

UNDERSTANDING ℓ^4 -BASED DICTIONARY LEARNING: INTERPRETATION, STABILITY, AND ROBUSTNESS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recently ℓ^4 -norm maximization has been proposed to solve the sparse dictionary learning (SDL) problem. The simple MSP (matching, stretching, and projection) algorithm proposed by Zhai et al. (2019a) has shown to be surprisingly efficient and effective. This paper aims to better understand this algorithm from its strong geometric and statistical connections with the classic PCA and ICA, as well as their associated fixed-point style algorithms. Such connections provide a unified way of viewing problems that pursue *principal*, *independent*, or *sparse* components of high-dimensional data. Our studies reveal additional good properties of the ℓ^4 -maximization: not only is the MSP algorithm for sparse coding insensitive to small noise, but also robust to outliers, and resilient to sparse corruptions. We provide preliminary statistical justification for such inherently nice properties. To corroborate the theoretical analysis, we also provide extensive and compelling experimental evidence with both synthetic data and real images.

1 INTRODUCTION

The explosion of massive amounts of high-dimensional data has become the modern-day norm for a large number of scientific and engineering disciplines and hence presents a daunting challenge for both computation and learning. Rising to this challenge, *sparse dictionary learning* (SDL) provides a potent framework in representation learning that exploits the blessing of dimensionality: real data tends to lie in or near some low-dimensional subspaces or manifolds, even though the ambient dimension is often extremely large (e.g. the number of raw pixels in an image). More specifically, SDL (Olshausen & Field (1997); Mairal et al. (2008; 2012; 2014); Spielman et al. (2012); Sun et al. (2015); Bai et al. (2018)) concerns the problem of learning a compact, sparse representation from raw, unlabelled data: given a data matrix $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_p] \in \mathbb{R}^{n \times p}$ that contains p n -dimensional samples, one aims to find a linear transformation (i.e. a *dictionary*) $\mathbf{D} \in \mathbb{R}^{n \times m}$ and an associated maximally sparse representation $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p] \in \mathbb{R}^{m \times p}$ that satisfies

$$\mathbf{Y} = \mathbf{D}\mathbf{X}. \tag{1}$$

As the data matrix \mathbf{Y} can represent a variety of signals (e.g. images, audios, languages, and genetics etc) in practical applications, SDL provides a versatile structure-seeking formulation that has found widespread applications in computational neuroscience, image processing and computer vision, and machine learning at large (Olshausen & Field (1996; 1997); Argyriou et al. (2008); Ranzato et al. (2007); Elad & Aharon (2006); Wright et al. (2008); Yang et al. (2010); Mairal et al. (2014)).

Related Work. Motivated by this practical significance, there has been a growing surge of interest recently (e.g. Rambhatla et al. (2019); Bai et al. (2018); Gilboa et al. (2018); Nguyen et al. (2018); Chatterji & Bartlett (2017); Mensch et al. (2016)) that aims to tackle SDL. In attempts to recover the sparse signals \mathbf{X} , these existing work adopt an ℓ^0 - or ℓ^1 -penalty function to promote the underlying sparsity and give various optimization algorithms for the resulting objectives (some of those are heuristics while a few others have theoretical convergence guarantees). Although these penalty functions are indeed sparsity-promoting, the resulting optimization problems must be solved one row at a time, hence resulting as many optimization problems as the ambient dimension n . Consequently, ℓ^0 - or ℓ^1 -based objectives result only in local methods (i.e. cannot yield the entire solution at once) and hence entail prohibitive computational burden. Another prominent approach in SDL is Sum-of-Squares (SOS), proposed by and articulated in a series of recent work Barak et al. (2015); Ma et al. (2016); Schramm & Steurer (2017). The key idea there is to utilize the properties of higher order SOS polynomials to correctly recover one column of the dictionary at a time, for which there are m columns in total. However, the computational complexity of these recovery methods are quasi-polynomial, hence again resulting in large computational expense.

Very recently, in the *complete dictionary learning*¹ setting, a novel global approach has been suggested in Zhai et al. (2019a;b) that presents a formulation that can efficiently recover the sparse signal matrix \mathbf{X} once for all. In particular, Zhai et al. (2019b) has shown that if the generative model for $\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o \in \mathbb{R}^{n \times p}$ satisfies that $\mathbf{D}_o \in \mathcal{O}(n; \mathbb{R})$ is orthonormal and $\mathbf{X}_o \in \mathbb{R}^{n \times p}$ is Bernoulli-Gaussian sparse,² then maximizing the ℓ^4 -norm³ of $\mathbf{A}\mathbf{Y}$ over $\mathcal{O}(n; \mathbb{R})$:

$$\max_{\mathbf{A}} \frac{1}{4} \|\mathbf{A}\mathbf{Y}\|_4^4 \quad \text{subject to } \mathbf{A} \in \mathcal{O}(n; \mathbb{R}) \quad (\text{or } \mathbf{A}\mathbf{A}^* = \mathbf{I}), \quad (2)$$

is able to find the ground truth dictionary \mathbf{D}_o up to an arbitrary signed permutation. Moreover, Zhai et al. (2019b) has proposed the simple “*Matching, Stretching, and Projection*” (MSP) algorithm, which is shown to be experimentally efficient and effective, for solving the program in equation 2:

$$\text{MSP: } \mathbf{A}_{t+1} = \mathcal{P}_{\mathcal{O}(n; \mathbb{R})} [(\mathbf{A}_t \mathbf{Y})^{\circ 3} \mathbf{Y}^*] = \mathbf{U}_t \mathbf{V}_t^*, \quad (3)$$

where $\mathbf{U}_t \mathbf{V}_t^*$ are from the singular value decomposition: $\mathbf{U}_t \Sigma_t \mathbf{V}_t^* = \text{SVD}[(\mathbf{A}_t \mathbf{Y})^{\circ 3} \mathbf{Y}^*]$.

To serve purposes of this paper, we here give an alternative (arguably simpler and more revealing) derivation of the MSP Algorithm 3. Consider the *Lagrangian* formulation of the constrained optimization equation 2, the necessary condition of critical points $\nabla_{\mathbf{A}} \frac{1}{4} \|\mathbf{A}\mathbf{Y}\|_4^4 = \nabla_{\mathbf{A}} \langle \mathbf{A}, \mathbf{A}\mathbf{A}^* - \mathbf{I} \rangle$ for some Lagrangian multipliers $\mathbf{\Lambda}$ implies:

$$(\mathbf{A}\mathbf{Y})^{\circ 3} \mathbf{Y}^* = (\mathbf{\Lambda} + \mathbf{\Lambda}^*) \mathbf{A}. \quad (4)$$

As the optimization is over the orthogonal group $\mathcal{O}(n; \mathbb{R})$, restricting the condition in equation 4 onto the orthogonal group yields a necessary condition for any critical point \mathbf{A} :⁴

$$\mathcal{P}_{\mathcal{O}(n; \mathbb{R})} [(\mathbf{A}\mathbf{Y})^{\circ 3} \mathbf{Y}^*] = \mathbf{A}. \quad (5)$$

Hence the critical point \mathbf{A} can be viewed as a “*fixed point*” of the map: $\mathcal{P}_{\mathcal{O}(n; \mathbb{R})} [((\cdot)\mathbf{Y})^{\circ 3} \mathbf{Y}^*]$ from $\mathcal{O}(n; \mathbb{R})$ to itself. The MSP algorithm in equation 3 is to find the fixed point(s) of this map.

Notice that the orthonormal constraint $\mathbf{A} \in \mathcal{O}(n; \mathbb{R})$ in equation 2 can be viewed as enforcing the orthogonality of n unit vectors simultaneously. So, more flexibly and generally, one may choose to compute any k , for $1 \leq k \leq n$, leading orthonormal bases of \mathbf{D}_o by solving the program:

$$\max_{\mathbf{W}} \frac{1}{4} \|\mathbf{W}^* \mathbf{Y}\|_4^4 \quad \text{subject to } \mathbf{W} \in \text{St}(k, n; \mathbb{R}) \subset \mathbb{R}^{n \times k}, \quad (6)$$

where $\text{St}(k, n; \mathbb{R})$ is the *Stiefel manifold*.⁵ The orthogonal group $\mathcal{O}(n; \mathbb{R})$ and the unit sphere \mathbb{S}^{n-1} can be viewed as two special cases of the Stiefel manifold $\text{St}(k, n; \mathbb{R})$, with $k = n$ and $k = 1$, respectively. In some specific tasks such as dictionary learning and blind deconvolution, optimization over the unit sphere has been widely practiced, see Sun et al. (2015); Bai et al. (2018); Zhang et al. (2018); Kuo et al. (2019). The more general setting of maximizing a convex function over any compact set also has been studied by Journée et al. (2010) in the context of sparse PCA, which has provided convergence guarantees for this class of programs.

Our Contributions. Our contributions are twofold. First, by taking a suitable analytical angle, we reveal novel geometric and statistical connections between PCA, ICA and the ℓ^4 -norm maximization based SDL. We then show that algorithm-wise, the fixed-point type MSP algorithm for ℓ^4 -norm maximization has the same nature as the classic *power-iteration method* for PCA Jolliffe (2011) and the FastICA algorithm for ICA Hyvärinen & Oja (1997). This interpretation gives a unified view for problems that pursue *principal, independent, or sparse* components from high-dimensional data and enriches our understanding of low-dimensional structure recovery frameworks, classical and new, at both formulation and algorithmic fronts.

Second, and more importantly from a practical perspective, we examine how MSP performs under a variety of more realistic conditions, when the measurements \mathbf{Y} could be contaminated with noise,

¹Complete dictionary learning requires the learned dictionary \mathbf{D} in equation 1 to be square and invertible.

²Each entry $x_{i,j}$ of \mathbf{X} can be represented as the product of a Bernoulli variable and a normal Gaussian variable: $x_{i,j} = \Omega_{i,j} V_{i,j}$, where $\Omega_{i,j} \sim_{iid} \text{Ber}(\theta)$ and $V_{i,j} \sim_{iid} \mathcal{N}(0, 1)$, similar for vectors or scalars. This is the standard setting adopted in Spielman et al. (2012); Sun et al. (2015); Bai et al. (2018).

³We abuse the notation a bit, by denoting $\|\cdot\|_4^4$ as the sum of element-wise 4th power of all entries of a vector and matrix, that is, $\forall \mathbf{a} \in \mathbb{R}^n$, $\|\mathbf{a}\|_4^4 = \sum_{i=1}^n a_i^4$ and $\forall \mathbf{A} \in \mathbb{R}^{n \times m}$, $\|\mathbf{A}\|_4^4 = \sum_{i,j} a_{i,j}^4$.

⁴For any symmetric matrix $\mathbf{S} \in \mathbb{R}^{n \times n}$ and an orthogonal matrix $\mathbf{A} \in \mathcal{O}(n; \mathbb{R})$, the projection of $\mathbf{S}\mathbf{A}$ onto the orthogonal group is \mathbf{A} : $\mathcal{P}_{\mathcal{O}(n; \mathbb{R})} [\mathbf{S}\mathbf{A}] = \mathbf{A}$, one may see Absil & Malick (2012) for details.

⁵For any $1 \leq k \leq n$, $\text{St}(k, n; \mathbb{R}) \doteq \{\mathbf{W} \in \mathbb{R}^{n \times k} : \mathbf{W}^* \mathbf{W} = \mathbf{I}_k\}$.

outliers, or sparse corruptions. We show that, similar to PCA, ℓ^4 -norm maximization and the MSP algorithm are inherently stable to small noise. Somewhat surprisingly though, unlike PCA, the MSP algorithm is further robust to outliers and resilient to sparse gross errors! We provide characterizations of these desirable properties of MSP. The claims are further corroborated with extensive experiments on both synthetic data and real images. Taken as a whole, our results contribute to the broad landscape of dictionary learning by affirming that ℓ^4 -maximization based SDL and the corresponding global algorithm MSP provide a valuable toolkit to the existing literature.

2 SDL VERSUS PCA AND ICA

2.1 PURSUIT OF PRINCIPAL, INDEPENDENT, OR SPARSE COMPONENTS

Relation with the Geometric Interpretation of PCA. For a data matrix $\mathbf{Y} \in \mathbb{R}^{n \times p}$, *Principal Component Analysis* (PCA) aims to find the top (top k) left singular vector (vectors) of \mathbf{Y} :

$$\max_{\mathbf{W}} \frac{1}{2} \|\mathbf{W}^* \mathbf{Y}\|_F^2 \quad \text{subject to } \mathbf{W} \in \text{St}(k, n; \mathbb{R}), \quad (7)$$

can be considered as finding a direction (a k -dimensional subspace) in $\text{row}(\mathbf{Y})$ in which \mathbf{Y} has the largest ℓ^2 -norm (Frobenius norm). For instance, finding the direction with the largest ℓ^2 -norm over the unit sphere can be viewed as calculating the spectral norm (or the largest singular value), of matrix \mathbf{Y} . In comparison, we may view equation 6

$$\max_{\mathbf{W}} \frac{1}{4} \|\mathbf{W}^* \mathbf{Y}\|_4^4 \quad \text{subject to } \mathbf{W} \in \text{St}(k, n; \mathbb{R})$$

as to find a direction, or a k -dimensional subspace, in $\text{row}(\mathbf{Y})$ where the projection of \mathbf{Y} has the largest ℓ^4 -norm. For instance, finding the direction with the largest ℓ^4 -norm over the unit sphere equation 6 can be viewed as calculating the induced $\|\cdot\|_{2,4}$ norm of matrix \mathbf{Y} : $\|\mathbf{Y}\|_{2,4} \doteq \max_{\mathbf{a} \in \mathbb{S}^{n-1}} \|\mathbf{a}^* \mathbf{Y}\|_4$.

Relation with the Statistical Interpretation of PCA. If we view each column $\mathbf{y}_j, j \in [p]$ of data matrix $\mathbf{Y} \in \mathbb{R}^{n \times p}$ as an n dimensional random vector that is i.i.d. drawn from a distribution of random variable \mathbf{y} and let \mathbf{Y}_c denote the centered \mathbf{Y} : $\mathbf{Y}_c \doteq \mathbf{Y}[\mathbf{I} - \frac{1}{p} \mathbf{1}\mathbf{1}^*]$, where $\mathbf{1} \in \mathbb{R}^p$ is a vector of all 1's. Then finding the top k principal components of \mathbf{Y}_c : $\max_{\mathbf{W} \in \text{St}(k, n; \mathbb{R})} \frac{1}{2} \|\mathbf{W}^* \mathbf{Y}_c\|_2^2$ is to find k *uncorrelated* projections of $\mathbf{y} \in \mathbb{R}^n$ that has the top k sample variance, i.e. 2^{nd} order moment, Jolliffe (2011); Helwig (2017). Similar to PCA, the ℓ^4 -norm maximization of centered data matrix \mathbf{Y}_c : $\max_{\mathbf{W} \in \text{St}(k, n; \mathbb{R})} \frac{1}{4} \|\mathbf{W}^* \mathbf{Y}_c\|_4^4$ can be viewed as finding k *uncorrelated* projections of \mathbf{y} that have the top k sample 4^{th} order moment, whose statistical meaning is better revealed below.

Relation with ICA and Nonnormality. The ℓ^4 -norm maximization over the Stiefel manifold is strongly related to finding the maximal or minimal *kurtosis* in *Independent Component Analysis* (ICA) Hyvärinen & Oja (1997; 2000): In order to identify one component of a given random vector $\mathbf{y} \in \mathbb{R}^n$, ICA aims at finding a unit vector (a direction) $\mathbf{w} \in \mathbb{S}^{n-1}$ that maximizes or minimizes the kurtosis of $\mathbf{w}^* \mathbf{y}$, defined as:

$$\text{kurt}(\mathbf{w}^* \mathbf{y}) = \mathbb{E}(\mathbf{w}^* \mathbf{y})^4 - 3 \|\mathbf{w}\|_2^4.$$

Kurtosis is widely used for evaluating the *nonnormality* of a random variable, see DeCarlo (1997); Hyvärinen & Oja (1997; 2000). According to Huber (1985), the nonnormality of data carries “abnormal” hence interesting information in real data for many applications, e.g. Lee et al. (2003); Cain et al. (2017). Hence, extractin the 4^{th} order moment helps understand such statistics of real datasets Hyvärinen et al. (2009) and even their topology Carlsson (2009).

2.2 FIXED-POINT STYLE ALGORITHMS

In optimization, the ℓ^4 -norm maximization in equation 6 over the Stiefel manifold $\text{St}(k, n; \mathbb{R})$ is a special type of nonconvex optimization problem – convex maximization over a compact set. Although the work of Journée et al. (2010); Zhai et al. (2019b) have shown that the MSP algorithm is guaranteed to find critical points, the experiments in Zhai et al. (2019b) suggest that the MSP algorithm finds global maxima of the ℓ^4 -norm efficiently and effectively. To understand this phenomenon better, in this section we illustrate some striking similarities between the MSP algorithm and the “power-iteration” type algorithms for solving PCA as well as ICA.

Fixed-point Perspective of Power Iteration. For a general data matrix $\mathbf{Y} \in \mathbb{R}^{n \times p}$, finding the top singular value of \mathbf{Y} is equivalent to solving the following optimization problem:

$$\max_{\mathbf{w}} \varphi(\mathbf{w}) \doteq \frac{1}{2} \|\mathbf{w}^* \mathbf{Y}\|_2^2 \quad \text{subject to } \mathbf{w} \in \mathbb{S}^{n-1}. \quad (8)$$

	Objectives	Constraint Sets	Algorithms
Power Iteration	$\varphi(\mathbf{w}) \doteq \frac{1}{2} \ \mathbf{w}^* \mathbf{Y}\ _2^2$	$\mathbf{w} \in \mathbb{S}^{n-1}$	$\mathbf{w}_{t+1} = \mathcal{P}_{\mathbb{S}^{n-1}} [\nabla_{\mathbf{w}} \varphi(\mathbf{w}_t)]$
FastICA	$\psi(\mathbf{w}) \doteq \frac{1}{4} \text{kurt}[\mathbf{w}^* \mathbf{y}]$	$\mathbf{w} \in \mathbb{S}^{n-1}$	$\mathbf{w}_{t+1} = \mathcal{P}_{\mathbb{S}^{n-1}} [\nabla_{\mathbf{w}} \psi(\mathbf{w}_t)]$
MSP	$\phi(\mathbf{W}) \doteq \frac{1}{4} \ \mathbf{W}^* \mathbf{Y}\ _4^4$	$\mathbf{W} \in \text{St}(k, n; \mathbb{R})$	$\mathbf{W}_{t+1} = \mathcal{P}_{\text{St}(k, n; \mathbb{R})} [\nabla_{\mathbf{W}} \phi(\mathbf{W}_t)]$

Table 1: Similarities among fixed-point algorithms for Power Iteration, FastICA, and MSP.

For this constrained optimization, the Lagrangian multiplier method gives the necessary condition: $\nabla_{\mathbf{w}} \varphi(\mathbf{w}) = \mathbf{Y} \mathbf{Y}^* \mathbf{w} = \lambda \mathbf{w}$, similar to equation 4. If we restrict this condition onto the sphere, we obtain the fixed point condition $\mathbf{w} = \mathcal{P}_{\mathbb{S}^{n-1}} [\nabla_{\mathbf{w}} \varphi(\mathbf{w})]$. The classic power-iteration method

$$\mathbf{w}_{t+1} = \mathcal{P}_{\mathbb{S}^{n-1}} [\nabla_{\mathbf{w}} \varphi(\mathbf{w}_t)] = \frac{\mathbf{Y} \mathbf{Y}^* \mathbf{w}_t}{\|\mathbf{Y} \mathbf{Y}^* \mathbf{w}_t\|_2}, \quad (9)$$

is precisely to compute this fixed point, which is arguably the most efficient and widely used algorithm to solve equation 8, for PCA (or computing SVD of \mathbf{Y}).

Fixed-point Perspective of FastICA. In order to maximize (or minimize) the kurtosis over \mathbb{S}^{n-1} ,

$$\max_{\mathbf{w}} \psi(\mathbf{w}) \doteq \frac{1}{4} \text{kurt}[\mathbf{w}^* \mathbf{y}] = \frac{1}{4} \mathbb{E}[\mathbf{w}^* \mathbf{y}]^4 - \frac{3}{4} \|\mathbf{w}\|_2^4 \quad \text{subject to } \mathbf{w} \in \mathbb{S}^{n-1}, \quad (10)$$

Hyvärinen & Oja (1997) has proposed the following fixed-point type iteration:

$$\mathbf{w}_{t+1} = \mathcal{P}_{\mathbb{S}^{n-1}} [\nabla_{\mathbf{w}} \psi(\mathbf{w}_t)] = \frac{\mathbb{E}[\mathbf{y}(\mathbf{y}^* \mathbf{w}_t)^3] - 3 \|\mathbf{w}_t\|_2^2 \mathbf{w}_t}{\|\mathbb{E}[\mathbf{y}(\mathbf{y}^* \mathbf{w}_t)^3] - 3 \|\mathbf{w}_t\|_2^2 \mathbf{w}_t\|_2}, \quad (11)$$

that enjoys cubic (at least quadratic) rate of convergence, under the ICA model assumption.

Fixed-point Perspective of MSP. For the ℓ^4 -norm maximization program:

$$\max_{\mathbf{W}} \phi(\mathbf{W}) \doteq \frac{1}{4} \|\mathbf{W}^* \mathbf{Y}\|_4^4 \quad \text{subject to } \mathbf{W} \in \text{St}(k, n; \mathbb{R}),$$

through a similar derivation to that in Section 1 one can show that the MSP iteration in equation 5 for the orthogonal group generalizes to the Stiefel manifold case as:

$$\mathbf{W}_{t+1} = \mathcal{P}_{\text{St}(k, n; \mathbb{R})} [\nabla_{\mathbf{W}} \phi(\mathbf{W}_t)] = \mathbf{U}_t \mathbf{V}_t^*, \quad (12)$$

where $\mathbf{U}_t \boldsymbol{\Sigma}_t \mathbf{V}_t^* = \text{SVD}[\mathbf{Y}(\mathbf{Y}^* \mathbf{W}_t)^{\circ 3}]$. The above iteration has the same nature as the power iteration in equation 9 and equation 11, since they all solve a fixed-point type problem, by projecting gradient of the objective function $\nabla \varphi(\cdot)$, $\nabla \psi(\cdot)$, $\nabla \phi(\cdot)$ onto the constraint manifold \mathbb{S}^{n-1} and $\text{St}(k, n; \mathbb{R})$, respectively. Table 1 summarizes such striking similarities.

3 STABILITY AND ROBUSTNESS OF ℓ^4 -NORM MAXIMIZATION

Even though the MSP algorithm for ℓ^4 -norm maximization is similar to power-iteration for PCA, in real applications, PCA often requires modification to improve its robustness Candès et al. (2011); Xu et al. (2010; 2012). In this section, we want to examine the stability and robustness of the ℓ^4 -maximization for different types of imperfect measurement models: small noise, outliers, and sparse corruptions of large magnitude.

3.1 DIFFERENT MODELS FOR IMPERFECT MEASUREMENTS

We adopt the same Bernoulli-Gaussian model as in prior works Spielman et al. (2012); Sun et al. (2015); Bai et al. (2018); Zhai et al. (2019b) to test the stability and robustness of the ℓ^4 -maximization framework. Assume our clean observation matrix $\mathbf{Y} \in \mathbb{R}^{n \times p}$ is produced by the product of a ground truth orthogonal dictionary \mathbf{D}_o and a Bernoulli-Gaussian matrix $\mathbf{X}_o \in \mathbb{R}^{n \times p}$:

$$\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o, \quad \mathbf{D}_o \in \text{O}(n; \mathbb{R}), \quad \{\mathbf{X}_o\}_{i,j} \sim_{iid} \text{BG}(\theta). \quad (13)$$

Now let us assume we only observe different types of imperfect measurements of \mathbf{Y} :

Noisy Measurements: $\mathbf{Y}_N := \mathbf{Y} + \mathbf{G}$, where $\mathbf{G} \in \mathbb{R}^{n \times p}$ is matrix that satisfies $g_{i,j} \sim_{iid} \mathcal{N}(0, \eta^2)$, where $\eta > 0$ controls the variance of the noise.

Measurements with Outliers: $\mathbf{Y}_O := [\mathbf{Y}, \mathbf{G}']$, where \mathbf{Y}_O contains extra columns ($\mathbf{G}' \in \mathbb{R}^{n \times \tau p}$)⁶ that is generated from an independent Gaussian process $g'_{i,j} \sim_{iid} \mathcal{N}(0, 1)$, and τ controls the portion of the outliers, w.r.t. the clean data size p .

⁶In case τp is not an integer, we round τp to its closest integer in our implementation.

Measurements with Sparse Corruptions: $\mathbf{Y}_C := \mathbf{Y} + \sigma \mathbf{B} \circ \mathbf{S}$, where $\sigma > 0$ controls the scale of corrupting entries,⁷ $\mathbf{B} \in \mathbb{R}^{n \times p}$ is a Bernoulli matrix with $b_{i,j} \sim_{iid} \text{Ber}(\beta)$, where $\beta \in (0, 1)$ controls the ratio of the sparse corruptions, and entries $s_{i,j}$ of $\mathbf{S} \in \mathbb{R}^{n \times p}$ are i.i.d. drawn from a *Rademacher* distribution:

$$s_{i,j} = \begin{cases} 1 & \text{with probability } 1/2 \\ -1 & \text{with probability } 1/2 \end{cases}. \quad (14)$$

3.2 STATISTICAL ANALYSIS AND JUSTIFICATION

The analysis for the stability and robustness of the ℓ^4 -norm maximization follows similar statistical analysis techniques in Zhai et al. (2019b) that establish the global maximum of

$$\mathbf{W}_\star \in \arg \max_{\mathbf{W}} \mathbb{E} \|\mathbf{W}^* \mathbf{Y}_\diamond\|_4^4, \quad \text{subject to } \mathbf{W} \in \mathcal{O}(n; \mathbb{R}) \quad (15)$$

satisfying $\mathbf{W}_\star^* \mathbf{D}_o \in \text{SP}(n)$.⁸ We use \mathbf{Y}_\diamond here to denote different imperfect measurements (noisy \mathbf{Y}_N , with outliers \mathbf{Y}_O , and with sparse corruptions \mathbf{Y}_C , respectively). Below we calculate the expectation of $\|\mathbf{W}^* \mathbf{Y}_\diamond\|_4^4$ over the data distribution, since in general the objective function concentrates on its expectation.⁹ We show that $\mathbb{E} \|\mathbf{W}^* \mathbf{Y}_\diamond\|_4^4$ is largely determined by $\|\mathbf{W} \mathbf{D}_o\|_4^4$, a quantity that indicates a ‘‘distance’’ between $\mathbf{W}^* \mathbf{D}_o$ to $\text{SP}(n)$. As shown in Lemma 2.3 and Lemma 2.4 in Zhai et al. (2019b), the *only global maximizers* of $\|\mathbf{W}^* \mathbf{D}_o\|_4^4$ are signed permutation matrices, and $\mathbf{W}^* \mathbf{D}_o$ converges to a signed permutation matrix as $\|\mathbf{W}^* \mathbf{D}_o\|_4^4$ reaches its global maximum.

Proposition 3.1 (Expectation of Objective with Small Noise) $\forall \theta \in (0, 1)$, let $\mathbf{X}_o \in \mathbb{R}^{n \times p}$, $x_{i,j} \sim_{iid} \text{BG}(\theta)$, $\mathbf{D}_o \in \mathcal{O}(n; \mathbb{R})$ is any orthogonal matrix, and $\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o$. For any orthogonal matrix $\mathbf{W} \in \mathcal{O}(n; \mathbb{R})$ and any random Gaussian matrix $\mathbf{G} \in \mathbb{R}^{n \times p}$, $g_{i,j} \sim_{iid} \mathcal{N}(0, \eta^2)$ independent of \mathbf{X}_o , let $\mathbf{Y}_N = \mathbf{Y} + \mathbf{G}$ denote the data with noise. Then the expectation of $\|\mathbf{W}^* \mathbf{Y}_N\|_4^4$ is:

$$\frac{1}{np} \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \|\mathbf{W}^* \mathbf{Y}_N\|_4^4 = 3\theta(1 - \theta) \frac{\|\mathbf{W}^* \mathbf{D}_o\|_4^4}{n} + C_{\theta, \eta}, \quad (16)$$

where $C_{\theta, \eta}$ is a constant depending on θ and η .

Proof See Appendix A.1. ■

Proposition 3.2 (Expectation of Objective with Outliers) $\forall \theta \in (0, 1)$, let $\mathbf{X}_o \in \mathbb{R}^{n \times p}$, $x_{i,j} \sim_{iid} \text{BG}(\theta)$, $\mathbf{D}_o \in \mathcal{O}(n; \mathbb{R})$ is any orthogonal matrix and $\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o$. For any orthogonal matrix $\mathbf{W} \in \mathcal{O}(n; \mathbb{R})$ and any random Gaussian matrix $\mathbf{G}' \in \mathbb{R}^{n \times \tau p}$, $g'_{i,j} \sim_{iid} \mathcal{N}(0, 1)$ independent of \mathbf{X}_o , let $\mathbf{Y}_O = [\mathbf{Y}, \mathbf{G}']$ denote the data with outliers \mathbf{G}' . Then the expectation of $\|\mathbf{W}^* \mathbf{Y}_O\|_4^4$ is:

$$\frac{1}{np} \mathbb{E}_{\mathbf{X}_o, \mathbf{G}'} \|\mathbf{W}^* \mathbf{Y}_O\|_4^4 = 3\theta(1 - \theta) \frac{\|\mathbf{W}^* \mathbf{D}_o\|_4^4}{n} + C_\theta, \quad (17)$$

where C_θ is a constant depending on θ .

Proof See Appendix A.2. ■

In the above results, Proposition 3.1, Proposition 3.2 reveal that both normalized $\frac{1}{np} \mathbb{E} \|\mathbf{W}^* \mathbf{Y}_N\|_4^4$, $\frac{1}{np} \mathbb{E} \|\mathbf{W}^* \mathbf{Y}_O\|_4^4$ are only determined by $\|\mathbf{W}^* \mathbf{D}_o\|_4^4$, therefore, the ℓ^4 -norm maximization formulation is stable to $\|\mathbf{W}^* \mathbf{Y}_\diamond\|_4^4$ with dense Gaussian noise and Gaussian outliers – maximizing the ℓ^4 -norm of $\mathbb{E} \|\mathbf{W}^* \mathbf{Y}_\diamond\|_4^4$ will recover the ground truth orthogonal matrix \mathbf{D}_o .

Proposition 3.3 (Expectation of Objective with Sparse Corruptions) $\forall \theta \in (0, 1)$, let $\mathbf{X}_o \in \mathbb{R}^{n \times p}$, $x_{i,j} \sim_{iid} \text{BG}(\theta)$, $\mathbf{D}_o \in \mathcal{O}(n; \mathbb{R})$ is any orthogonal matrix and $\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o$. For any orthogonal matrix $\mathbf{W} \in \mathcal{O}(n; \mathbb{R})$ and any random Bernoulli matrix $\mathbf{B} \in \mathbb{R}^{n \times p}$, $b_{i,j} \sim_{iid} \text{Ber}(\beta)$ independent of \mathbf{X}_o , let $\mathbf{Y}_C = \mathbf{Y} + \sigma \mathbf{B} \circ \mathbf{S}$ denote the data with sparse corruptions, and $\mathbf{S} \in \mathbb{R}^{n \times p}$ is defined as equation 14. Then the expectation of $\|\mathbf{W}^* \mathbf{Y}_C\|_4^4$ is:

$$\frac{1}{np} \mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \|\mathbf{W}^* \mathbf{Y}_C\|_4^4 = 3\theta(1 - \theta) \frac{\|\mathbf{W}^* \mathbf{D}_o\|_4^4}{n} + \sigma^4 \beta(1 - 3\beta) \frac{\|\mathbf{W}\|_4^4}{n} + C_{\theta, \sigma, \beta}, \quad (18)$$

where $C_{\theta, \sigma, \beta}$ is a constant depending on θ , σ and β .

Proof See Appendix A.3. ■

⁷In our context, $\sigma = 1$ is already corruption of large magnitude, since the variance of the sparse signal is 1.

⁸ $\text{SP}(n)$ is the signed permutation group, a group of orthogonal matrices that only contain 0, ± 1 .

⁹One can prove concentration bounds similar to that of Lemma 2.2 in Zhai et al. (2019b).

Unlike the cases with noise and outlier, Proposition 3.3 indicates $\frac{1}{np} \mathbb{E} \|\mathbf{W}^* \mathbf{Y}_C\|_4^4$ depends on not only $\|\mathbf{W}^* \mathbf{D}_o\|_4^4$ but also $\|\mathbf{W}\|_4^4$. Nevertheless, when the magnitude of $\sigma^4 \beta(1 - 3\beta)$ is significantly smaller than $3\theta(1 - \theta)$, the landscape of the objective $\|\mathbf{W}^* \mathbf{Y}_C\|_4^4$ would largely be determined by $\|\mathbf{W}^* \mathbf{D}_o\|_4^4$ only. As shown in Figure 1, this is indeed the case whenever: a) the sparsity level θ of ground truth signal \mathbf{X}_o , is “reasonably” small (neither diminishing to 0 nor larger than 0.5); b) β , the sparsity level of the corruption, is small (smaller than 0.5); c) σ , the magnitude of the sparse errors, is not significantly larger than the intrinsic variance of the sparse signal (the intrinsic variance of the sparse signal Bernoulli-Gaussian model is 1).

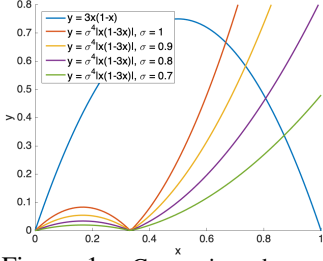


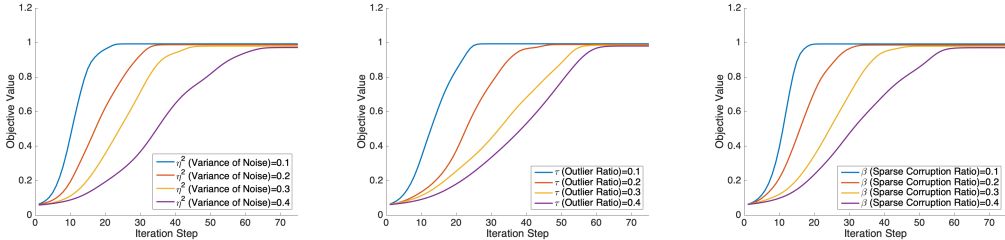
Figure 1: Comparison between $y = 3x(1 - x)$ and $y = \sigma^4|x(1 - 3x)|$ when $x \in [0, 1]$ with different σ .

Hence, the above analysis shows that the dictionary estimated from ℓ^4 -maximization should be insensitive to small noise, robust to fraction of outliers, and resilient to sparse corruptions.

4 SIMULATIONS AND EXPERIMENTS

4.1 QUANTITATIVE EVALUATION: SIMULATIONS ON SYNTHETIC DATA

Single Trial of MSP. In this simulation, we run the MSP Algorithm from equation 3, using the imperfect measurements \mathbf{Y}_o of different types ($\mathbf{Y}_N, \mathbf{Y}_O, \mathbf{Y}_C$). As shown in Figure 2, the normalized value of $\|\mathbf{W}^* \mathbf{D}_o\|_4^4/n$ reaches global maximum with all types of inputs when varying the level of noise, outliers, and sparse corruptions. Moreover, as the level of noise increases, Figure 2 shows that a) the iterations for convergence increases and b) the final objective value $\|\mathbf{W}^* \mathbf{D}_o\|_4^4$ decreases almost negligibly. This numerical experiment suggests that the MSP Algorithm is able to identify the ground truth orthogonal transformation \mathbf{D}_o despite different types of imperfect measurement.



(a) $n = 50, p = 20,000, \theta = 0.3$, (b) $n = 50, p = 20,000, \theta = 0.3$, (c) $n = 50, p = 20,000, \theta = 0.3$, varying η^2 from 0.1 to 0.4, varying τ from 0.1 to 0.4, $\sigma = 1$, varying β from 0.1 to 0.4

Figure 2: Normalized $\|\mathbf{W}^* \mathbf{D}_o\|_4^4/n$ of the MSP algorithm for dictionary learning, using imperfect measurements $\mathbf{Y}_N, \mathbf{Y}_O, \mathbf{Y}_C$, respectively.

Phase Transition. Next, we conduct extensive simulations to study the relation between recovery accuracy and sample size p . We run the experiments by increasing the sample size p w.r.t. the levels of noises and corruptions η, τ, β , respectively. As shown in Figure 3, the MSP Algorithm 3 demonstrates a clear phase transition behavior w.r.t. noise, outliers, and sparse corruptions. Such phenomena suggest that the algorithm is inherently stable and robust to certain amounts of noise, outliers, and sparse corruptions. The results also indicate that a larger sample size p increases the accuracy and robustness of the MSP Algorithm 3 for all types of nuisances.

4.2 QUALITATIVE EVALUATION: EXPERIMENTS ON REAL IMAGES AND PATCHES

Besides simulations, we also conduct extensive experiments to verify stability and robustness of the MSP Algorithm with real imagery data, at both image level and patch level. Throughout these experiments, rather than visualize all bases, we routinely show the top bases learned – heuristically, top bases are those with the largest coefficients (here, in terms of ℓ^1 -norm).

Image Level. At image level, we first vectorize all 60,000 images in the MNIST dataset (LeCun et al., 1998) into a single matrix $\mathbf{Y} \in \mathbb{R}^{784 \times 60,000}$, then create imperfect measurements based on models specified in Section 3: \mathbf{Y}_N (MNIST with noise), \mathbf{Y}_O (MNIST with outliers), \mathbf{Y}_C (MNIST

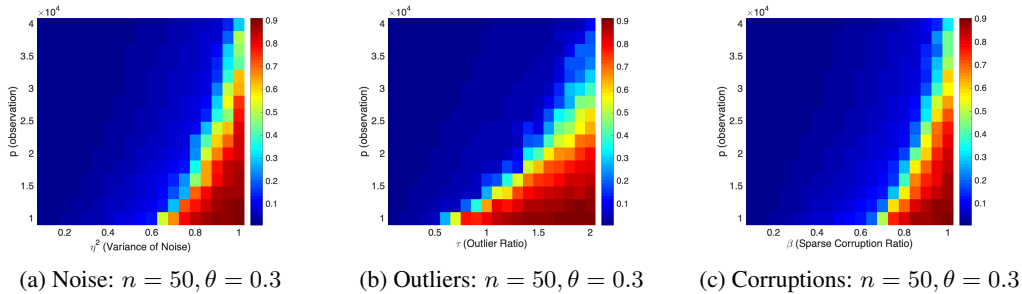


Figure 3: Average normalized error $|1 - \|\mathbf{W}^* \mathbf{D}_o\|_4^4 / n|$ of 10 random trials for the MSP Algorithm: (a) Varying sample size p and variance of noise η^2 ; (b) Varying sample size p and Gaussian Outlier ratio τ ; (c) Varying sample size p and sparse corruption ratio β , with fixed $\sigma = 1$.

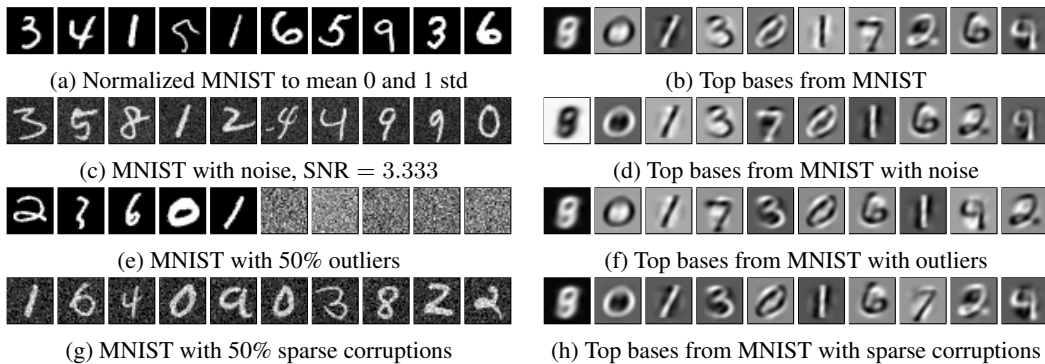


Figure 4: **Left:** Examples of MNIST and its different imperfect measurements. **Right:** Learned bases from MNIST and its different imperfect measurements using the MSP Algorithm 3.

with sparse corruptions). We run the MSP Algorithm 3 with $\mathbf{Y}, \mathbf{Y}_N, \mathbf{Y}_C, \mathbf{Y}_O$ and compare the bases learned. Figure 4(a), (c), (e), and (g) show examples of $\mathbf{Y}, \mathbf{Y}_N, \mathbf{Y}_O,$ and $\mathbf{Y}_C,$ and Figure 4(b), (d), (f), and (h) show top 10 bases learned from $\mathbf{Y}, \mathbf{Y}_N, \mathbf{Y}_C, \mathbf{Y}_O,$ respectively. Despite that we use different types of imperfect measurements of MNIST, the top bases learned from MSP Algorithm 3 are very much the same.¹⁰ This result corroborates with our analysis: the ℓ^4 -maximization and the MSP algorithm is inherently insensitive to noise, robust to outliers, and resilient to sparse corruptions.

Patch Level. A classic application of dictionary learning involves learning sparse representations of image patches (Elad & Aharon, 2006; Mairal et al., 2007). In this section, we extend the experiments of Zhai et al. (2019b) to learn patches from grayscale and color images. First, we construct a data matrix \mathbf{Y} by vectorizing each 8×8 patch from the 512×512 grayscale image, “Barbara” (see Figure 5). We then run the MSP algorithm with 100 iterations on both \mathbf{Y} and a noisy version $\mathbf{Y}_N,$ and the learned top bases are visualized in Figure 5.

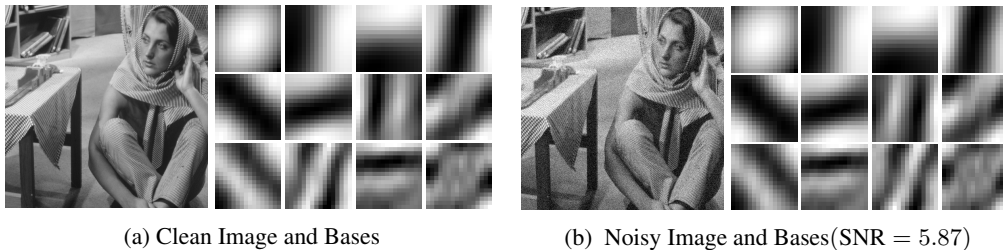


Figure 5: The top 12 bases learned from all 16×16 patches of Barbara, both with (b) and without (a) noise. The noisy image is produced by adding Gaussian noise to the clean image, resulting in a signal-to-noise (SNR) ratio of 5.87. We observed a similar effect when using an 8×8 patch size.

¹⁰Bases with opposite intensity are considered as the same base, since they only differ by a sign.

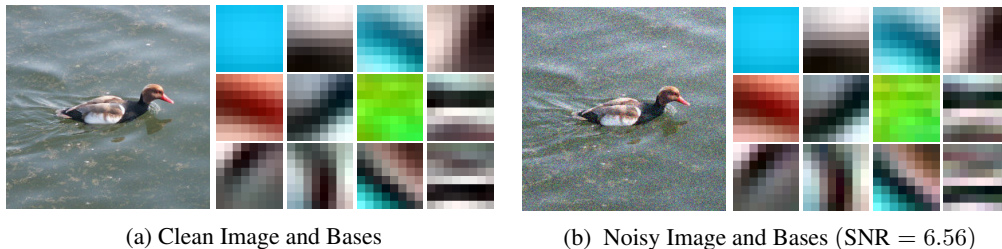


Figure 6: The top 12 bases learned from all $8 \times 8 \times 3$ color patches of the clean and noisy image, respectively. Here, the SNR of the noisy image is 6.56.

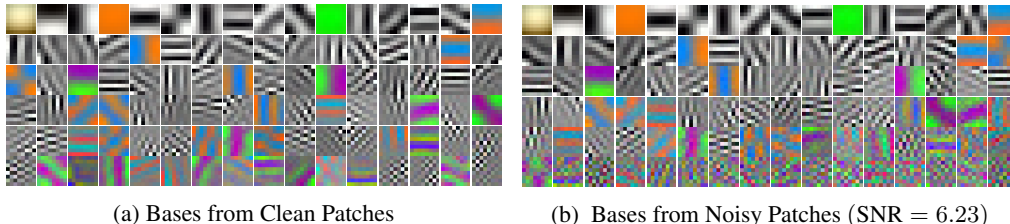


Figure 7: Top half (96) bases learned from 100,000 random $8 \times 8 \times 3$ patches sampled from CIFAR-10, before and after adding Gaussian noise, with SNR 6.23.

Analogously, we apply the same scheme to a 256×256 color image, “Duck” (see Figure 6), converting each $8 \times 8 \times 3$ patch into a column vector (in \mathbb{R}^{192}) of \mathbf{Y} . Notice this forces the algorithm to learn bases for all three channels, rather than one at a time. After running the MSP algorithm for 100 iterations, we visualize the top bases learned from both \mathbf{Y} and corresponding \mathbf{Y}_N in Figure 6.

We next consider the problem of learning a “global dictionary” (Mairal et al., 2007) for patches from many different images. To construct our data matrix, \mathbf{Y} , we randomly sample 100,000 $8 \times 8 \times 3$ patches from the CIFAR-10 data-set (Krizhevsky et al., 2009). A noisy, \mathbf{Y}_N , is then generated by adding Gaussian noise. Again, we apply the MSP algorithm with 200 iterations to learn 192 bases and visualize the results in Figure 7. We leave the experiments of CIFAR-10 with outliers and sparse corruptions in the Appendix due to limited space.

In each of these experiments, the top bases in the learned dictionary remain relatively unchanged with the addition of noise. To quantify this similarity, we take the top bases from the noisy dictionary and find the closest top clean base for each. If the bases are nearly identical, then the inner product of each of these pairs should be close to 1. Table 2 reports the statistics.

	Minimum	Lower Quartile	Median	Upper Quartile	Maximum
Barbara	0.3048	0.8471	0.9941	0.9993	1.0000
Duck	0.2510	0.9782	0.9891	0.9971	1.0000
CIFAR-10	0.5147	0.7203	0.9892	0.9998	1.0000

Table 2: Statistics about the inner products between the top 20 noisy bases and their corresponding closest top-20 clean bases.

5 CONCLUSION AND DISCUSSIONS

In this paper, we find the ℓ^4 -norm maximization based dictionary learning and the MSP algorithm introduced by Zhai et al. (2019b) have strong geometric and statistical connections to classic data analysis methods PCA and ICA. Such connections seem to be the reasons why they all admit similarly simple and efficient algorithms.

Empirically, we have observed that ℓ^4 -norm maximization is surprisingly insensitive to noise, robust to outliers, and resilient to sparse corruptions. Our preliminary analysis supports such phenomena but also suggests that the formulation could still be improved in the case of sparse corruptions.

From experiments on real images, we observed that top bases learned are rather stable but tail bases can be less stable (see Figure 10 in Appendix B). This is largely due to the fact that real images do not follow the uniformly sparse Bernoulli Gaussian model (of equation 13). Generalizing dictionary learning to non-uniformly sparsely generated data would be a good topic for future study.

REFERENCES

- P-A Absil and Jérôme Malick. Projection-like retractions on matrix manifolds. *SIAM Journal on Optimization*, 22(1):135–158, 2012.
- Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. Convex multi-task feature learning. *Machine Learning*, 73(3):243–272, 2008.
- Yu Bai, Qijia Jiang, and Ju Sun. Subgradient descent learns orthogonal dictionaries. *arXiv preprint arXiv:1810.10702*, 2018.
- Boaz Barak, Jonathan Kelner, and David Steurer. Dictionary learning and tensor decomposition via the sum-of-squares method. In *STOC*, 2015.
- Meghan K Cain, Zhiyong Zhang, and Ke-Hai Yuan. Univariate and multivariate skewness and kurtosis for measuring nonnormality: Prevalence, influence and estimation. *Behavior Research Methods*, 49(5):1716–1735, 2017.
- Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3):11, 2011.
- Gunnar Carlsson. Topology and data. *Bulletin of the American Mathematical Society*, 46(2):255–308, 2009.
- Niladri Chatterji and Peter L Bartlett. Alternating minimization for dictionary learning with random initialization. In *Advances in Neural Information Processing Systems*, pp. 1997–2006, 2017.
- Lawrence T DeCarlo. On the meaning and use of kurtosis. *Psychological methods*, 2(3):292, 1997.
- Michael Elad and Michal Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image processing*, 15(12):3736–3745, 2006.
- Dar Gilboa, Sam Buchanan, and John Wright. Efficient dictionary learning with gradient descent. *arXiv preprint arXiv:1809.10313*, 2018.
- Nathaniel E Helwig. Principal components analysis. 2017.
- Peter J Huber. Projection pursuit. *The annals of Statistics*, pp. 435–475, 1985.
- Aapo Hyvärinen and Erkki Oja. A fast fixed-point algorithm for independent component analysis. *Neural Computation*, 9:1483–1492, 1997.
- Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural networks*, 13(4-5):411–430, 2000.
- Aapo Hyvärinen, Jarmo Hurri, and Patrick O Hoyer. *Natural image statistics: A probabilistic approach to early computational vision.*, volume 39. Springer Science & Business Media, 2009.
- Ian Jolliffe. *Principal component analysis*. Springer, 2011.
- Michel Journée, Yurii Nesterov, Peter Richtárik, and Rodolphe Sepulchre. Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research*, 11 (Feb):517–553, 2010.
- Alex Krizhevsky et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Han-Wen Kuo, Yenson Lau, Yuqian Zhang, and John Wright. Geometry and symmetry in short-and-sparse deconvolution. *arXiv preprint arXiv:1901.00256*, 2019.
- Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Ann B Lee, Kim S Pedersen, and David Mumford. The nonlinear statistics of high-contrast patches in natural images. *International Journal of Computer Vision*, 54(1-3):83–103, 2003.

- Tengyu Ma, Jonathan Shi, and David Steurer. Polynomial-time tensor decompositions with sum-of-squares. In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 438–446. IEEE, 2016.
- Julien Mairal, Michael Elad, and Guillermo Sapiro. Sparse representation for color image restoration. *IEEE Transactions on image processing*, 17(1):53–69, 2007.
- Julien Mairal, Francis Bach, Jean Ponce, Guillermo Sapiro, and Andrew Zisserman. Discriminative learned dictionaries for local image analysis. Technical report, MINNESOTA UNIV MINNEAPOLIS INST FOR MATHEMATICS AND ITS APPLICATIONS, 2008.
- Julien Mairal, Francis Bach, and Jean Ponce. Task-driven dictionary learning. *IEEE transactions on pattern analysis and machine intelligence*, 34(4):791–804, 2012.
- Julien Mairal, Francis Bach, and Jean Ponce. Sparse modeling for image and vision processing. *Foundations and Trends in Computer Graphics and Vision*, 8(2-3):85–283, 2014. ISSN 1572-2740. doi: 10.1561/06000000058.
- Arthur Mensch, Julien Mairal, Bertrand Thirion, and Gaël Varoquaux. Dictionary learning for massive matrix factorization. In *International Conference on Machine Learning*, pp. 1737–1746, 2016.
- Thanh Nguyen, Akshay Soni, and Chinmay Hegde. On learning sparsely used dictionaries from incomplete samples. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 3769–3778, Stockholm, Sweden, 10–15 Jul 2018. PMLR.
- Bruno A. Olshausen and David J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.
- Bruno A. Olshausen and David J. Field. Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision research*, 37(23):3311–3325, 1997.
- Sirisha Rambhatla, Xingguo Li, and Jarvis Haupt. Noodl: Provable online dictionary learning and sparse coding. In *International Conference on Learning Representations*, 2019.
- Marc’Aurelio Ranzato, Y-Lan Boureau, and Yann LeCun. Sparse feature learning for deep belief networks. In *Advances in Neural Information Processing Systems (NIPS 2007)*, volume 20, 2007.
- Tselil Schramm and David Steurer. Fast and robust tensor decomposition with applications to dictionary learning. *arXiv preprint arXiv:1706.08672*, 2017.
- Daniel A Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. In *Conference on Learning Theory*, pp. 37–1, 2012.
- Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere. *arXiv preprint arXiv:1504.06785*, 2015.
- John Wright, Allen Y Yang, Arvind Ganesh, S Shankar Sastry, and Yi Ma. Robust face recognition via sparse representation. *IEEE transactions on pattern analysis and machine intelligence*, 31(2): 210–227, 2008.
- Huan Xu, Constantine Caramanis, and Sujay Sanghavi. Robust PCA via outlier pursuit. In *Advances in Neural Information Processing Systems*, pp. 2496–2504, 2010.
- Huan Xu, Constantine Caramanis, and Shie Mannor. Outlier-robust PCA: The high-dimensional case. *IEEE transactions on information theory*, 59(1):546–572, 2012.
- Jianchao Yang, John Wright, Thomas S Huang, and Yi Ma. Image super-resolution via sparse representation. *IEEE transactions on image processing*, 19(11):2861–2873, 2010.
- Yuxiang Zhai, Zitong Yang, Zhenyu Liao, John Wright, and Yi Ma. A fast holistic algorithm for complete dictionary learning via ℓ^4 norm maximization. *Signal Processing with Adaptive Sparse Structured Representations (SPARS)*, 2019a.

Yuexiang Zhai, Zitong Yang, Zhenyu Liao, John Wright, and Yi Ma. Complete dictionary learning via ℓ^4 -norm maximization over the orthogonal group. *arXiv preprint arXiv:1906.02435*, 2019b.

Yuqian Zhang, Han-Wen Kuo, and John Wright. Structured local optima in sparse blind deconvolution. *arXiv preprint arXiv:1806.00338*, 2018.

A PROOF OF SECTION 3

A.1 PROOF OF PROPOSITION 3.1

Claim A.1 (Expectation of Objective with Small Noise) $\forall \theta \in (0, 1)$, let $\mathbf{X}_o \in \mathbb{R}^{n \times p}$, $x_{i,j} \sim_{iid} BG(\theta)$, $\mathbf{D}_o \in \mathcal{O}(n; \mathbb{R})$ is any orthogonal matrix and $\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o$. For any orthogonal matrix $\mathbf{W} \in \mathcal{O}(n; \mathbb{R})$ and any random Gaussian matrix $\mathbf{G} \in \mathbb{R}^{n \times p}$, $g_{i,j} \sim_{iid} \mathcal{N}(0, \eta^2)$ independent of \mathbf{X}_o , let $\mathbf{Y}_N = \mathbf{Y} + \mathbf{G}$ denote the input with noise. The expectation of $\|\mathbf{W}^* \mathbf{Y}_N\|_4^4$ satisfies the following property:

$$\frac{1}{np} \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \|\mathbf{W}^* \mathbf{Y}_N\|_4^4 = 3\theta(1-\theta) \frac{\|\mathbf{W}^* \mathbf{D}_o\|_4^4}{n} + 3\theta^2 + 6\theta\eta^2 + 3\eta^4. \quad (19)$$

Proof Let $\mathbf{W}^* \mathbf{D}_o = \mathbf{M} \in \mathcal{O}(n; \mathbb{R})$, notice that the orthogonal transformation ($\mathbf{W}^* \mathbf{G}$) of a Gaussian matrix (\mathbf{G}) is still a Gaussian matrix and satisfies $\{\mathbf{W}^* \mathbf{G}\}_{i,j} \sim \mathcal{N}(0, 1)$, and it is independent of $\mathbf{Y}(\mathbf{X}_o)$. We abuse the notation a bit let $\mathbf{G} = \mathbf{W}^* \mathbf{G}$ in the following calculation, since they are independent in calculation.

$$\begin{aligned} & \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \|\mathbf{W}^* \mathbf{Y}_N\|_4^4 = \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \|\mathbf{M} \mathbf{X}_o + \mathbf{G}\|_4^4 \\ &= \sum_{j=1}^p \sum_{i=1}^n \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \left(\sum_{k=1}^n m_{i,k} x_{k,j} + g_{i,j} \right)^4 \\ &= \sum_{j=1}^p \sum_{i=1}^n \left\{ \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)^4 + 6 \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \left[n_{i,j}^2 \left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)^2 \right] + \mathbb{E}_{\mathbf{G}} g_{i,j}^4 \right\} \\ &= \sum_{j=1}^p \sum_{i=1}^n \left\{ \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)^4 + 6\eta^2 \mathbb{E}_{\mathbf{X}_o} \left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)^2 \right\} + 3np\eta^4 \\ &= \sum_{j=1}^p \sum_{i=1}^n \left\{ \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)^4 \right\} + 6np\theta\eta^2 + 3np\eta^4 \\ &= \mathbb{E}_{\mathbf{X}_o} \|\mathbf{M} \mathbf{X}_o\|_4^4 + 6np\theta\eta^2 + 3np\eta^4 \\ &= 3p\theta(1-\theta) \|\mathbf{M}\|_4^4 + 3np\theta^2 + 6np\theta\eta^2 + 3np\eta^4, \end{aligned} \quad (20)$$

therefore,

$$\frac{1}{np} \mathbb{E}_{\mathbf{X}_o, \mathbf{G}} \|\mathbf{W}^* \mathbf{Y}_N\|_4^4 = 3\theta(1-\theta) \frac{\|\mathbf{W}^* \mathbf{D}_o\|_4^4}{n} + 3\theta^2 + 6\theta\eta^2 + 3\eta^4, \quad (21)$$

which completes the proof. \blacksquare

A.2 PROOF OF PROPOSITION 3.2

Claim A.2 (Expectation of Objective with Outliers) $\forall \theta \in (0, 1)$, let $\mathbf{X}_o \in \mathbb{R}^{n \times p}$, $x_{i,j} \sim_{iid} BG(\theta)$, $\mathbf{D}_o \in \mathcal{O}(n; \mathbb{R})$ is any orthogonal matrix and $\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o$. For any orthogonal matrix $\mathbf{W} \in \mathcal{O}(n; \mathbb{R})$ and any random Gaussian matrix $\mathbf{G}' \in \mathbb{R}^{n \times \tau p}$, $g'_{i,j} \sim_{iid} \mathcal{N}(0, 1)$ independent of \mathbf{X}_o , let $\mathbf{Y}_O = [\mathbf{Y}, \mathbf{G}']$ denote the input with outlier \mathbf{G}' . The expectation of $\|\mathbf{W}^* \mathbf{Y}_O\|_4^4$ satisfies the following property:

$$\frac{1}{np} \mathbb{E}_{\mathbf{X}_o, \mathbf{G}'} \|\mathbf{W}^* \mathbf{Y}_O\|_4^4 = 3\theta(1-\theta) \frac{\|\mathbf{W}^* \mathbf{D}_o\|_4^4}{n} + 3\theta^2 + 3. \quad (22)$$

Proof Notice that

$$\mathbb{E}_{\mathbf{X}_o, \mathbf{G}'} \|\mathbf{W}^* \mathbf{Y}_O\|_4^4 = \mathbb{E}_{\mathbf{X}_o} \|\mathbf{W}^* \mathbf{Y}\|_4^4 + \mathbb{E}_{\mathbf{G}'} \|\mathbf{W}^* \mathbf{G}'\|_4^4, \quad (23)$$

and

$$\mathbb{E}_{\mathbf{X}_o} \|\mathbf{W}^* \mathbf{Y}\|_4^4 = \mathbb{E}_{\mathbf{X}_o} \|\mathbf{W}^* \mathbf{D}_o \mathbf{X}_o\|_4^4 = 3p\theta(1-\theta) \|\mathbf{W}^* \mathbf{D}_o\|_4^4 + 3np\theta^2. \quad (24)$$

Moreover, the orthogonal rotation ($\mathbf{W}^* \mathbf{G}'$) of a standard Gaussian matrix \mathbf{G}' is also a standard Gaussian matrix, therefore,

$$\mathbb{E}_{\mathbf{G}'} \|\mathbf{W}^* \mathbf{G}'\|_4^4 = 3np. \quad (25)$$

Hence,

$$\frac{1}{np} \mathbb{E}_{\mathbf{X}_o, \mathbf{G}'} \|\mathbf{W}^* \mathbf{Y}_O\|_4^4 = 3\theta(1-\theta) \frac{\|\mathbf{W}^* \mathbf{D}_o\|_4^4}{n} + 3\theta^2 + 3, \quad (26)$$

which completes the proof. \blacksquare

A.3 PROOF OF PROPOSITION 3.3

Claim A.3 (Expectation of Objective with Sparse Corruptions) $\forall \theta \in (0, 1)$, let $\mathbf{X}_o \in \mathbb{R}^{n \times p}$, $x_{i,j} \sim_{iid} BG(\theta)$, $\mathbf{D}_o \in O(n; \mathbb{R})$ is any orthogonal matrix and $\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o$. For any orthogonal matrix $\mathbf{W}^* \in O(n; \mathbb{R})$ and any random Bernoulli matrix $\mathbf{B} \in \mathbb{R}^{n \times p}$, $b_{i,j} \sim_{iid} Ber(\beta)$ independent of \mathbf{X}_o , let $\mathbf{Y}_C = \mathbf{Y} + \sigma \mathbf{B} \circ \mathbf{S}$ denote the input with sparse corruptions, and $\mathbf{S} \in \mathbb{R}^{n \times p}$ is defined as equation 14. The expectation of $\|\mathbf{W}^* \mathbf{Y}_C\|_4^4$ satisfies the following property:

$$\frac{1}{np} \mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \|\mathbf{W}^* \mathbf{Y}_C\|_4^4 = 3\theta(1-\theta) \frac{\|\mathbf{W}^* \mathbf{D}_o\|_4^4}{n} + \sigma^4 \beta(1-3\beta) \frac{\|\mathbf{W}\|_4^4}{n} + 3\theta^2 + 6\sigma^2 \theta \beta + 3\sigma^4 \beta^2 \quad (27)$$

Proof Let $\mathbf{W}^* \mathbf{D}_o = \mathbf{M} \in O(n; \mathbb{R})$, notice that

$$\|\mathbf{W}^* \mathbf{Y}_C\|_4^4 = \|\mathbf{M} \mathbf{X}_o + \sigma \mathbf{W}^* (\mathbf{B} \circ \mathbf{S})\|_4^4, \quad (28)$$

hence

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \|\mathbf{W}^* \mathbf{Y}_C\|_4^4 &= \mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \sum_{j=1}^p \sum_{i=1}^n \left(\sum_{k=1}^n m_{i,k} x_{k,j} + \sigma \sum_{k=1}^n w_{k,i} b_{k,j} s_{k,j} \right)^4 \\ &= \sum_{j=1}^p \sum_{i=1}^n \mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \left[\underbrace{\left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)^4}_{\Gamma_1} + 6 \underbrace{\left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)^2}_{\Gamma_2} \underbrace{\left(\sigma \sum_{k=1}^n w_{k,i} b_{k,j} s_{k,j} \right)^2}_{\Gamma_2} \right. \\ &\quad + \underbrace{\left(\sigma \sum_{k=1}^n w_{k,i} b_{k,j} s_{k,j} \right)^4}_{\Gamma_3} + 4 \underbrace{\left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)^3}_{\Gamma_4} \underbrace{\left(\sigma \sum_{k=1}^n w_{k,i} b_{k,j} s_{k,j} \right)}_{\Gamma_4} \\ &\quad \left. + 4 \underbrace{\left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)}_{\Gamma_5} \underbrace{\left(\sigma \sum_{k=1}^n w_{k,i} b_{k,j} s_{k,j} \right)^3}_{\Gamma_5} \right]. \end{aligned} \quad (29)$$

Moreover,

- Γ_1 :

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \Gamma_1 &= \mathbb{E}_{\mathbf{X}_o} \Gamma_1 = 3\theta \sum_{k=1}^n m_{i,k}^4 + 6\theta^2 \sum_{1 \leq k_1 < k_2 \leq n} m_{i,k_1}^2 m_{i,k_2}^2 \\ &= 3\theta(1-\theta) \sum_{k=1}^n m_{i,k}^4 + 3\theta^2, \end{aligned} \quad (30)$$

- Γ_2 :

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \Gamma_2 &= \left[\mathbb{E}_{\mathbf{X}_o} \left(\sum_{k=1}^n m_{i,k} x_{k,j} \right)^2 \right] \left[\mathbb{E}_{\mathbf{B}, \mathbf{S}} \left(\sigma \sum_{k=1}^n w_{k,i} b_{k,j} s_{k,j} \right)^2 \right] \\ &= \theta \left(\sum_{k=1}^n m_{i,k}^2 \right) \sigma^2 \beta \left(\sum_{k=1}^n w_{k,i}^2 \right) = \sigma^2 \theta \beta, \end{aligned} \quad (31)$$

- Γ_3 :

$$\begin{aligned}\mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \Gamma_3 &= \mathbb{E}_{\mathbf{B}, \mathbf{S}} \Gamma_3 = \sigma^4 \beta \sum_{k=1}^n w_{k,i}^4 + 6\sigma^4 \beta^2 \sum_{1 \leq k_1 < k_2 \leq n} w_{k_1,i}^2 w_{k_2,i}^2 \\ &= \sigma^4 \beta (1 - 3\beta) \sum_{k=1}^n w_{k,i}^4 + 3\sigma^4 \beta^2,\end{aligned}\quad (32)$$

- Γ_4, Γ_5 :

$$\mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \Gamma_4 = 0, \quad \mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \Gamma_5 = 0. \quad (33)$$

Substitute $\mathbb{E}\Gamma_1, \mathbb{E}\Gamma_2, \mathbb{E}\Gamma_3, \mathbb{E}\Gamma_4, \mathbb{E}\Gamma_5$ back to equation 29, yields

$$\frac{1}{np} \mathbb{E}_{\mathbf{X}_o, \mathbf{B}, \mathbf{S}} \|\mathbf{A}\mathbf{Y}_C\|_4^4 = 3\theta(1-\theta) \frac{\|\mathbf{W}^* \mathbf{D}_o\|_4^4}{n} + \sigma^4 \beta (1 - 3\beta) \frac{\|\mathbf{W}\|_4^4}{n} + 3\theta^2 + 6\sigma^2 \theta \beta + 3\sigma^4 \beta^2. \quad (34)$$

■

B ADDITIONAL EXPERIMENTAL RESULTS

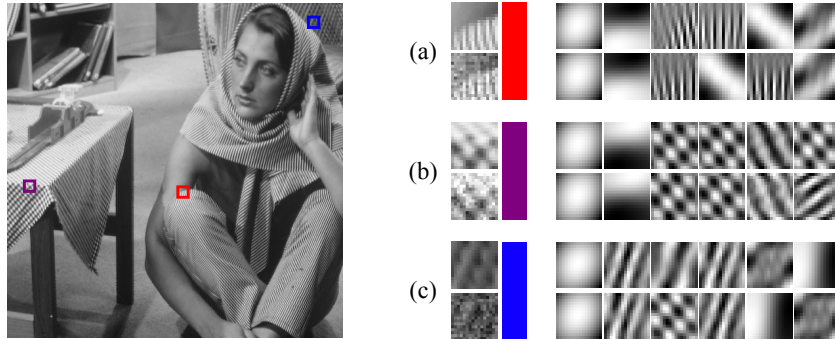


Figure 8: Representations of three 16×16 textured patches in both the clean and noisy images. Each selected patch is visualized, both with and without noise, and the 6 corresponding bases with largest absolute coefficients are shown.

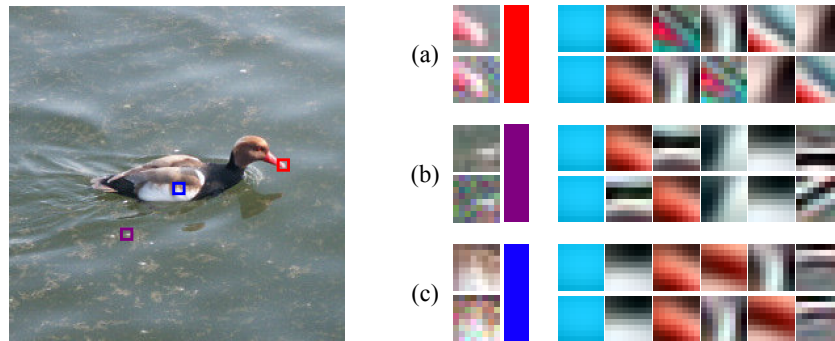


Figure 9: Representations of three $8 \times 8 \times 3$ patches from the colored images. Once again, these the patch is shown in the clean and noisy image, along with the corresponding top bases in the learned sparse representations.

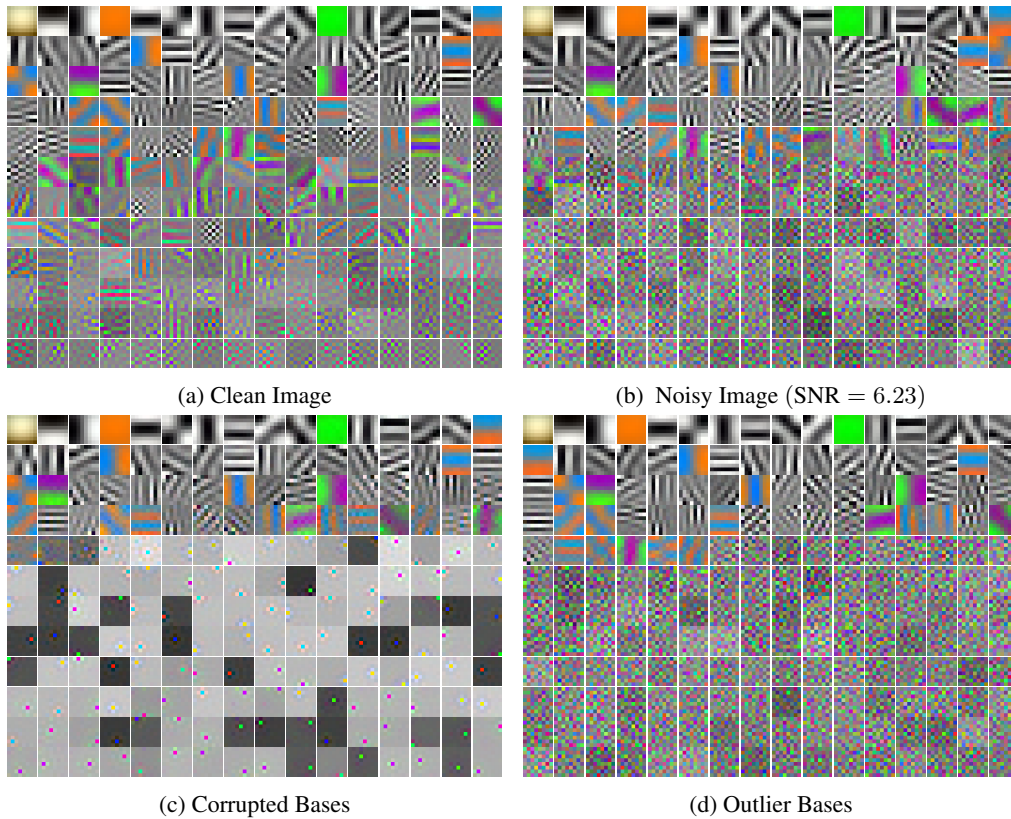


Figure 10: All $8 \times 8 \times 3 = 192$ bases learned from 100,000 random 8×8 colored patches sampled from the CIFAR-10 data-set. (a) Learned Bases from clean CIFAR-10; (b) Learned Bases from CIFAR-10 with Gaussian noise, SNR = 6.23; (c) Learned Bases from CIFAR-10 with 50% of sparse corruptions; (d) Learned Bases from CIFAR-10 with 20% of Gaussian outliers. All learned bases the resulting atoms are sorted according to the ℓ^1 -norm of their coefficients in the sparse code.