h}})\|_\infty \leq \|P(\bar{\mathbf{h}}')\|_\infty$. If it is not the case, we simply switch $\bar{\mathbf{h}}$ and $\bar{\mathbf{h}}'$ in the definition of X and e below. We denote

$$X = \mathcal{A}P(\bar{\mathbf{h}}) \quad \text{and} \quad e = \mathcal{A}P(\bar{\mathbf{h}}) - \mathcal{A}P(\bar{\mathbf{h}}').$$

We have $X = \mathcal{A}P(\bar{\mathbf{h}}') + e$ and $\|e\| \leq \delta$. Therefore, the hypothesis of the theorem (applied with $\mathbf{h}^* = \bar{\mathbf{h}}$ and $L^* = \bar{L}$) guarantees that

$$d_2(\langle \bar{\mathbf{h}}, \bar{\mathbf{h}}' \rangle) \leq C \|P(\bar{\mathbf{h}}')\|_\infty^{\frac{1}{K}-1} \|e\|.$$

Using the fact that $e = \mathcal{A}P(\bar{\mathbf{h}}) - \mathcal{A}P(\bar{\mathbf{h}}')$ and $T' \in \text{Ker}(\mathcal{A})$ we obtain

$$\|e\| = \|\mathcal{A}(P(\bar{\mathbf{h}}) - P(\bar{\mathbf{h}}') - T')\| \leq \sigma_{\max} \|P(\bar{\mathbf{h}}) - P(\bar{\mathbf{h}}') - T'\|.$$

where σ_{\max} is the spectral radius of \mathcal{A} . Therefore

$$d_2(\langle \bar{\mathbf{h}}, \bar{\mathbf{h}}' \rangle) \leq C \|P(\bar{\mathbf{h}}')\|_\infty^{\frac{1}{K}-1} \sigma_{\max} \|P(\bar{\mathbf{h}}) - P(\bar{\mathbf{h}}') - T'\|.$$

Finally, using [Theorem 4.5](#) and the fact that $\|P(\bar{\mathbf{h}})\|_\infty \leq \|P(\bar{\mathbf{h}}')\|_\infty$, we obtain

$$\begin{aligned} \|P(\bar{\mathbf{h}}') - P(\bar{\mathbf{h}})\| &\leq S^{\frac{K-1}{2}} K^{1-\frac{1}{2}} \|P(\bar{\mathbf{h}}')\|_\infty^{1-\frac{1}{K}} d_2(\langle \bar{\mathbf{h}}', \bar{\mathbf{h}} \rangle) \\ &\leq CS^{\frac{K-1}{2}} \sqrt{K} \sigma_{\max} \|P(\bar{\mathbf{h}}) - P(\bar{\mathbf{h}}') - T'\| \\ &= \gamma \|P(\bar{\mathbf{h}}) - P(\bar{\mathbf{h}}') - T'\| \end{aligned}$$

for $\gamma = CS^{\frac{K-1}{2}} \sqrt{K} \sigma_{\max}$.

Summarizing, we conclude that under the hypothesis of the theorem: For any $T \in P(\mathcal{M}^{\bar{L}}) - P(\mathcal{M}^{\bar{L}'})$ such that $\|\mathcal{A}T\| \leq \delta$ we have for any $T' \in \text{Ker}(\mathcal{A})$

$$\|T\| \leq \gamma \|T - T'\|.$$

\square

10.12. Proof of [Proposition 8.2](#). Throughout the proof, we define, for any $\mathbf{i} \in [S]^K$, $\mathbf{h}^{\mathbf{i}} \in \mathbb{R}^{S \times K}$ by

$$(10.13) \quad \mathbf{h}_{k,j}^{\mathbf{i}} = \begin{cases} 1 & \text{if } j = \mathbf{i}_k \\ 0 & \text{otherwise} \end{cases}, \quad \text{for all } k \in [K] \text{ and } j \in [S].$$

This notation shall not be confused with $\mathbf{h}^{\mathbf{p}}$, with $\mathbf{p} \in \mathcal{P}$.

- Let us first prove the first statement: We can easily check that $(P(\mathbf{h}^{\mathbf{i}}))_{\mathbf{i} \in \mathbf{I}}$ forms a basis of $\{T \in \mathbb{R}^{S^K} \mid \forall \mathbf{i} \in \mathbf{I}, T_{\mathbf{i}} = 0\}$. We can also easily check using (8.2) that, for any $\mathbf{i} \notin \mathbf{I}$,

$$\mathcal{A}P(\mathbf{h}^{\mathbf{i}}) = M_1(\mathbf{h}_1^{\mathbf{i}}) \dots M_K(\mathbf{h}_K^{\mathbf{i}}) = 0.$$

Therefore, $\{T \in \mathbb{R}^{S^K} \mid \forall \mathbf{i} \in \mathbf{I}, T_{\mathbf{i}} = 0\} \subset \text{Ker}(\mathcal{A})$.

Conversely, for any $\mathbf{i} \in \mathbf{I}$, we can deduce from (8.2) and the hypotheses of the proposition that all the entries of $\mathcal{A}P(\mathbf{h}^{\mathbf{i}})$ are in $\{0, 1\}$. We denote $D_{\mathbf{i}} = \{(i, j) \in [N] \times [N|\mathcal{F}|] \mid \mathcal{A}P(\mathbf{h}^{\mathbf{i}})_{i,j} = 1\}$. Using (again) the hypothesis of the proposition and (8.2), we can prove that, for any distinct \mathbf{i} and $\mathbf{j} \in \mathbf{I}$, we have $D_{\mathbf{i}} \cap D_{\mathbf{j}} = \emptyset$. This easily leads to the Item 1 of the first statement. We also deduce that

$$\text{rk}(\mathcal{A}) \geq |\mathbf{I}| = S^K - \dim(\{T \in \mathbb{R}^{S^K} \mid \forall \mathbf{i} \in \mathbf{I}, T_{\mathbf{i}} = 0\}).$$

Finally, we deduce that $\dim(\text{Ker}(\mathcal{A})) \leq \dim(\{T \in \mathbb{R}^{S^K} \mid \forall \mathbf{i} \in \mathbf{I}, T_{\mathbf{i}} = 0\})$ and therefore

$$\text{Ker}(\mathcal{A}) = \{T \in \mathbb{R}^{S^K} \mid \forall \mathbf{i} \in \mathbf{I}, T_{\mathbf{i}} = 0\}.$$

- Let us now prove the second statement: Using the hypothesis of the second statement and (8.2), we know that there is $f \in \mathcal{F}$ and $n \in [N]$ such that

$$\sum_{p \in \mathcal{P}(f)} \mathcal{T}^p(1)_n \geq 2.$$

As a consequence, there is \mathbf{i} and $\mathbf{j} \in [S]^K$ with $\mathbf{i} \neq \mathbf{j}$ and

$$\mathcal{T}^{\mathbf{p}_i}(\mathbf{h}^{\mathbf{i}})_n = \mathcal{T}^{\mathbf{p}_j}(\mathbf{h}^{\mathbf{j}})_n = 1.$$

Therefore, $\mathcal{A}P(\mathbf{h}^{\mathbf{i}}) = \mathcal{A}P(\mathbf{h}^{\mathbf{j}})$ and the network is not identifiable. \square

10.13. Proof of Proposition 8.3. The fact that, under the hypotheses of the proposition, $\text{Ker}(\mathcal{A}) = \{0\}$ is a direct consequence of Proposition 8.2. The deep-NSP property and the value of γ also follow from the definition of the deep-NSP.

To calculate σ_{\min} , let us consider $T \in \mathbb{R}^{S^K}$ and express it under the form $T = \sum_{\mathbf{i} \in \mathbf{I}} T_{\mathbf{i}} P(\mathbf{h}^{\mathbf{i}})$, where $\mathbf{h}^{\mathbf{i}}$ is defined (10.13). Let us also remind that, applying Proposition 8.2, the supports of $\mathcal{A}P(\mathbf{h}^{\mathbf{i}})$ and $\mathcal{A}P(\mathbf{h}^{\mathbf{j}})$ are disjoint, when $\mathbf{i} \neq \mathbf{j}$. Let us finally add that, since $\mathcal{A}P(\mathbf{h}^{\mathbf{j}})$ is the matrix of a convolution with a Dirac mass, its support is of size N . We finally have

$$\begin{aligned} \|AT\|^2 &= \left\| \sum_{\mathbf{i} \in \mathbf{I}} T_{\mathbf{i}} \mathcal{A}P(\mathbf{h}^{\mathbf{i}}) \right\|^2, \\ &= N \sum_{\mathbf{i} \in \mathbf{I}} T_{\mathbf{i}}^2 = N \|T\|^2, \end{aligned}$$

from which we deduce the value of σ_{\min} . \square

10.14. Proof of Theorem 8.4. Considering Proposition 8.2, we only need to prove that the condition is sufficient to guarantee the parameter stability.

Let us consider a path $\mathbf{p} \in \mathcal{P}$, using (8.2), since all the entries of $M_1(1) \dots M_K(1)$ belong to $\{0, 1\}$, all the entries of $M_1(\mathbf{1}^{\mathbf{P}}) \dots M_K(\mathbf{1}^{\mathbf{P}})$ belong to $\{0, 1\}$. Therefore, we can apply Proposition 8.3 and Theorem 7.2 to the restriction of the convolutional linear network to \mathbf{p} and obtain

$$d_p(\langle (\mathbf{h}^*)^{\mathbf{P}} \rangle, \langle \bar{\mathbf{h}}^{\mathbf{P}} \rangle) \leq \frac{7(KS')^{\frac{1}{p}}}{\sqrt{N}} \min \left(\|P(\bar{\mathbf{h}}^{\mathbf{P}})\|_{\infty}^{\frac{1}{K}-1}, \|P((\mathbf{h}^*)^{\mathbf{P}})\|_{\infty}^{\frac{1}{K}-1} \right) (\delta^{\mathbf{P}} + \eta^{\mathbf{P}}),$$

where $\delta^{\mathbf{P}}$ and $\eta^{\mathbf{P}}$ are the restrictions of the errors on $\text{Supp}(\mathbf{p})$.

We therefore have

$$d_p(\langle (\mathbf{h}^*)^{\mathbf{P}} \rangle, \langle \bar{\mathbf{h}}^{\mathbf{P}} \rangle) \leq \frac{7(KS')^{\frac{1}{p}}}{\sqrt{N}} \varepsilon^{1-K} (\delta^{\mathbf{P}} + \eta^{\mathbf{P}}),$$

and finally

$$\begin{aligned} \delta_p(\{\mathbf{h}^*\}, \{\bar{\mathbf{h}}\}) &\leq \frac{7(KS')^{\frac{1}{p}} \varepsilon^{1-K}}{\sqrt{N}} \left(\sum_{\mathbf{p} \in \mathcal{P}} (\delta^{\mathbf{P}} + \eta^{\mathbf{P}})^p \right)^{\frac{1}{p}}, \\ &\leq \frac{7(KS')^{\frac{1}{p}} \varepsilon^{1-K}}{\sqrt{N}} (\delta + \eta). \end{aligned}$$

□

REFERENCES

- [1] Alekh Agarwal, Animashree Anandkumar, and Praneeth Netrapalli. Exact recovery of sparsely used overcomplete dictionaries. *stat*, 1050:8, 2013.
- [2] Arif Ahmed, Benjamin Recht, and Justin Romberg. Blind deconvolution using convex programming. *IEEE Transactions on Information Theory*, 60(3):1711–1732, 2014.
- [3] Alexandr Andoni, Rina Panigrahy, Gregory Valiant, and Li Zhang. Learning polynomials with neural networks. In *International Conference on Machine Learning*, pages 1908–1916, 2014.
- [4] Sanjeev Arora, Aditya Bhaskara, Rong Ge, and Tengyu Ma. Provable bounds for learning some deep representations. In *International Conference on Machine Learning*, pages 584–592, 2014.
- [5] Sanjeev Arora, Rong Ge, Ravindran Kannan, and Ankur Moitra. Computing a nonnegative matrix factorization—provably. In *Proceedings of the forty-fourth annual ACM symposium on Theory of computing*, pages 145–162. ACM, 2012.
- [6] Sanjeev Arora, Rong Ge, and Ankur Moitra. New algorithms for learning incoherent and overcomplete dictionaries. In *COLT*, pages 779–806, 2014.
- [7] Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. *arXiv preprint arXiv:1802.05296*, 2018.
- [8] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140, 2015.
- [9] Sohail Bahmani and Justin Romberg. Lifting for blind deconvolution in random mask imaging: Identifiability and convex relaxation. *SIAM Journal on Imaging Sciences*, 8(4):2203–2238, 2015.
- [10] Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.

- [11] Pierre F Baldi and Kurt Hornik. Learning in linear neural networks: A survey. *IEEE Transactions on neural networks*, 6(4):837–858, 1995.
- [12] Andrew R Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, 39(3):930–945, 1993.
- [13] Anthony Bourrier, Mike Davies, Tomer Peleg, Patrick Pérez, and Rémi Gribonval. Fundamental performance limits for ideal decoders in high-dimensional linear inverse problems. *IEEE Transactions on Information Theory*, 60(12):7928–7946, 2014.
- [14] Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with gaussian inputs. *arXiv preprint arXiv:1702.07966*, 2017.
- [15] Emmanuel Candès, Justin Romberg, and Terence Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on information theory*, 52(2):489–509, 2006.
- [16] Emmanuel J Candes, Yonina C Eldar, Thomas Strohmer, and Vladislav Voroninski. Phase retrieval via matrix completion. *SIAM review*, 57(2):225–251, 2015.
- [17] Emmanuel J Candès and Benjamin Recht. Exact matrix completion via convex optimization. *Foundations of Computational mathematics*, 9(6):717, 2009.
- [18] Emmanuel J Candes, Thomas Strohmer, and Vladislav Voroninski. Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 66(8):1241–1274, 2013.
- [19] Emmanuel J Candès and Terence Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080, 2010.
- [20] Olivier Chabiron, François Malgouyres, Jean-Yves Tournéret, and Nicolas Dobigeon. Toward fast transform learning. *International Journal of Computer Vision*, pages 1–22, 2014.
- [21] Olivier Chabiron, François Malgouyres, Herwig Wendt, and Jean-Yves Tournéret. Optimization of a fast transform structured as a convolutional tree. *preprint HAL*, (hal-01258514), 2016.
- [22] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. The loss surfaces of multilayer networks. In *Artificial Intelligence and Statistics*, pages 192–204, 2015.
- [23] Anna Choromanska, Yann LeCun, and Gérard Ben Arous. Open problem: The landscape of the loss surfaces of multilayer networks. In *Conference on Learning Theory*, pages 1756–1760, 2015.
- [24] Sunav Choudhary and Urbashi Mitra. Identifiability scaling laws in bilinear inverse problems. *arXiv preprint arXiv:1402.2637*, 2014.
- [25] Albert Cohen, Wolfgang Dahmen, and Ronald DeVore. Compressed sensing and best k -term approximation. *Journal of the American mathematical society*, 22(1):211–231, 2009.
- [26] Nadav Cohen, Or Sharir, and Amnon Shashua. On the expressive power of deep learning: A tensor analysis. In *Conference on Learning Theory*, pages 698–728, 2016.
- [27] Nadav Cohen and Amnon Shashua. Convolutional rectifier networks as generalized tensor decompositions. In *International Conference on Machine Learning*, pages 955–963, 2016.
- [28] Chandler Davis and William Morton Kahan. The rotation of eigenvectors by a perturbation. iii. *SIAM Journal on Numerical Analysis*, 7(1):1–46, 1970.
- [29] David Donoho and Victoria Stodden. When does non-negative matrix factorization give a correct decomposition into parts? In *Advances in neural information processing systems*, pages 1141–1148, 2004.
- [30] David L. Donoho. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006.
- [31] Ronen Eldan and Ohad Shamir. The power of depth for feedforward neural networks. In *Conference on Learning Theory*, pages 907–940, 2016.
- [32] M Fazel, E Candes, B Recht, and P Parrilo. Compressed sensing and robust recovery of low rank matrices. In *Signals, Systems and Computers, 2008 42nd Asilomar Conference on*, pages 1043–1047. IEEE, 2008.
- [33] Quan Geng, Huan Wangy, and John Wright. On the local correctness of ℓ_1 -minimization for dictionary learning. In *International Symposium on Information Theory (ISIT)*, 2014.
- [34] Surbhi Goel and Adam Klivans. Eigenvalue decay implies polynomial-time learnability for neural networks. In *Advances in Neural Information Processing Systems*, pages 2192–2202, 2017.
- [35] Rémi Gribonval, Rodolphe Jenatton, and Francis Bach. Sample complexity of dictionary learning and other matrix factorizations. *Information Theory, IEEE Transactions on*, 61(6):3469–3486, June 2015.
- [36] Rémi Gribonval and Karin Schnass. Dictionary identification - sparse matrix-factorisation via ℓ_1 -minimisation. *IEEE transaction on information theory*, 56(7):3523–3539, 2010.
- [37] Benjamin D Haefele and René Vidal. Global optimality in tensor factorization, deep learning,

- and beyond. *arXiv preprint arXiv:1506.07540*, 2015.
- [38] Joe Harris. *Algebraic geometry*, volume 133 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 1995. A first course, Corrected reprint of the 1992 original.
- [39] Jonathan D. Hauenstein and Andrew J. Sommese. Witness sets of projections. *Applied Mathematics and Computation*, 217(7):3349 – 3354, 2010.
- [40] Jonathan D. Hauenstein and Andrew J. Sommese. Membership tests for images of algebraic sets by linear projections. *Applied Mathematics and Computation*, 219(12):6809 – 6818, 2013.
- [41] Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. Generating visual explanations. In *European Conference on Computer Vision*, pages 3–19. Springer, 2016.
- [42] Majid Janzamin, Hanie Sedghi, and Anima Anandkumar. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods. *arXiv preprint arXiv:1506.08473*, 2015.
- [43] Rodolphe Jenatton, Rémi Gribonval, and Francis Bach. Local stability and robustness of sparse dictionary learning in the presence of noise. arxiv, 2012.
- [44] Kenji Kawaguchi. Deep learning without poor local minima. In *Advances in Neural Information Processing Systems*, pages 586–594, 2016.
- [45] Valentin Khruikov, Alexander Novikov, and Ivan Oseledets. Expressive power of recurrent neural networks. *arXiv preprint arXiv:1711.00811*, 2017.
- [46] Risi Kondor, Nedelina Teneva, and Vikas Garg. Multiresolution matrix factorization. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1620–1628, 2014.
- [47] Joseph M. Landsberg. *Tensors: Geometry and Applications*, volume 128 of *Graduate studies in mathematics*. American Mathematical Soc., 2012.
- [48] Hans Laurberg, Mads Græsbøll Christensen, Mark D Plumbley, Lars Kai Hansen, and Søren Holdt Jensen. Theorems on positive data: On the uniqueness of nmf. *Computational intelligence and neuroscience*, 2008, 2008.
- [49] Vadim Lebedev, Yaroslav Ganin, Maksim Rakhuba, Ivan Oseledets, and Victor Lempitsky. Speeding-up convolutional neural networks using fine-tuned cp-decomposition. *arXiv preprint arXiv:1412.6553*, 2014.
- [50] Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
- [51] Xiaodong Li, Shuyang Ling, Thomas Strohmer, and Ke Wei. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *CoRR*, abs/1606.04933, 2016.
- [52] Xiaodong Li and Vladislav Voroninski. Sparse signal recovery from quadratic measurements via convex programming. *SIAM Journal on Mathematical Analysis*, 45(5):3019–3033, 2013.
- [53] Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu activation. In *Advances in Neural Information Processing Systems*, pages 597–607, 2017.
- [54] Shuyang Ling and Thomas Strohmer. Self-calibration and biconvex compressive sensing. *Inverse Problems*, 31(11):115002, 2015.
- [55] Roi Livni, Shai Shalev-Shwartz, and Ohad Shamir. On the computational efficiency of training neural networks. In *Advances in Neural Information Processing Systems*, pages 855–863, 2014.
- [56] Siwei Lyu and Xin Wang. On algorithms for sparse multi-factor nmf. In *Proceedings of the 26th International Conference on Neural Information Processing Systems, NIPS’13*, pages 602–610, USA, 2013. Curran Associates Inc.
- [57] Luc Le Magoarou. *Matrices efficaces pour le traitement du signal et l’apprentissage automatique*. PhD thesis, Université Bretagne Loire, 2016.
- [58] Luc Le Magoarou and Rémi Gribonval. Are there approximate fast fourier transforms on graphs? In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4811–4815, 2016.
- [59] Luc Le Magoarou and Rémi Gribonval. Flexible multi-layer sparse approximations of matrices and applications. *IEEE Journal of Selected Topics in Signal Processing*, 10(4):688–700, 2016.
- [60] François Malgouyres and Joseph Landsberg. On the identifiability and stable recovery of deep/multi-layer structured matrix factorization. In *IEEE, Info. Theory Workshop*, Sept. 2016.
- [61] Stéphane Mallat. *A Wavelet Tour of Signal Processing*. Academic Press, Boston, 1998.
- [62] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 65:211–222, 2017.

- [63] Guido F Montufar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio. On the number of linear regions of deep neural networks. In *Advances in neural information processing systems*, pages 2924–2932, 2014.
- [64] Behnam Neyshabur and Rina Panigrahy. Sparse matrix factorization. *arXiv preprint arXiv:1311.3315*, 2013.
- [65] Alexander Novikov, Dmitrii Podoprikin, Anton Osokin, and Dmitry P Vetrov. Tensorizing neural networks. In *Advances in Neural Information Processing Systems*, pages 442–450, 2015.
- [66] Benjamin Recht, Maryam Fazel, and Pablo A Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review*, 52(3):471–501, 2010.
- [67] Ron Rubinstein, Alfred Bruckstein, and Michael Elad. Dictionaries for sparse representation modeling. *Proc. IEEE - Special issue on applications of sparse representation and compressive sensing*, 98(6):1045–1057, 2010.
- [68] Itay Safran and Ohad Shamir. On the quality of the initial basin in overspecified neural networks. In *International Conference on Machine Learning*, pages 774–782, 2016.
- [69] Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation function. *arXiv preprint arXiv:1708.06633*, 2017.
- [70] Karin Schnass. On the identifiability of overcomplete dictionaries via the minimisation principle underlying k-svd. *Applied and Computational Harmonic Analysis*, 37(3):464–491, 2014.
- [71] Hanie Sedghi and Anima Anandkumar. Provable methods for training neural networks with sparse connectivity. In *Deep Learning and representation learning workshop: NIPS*, 2014.
- [72] Igor R. Shafarevich. *Basic algebraic geometry. 1*. Springer, Heidelberg, third edition, 2013. Varieties in projective space.
- [73] Richard Socher, Danqi Chen, Christopher D Manning, and Andrew Ng. Reasoning with neural tensor networks for knowledge base completion. In *Advances in neural information processing systems*, pages 926–934, 2013.
- [74] Daniel Spielmana, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. In *COLT*, pages 37–1, 2012.
- [75] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of machine learning research*, 15(1):1929–1958, 2014.
- [76] Matus Telgarsky. Benefits of depth in neural networks. *arXiv preprint arXiv:1602.04485*, 2016.
- [77] Theodoros Tsiligkaridis, Alfred O Hero, and Shuheng Zhou. On convergence of kronecker graphical lasso algorithms. *IEEE transactions on signal processing*, 61(7):1743–1755, 2013.
- [78] Stephen A Vavasis. On the complexity of nonnegative matrix factorization. *SIAM Journal on Optimization*, 20(3):1364–1377, 2009.
- [79] L. Venturi, A. S. Bandeira, and J. Bruna. Neural Networks with Finite Intrinsic Dimension have no Spurious Valleys. *ArXiv e-prints*, February 2018.
- [80] Bo Xie, Yingyu Liang, and Le Song. Diversity leads to generalization in neural networks. *arXiv preprint Arxiv:1611.03131*, 2016.
- [81] Dong Yu, Li Deng, and Frank Seide. The deep tensor neural network with applications to large vocabulary speech recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, 21(2):388–396, 2013.
- [82] Chen Yunpeng, Jin Xiaojie, Kang Bingyi, Feng Jiashi, and Yan Shuicheng. Sharing residual units through collective tensor factorization in deep neural networks. *arXiv preprint arXiv:1703.02180*, 2017.
- [83] Kai Zhong, Zhao Song, Prateek Jain, Peter L. Bartlett, and Inderjit S. Dhillon. Recovery guarantees for one-hidden-layer neural networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 4140–4149, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.