

# PERFORMANCE LIMITS OF DICTIONARY LEARNING FOR SPARSE CODING

Alexander Jung<sup>a</sup>, Yonina C. Eldar<sup>b</sup>, Norbert Görtz<sup>a</sup>

<sup>a</sup>Institute of Telecommunications, Vienna University of Technology, Austria; {ajung, norbert.goertz}@nt.tuwien.ac.at

<sup>b</sup>Technion—Israel Institute of Technology, Israel; e-mail: yonina@ee.technion.ac.il

## ABSTRACT

We consider the problem of dictionary learning under the assumption that the observed signals can be represented as sparse linear combinations of the columns of a single large dictionary matrix. In particular, we analyze the minimax risk of the dictionary learning problem which governs the mean squared error (MSE) performance of any learning scheme, regardless of its computational complexity. By following an established information-theoretic method based on Fano's inequality, we derive a lower bound on the minimax risk for a given dictionary learning problem. This lower bound yields a characterization of the sample-complexity, i.e., a lower bound on the required number of observations such that consistent dictionary learning schemes exist. Our bounds may be compared with the performance of a given learning scheme, allowing to characterize how far the method is from optimal performance.

**Index Terms**—Dictionary Identification, Dictionary Learning, Big Data, Minimax Risk, Fano Inequality.

## 1. INTRODUCTION

Consider observing  $N$  signals  $\mathbf{y}_k \in \mathbb{R}^m$ ,  $k = 1, \dots, N$ , which are assumed to be sparse linear combinations of the columns of an underlying dictionary  $\mathbf{D} \in \mathbb{R}^{m \times p}$ . Each signal  $\mathbf{y}_k$  is an i.i.d. realization of the random vector

$$\mathbf{y} = \mathbf{D}\mathbf{x} + \mathbf{w}. \quad (1)$$

The matrix  $\mathbf{D} \in \mathbb{R}^{m \times p}$ , with  $p \geq m$ , is the underlying dictionary we wish to learn. The random vector  $\mathbf{x} \in \mathbb{R}^p$  is a sparse coefficient vector and  $\mathbf{w} \in \mathbb{R}^m$  denotes zero-mean additive white Gaussian noise with variance  $\sigma^2 > 0$ , i.e.,  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ . The dictionary learning problem is relevant to a wide range of applications and has been studied extensively. In particular, dictionary learning is applied to *Big Data* applications aiming at discovering an intrinsic low dimensional structure in very high-dimensional data, in order to make the processing of this data flood tractable.

*State of the Art:* A variety of (locally) efficient learning schemes have been proposed and analyzed in the literature (e.g., [1, 2, 3, 4, 5, 6]). In [6] the authors apply a variant of the *approximate message passing* scheme [7] to the dictionary

learning problem. The works in [2, 3, 4, 5, 8] consider estimates of the dictionary obtained by solving the (non-convex) minimization problem

$$\min_{\mathbf{D}, \mathbf{X}} \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_{\text{F}}^2 + \lambda \|\mathbf{X}\|_1, \quad (2)$$

where the  $k$ th columns of  $\mathbf{X}$  and  $\mathbf{Y}$  are given by the  $k$ th i.i.d. realizations  $\mathbf{y}_k$  and  $\mathbf{x}_k$ , respectively, and  $\|\mathbf{X}\|_1 := \sum_{k,l} |X_{k,l}|$ . The authors of [2, 3, 4] give upper bounds on the distance between the generating dictionary and the nearest local minimum of (2). Based on these characterizations of the local minima, it has been shown in [4], for the noiseless and square dictionary setting, that  $N \propto p \log(p)$  observations are sufficient to guarantee local identifiability of the generating dictionary. By contrast, the authors of [3] obtain a sample-complexity of  $N \propto p^3 m$  in the case of overcomplete dictionaries and noisy observations. The analysis presented in [2, 3, 4] is conceptually different from our analysis, since we focus on the (worst-case) MSE of learning schemes, whereas [2, 3, 4] characterize the existence of local minima (of (2)) close to the generating dictionary. We would also like to mention an exciting recent line of work [9, 10, 11, 12] presenting dictionary learning schemes that are proven to globally recover the generating dictionary.

*Contribution:* By now there have been proposed quite a few dictionary learning schemes, whose performance is theoretically analyzed in terms of a characterization of the sample size sufficient for (local) identification of the generating dictionary. However, an investigation of fundamental performance limits for the dictionary learning problem seems to be missing. Here, we close this gap and present a lower bound on the minimax risk for the dictionary learning problem, where the estimation quality is measured by the Frobenius norm. This bound applies to any algorithm, regardless of its computational complexity and seems to be the first analysis that targets directly the MSE of learning schemes. For the derivation of the lower bound, we make use of an established information-theoretic approach to minimax estimation, which is based on Fano's inequality [13]. Although this approach has been successfully applied to several other (sparse) estimation problems [14, 15, 16, 17], the adaptation of this method to the problem of dictionary learning for sparse coding seems to be new.

*Outline of the Paper:* We begin in Section 2 with a formalization of the problem setup and discuss the adaption of the information-theoretic proof method (for lower bounding the minimax risk) to this setting. A lower bound on the minimax risk for dictionary learning is presented in Section 3. A sketch of the proof is given in Section 4.

*Notation:* Given a natural number  $k \in \mathbb{N}$ , we define the set  $[k] \triangleq \{1, \dots, k\}$ . For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times p}$ , we denote its Frobenius norm by  $\|\mathbf{A}\|_F \triangleq \sqrt{\text{Tr}\{\mathbf{A}\mathbf{A}^T\}}$ . The  $k$ th column of the identity matrix is denoted by  $\mathbf{e}_k$ . The complementary Kronecker delta is denoted by  $\bar{\delta}_{l,l'}$ , where  $\bar{\delta}_{l,l'} = 0$  if  $l = l'$  and is equal to one otherwise. The determinant of a square matrix  $\mathbf{C}$  is denoted by  $|\mathbf{C}|$ . We denote by  $\mathbb{E}_{\mathbf{Z}}\{\cdot\}$  the expectation w.r.t. the distribution of the random vector or matrix  $\mathbf{Z}$ .

## 2. PROBLEM FORMULATION

### 2.1. The Dictionary Learning Problem

Consider the model (1). We collect the measurements into the observation matrix

$$\mathbf{Y} := (\mathbf{y}_1, \dots, \mathbf{y}_N) \in \mathbb{R}^{m \times N}, \quad (3)$$

where  $\mathbf{y}_k$  is an i.i.d. realization of the random vector given by (1). The underlying generating dictionary  $\mathbf{D}$  is modeled as deterministic but unknown. We assume the columns of  $\mathbf{D}$  to be normalized, i.e.,

$$\mathbf{D} \in \mathcal{D} \triangleq \{\mathbf{B} \in \mathbb{R}^{m \times p} \mid \|\mathbf{B}\mathbf{e}_j\|_2 = 1, \text{ for all } j \in [p]\}. \quad (4)$$

The set  $\mathcal{D}$  is known as the *oblique manifold* [3]. Moreover, we assume the true dictionary  $\mathbf{D}$  to be obtained as a small perturbation of a known ‘‘reference dictionary’’  $\mathbf{D}_0$ . In particular, for some small radius  $r > 0$ , we require

$$\mathbf{D} \in \mathcal{X}(\mathbf{D}_0, r) := \{\mathbf{D}' \in \mathcal{D} \mid \|\mathbf{D}' - \mathbf{D}_0\|_F \leq r\} \quad (5)$$

The statistics of the coefficient vector  $\mathbf{x}$  is modeled such that it is a strictly  $s$ -sparse vector. In particular, we introduce the random variable  $i$ , which is chosen uniformly at random (u.a.r.) from the set  $\binom{[p]}{s}$ . A specific value of  $i$  represents a certain index set  $\mathcal{S}(i) \subseteq [p]$  containing  $s$  different indices. More formally, the map

$$\mathcal{S}(\cdot) : \left[ \binom{[p]}{s} \right] \rightarrow \mathcal{E} \triangleq \{\mathcal{I} \subseteq [p], |\mathcal{I}| = s\} \quad (6)$$

is a bijection from the first  $\binom{[p]}{s}$  natural numbers to the set  $\mathcal{E}$  of all size- $s$  subsets  $\mathcal{I}$  of  $[p]$ .

The random variable  $i$  selects the active coefficients of  $\mathbf{x}$ , i.e.,

$$\text{supp}(\mathbf{x}) = \mathcal{S}(i), \text{ and } \mathbf{x}_{\mathcal{S}(i)} \sim \mathcal{N}(\mathbf{0}, \sigma_a^2 \mathbf{I}). \quad (7)$$

The (unconditional) covariance matrix of the sparse coefficient vector  $\mathbf{x}$  is given by

$$\Sigma_{\mathbf{x}} \triangleq \mathbb{E}\{\mathbf{x}\mathbf{x}^T\} = (s/p)\sigma_a^2 \mathbf{I}. \quad (8)$$

We define the signal to noise ratio of the observation model (1) as  $\text{SNR} := (\sigma_a/\sigma)^2$ .

Since the columns of  $\mathbf{Y}$  are i.i.d. realizations of the vector  $\mathbf{y}$  in (1), the conditional probability density function (pdf) of the observation  $\mathbf{Y}$ , given the  $N$  i.i.d. realizations  $\mathbf{i} = (i_1, \dots, i_N)$  of the random support index  $i$ , is

$$f_{\mathbf{D}}(\mathbf{Y}|\mathbf{i}) = \prod_{k \in [N]} \frac{\exp\left(- (1/2)\mathbf{y}_k^T \Sigma_{\mathbf{y}|i_k}^{-1} \mathbf{y}_k\right)}{(2\pi)^{m/2} |\Sigma_{\mathbf{y}|i_k}|^{1/2}}.$$

Here,  $\Sigma_{\mathbf{y}|i} \triangleq \mathbb{E}\{\mathbf{y}\mathbf{y}^T|i\}$  denotes the conditional covariance matrix of  $\mathbf{y}$ , given  $i$ , and reads explicitly as  $\Sigma_{\mathbf{y}|i} = \sigma_a^2 \mathbf{D}_{\mathcal{S}(i)} \mathbf{D}_{\mathcal{S}(i)}^T + \sigma^2 \mathbf{I}$ .

We note that any learning scheme based on the model (1) faces an intrinsic sign and permutation ambiguity for the dictionary  $\mathbf{D}$ . Indeed, by observing  $\mathbf{Y}$  only, one cannot distinguish between dictionaries which are related via column permutations and sign-flips of the columns [3, 4]. While we do not take this intrinsic ambiguity into account explicitly, our results are meaningful as they apply to dictionary learning problems where the true dictionary belongs to the (small) neighborhood  $\mathcal{X}(\mathbf{D}_0, r)$  of a known reference dictionary  $\mathbf{D}_0$ .

We investigate the fundamental limits on the accuracy achievable by any learning scheme producing an estimate  $\hat{\mathbf{D}}(\mathbf{Y})$  of the underlying dictionary based on the observation  $\mathbf{Y}$ . For the moment, suppose that we have access to the coefficients  $\mathbf{x}$  in (1) and the estimate  $\hat{\mathbf{D}}$  is held fixed, i.e., does not depend on the observation  $\mathbf{Y}$ . Then, we obtain for the prediction error, when using the estimate  $\hat{\mathbf{D}}$  instead of the generating dictionary  $\mathbf{D}$ ,

$$\mathbb{E}_{\mathbf{x}}\{\|\mathbf{D}\mathbf{x} - \hat{\mathbf{D}}\mathbf{x}\|^2\} \stackrel{(8)}{=} (s/p)\sigma_a^2 \|\mathbf{D} - \hat{\mathbf{D}}\|_F^2. \quad (9)$$

Therefore, the prediction error is proportional to the squared Frobenius norm of the estimation error  $\mathbf{D} - \hat{\mathbf{D}}$ . Based on (9), we measure the accuracy of a specific learning scheme  $\hat{\mathbf{D}}(\cdot)$  by the MSE  $\varepsilon(\mathbf{D}, \hat{\mathbf{D}}(\cdot)) \triangleq \mathbb{E}_{\mathbf{Y}}\{\|\hat{\mathbf{D}}(\mathbf{Y}) - \mathbf{D}\|_F^2\}$ . Note that the MSE depends on the underlying generating dictionary  $\mathbf{D}$  and the learning scheme  $\hat{\mathbf{D}}(\cdot)$ .

Define the minimax risk  $\varepsilon$  for the problem of learning the dictionary  $\mathbf{D}$  based on the observation of  $N$  i.i.d. realizations of  $\mathbf{y}$  in (1), as

$$\varepsilon \triangleq \inf_{\hat{\mathbf{D}}} \sup_{\mathbf{D} \in \mathcal{X}(\mathbf{D}_0, r)} \varepsilon(\mathbf{D}, \hat{\mathbf{D}}(\cdot)). \quad (10)$$

The minimax risk  $\varepsilon$  will in general depend on the number of observations  $N$ , the dimension  $m$  of the observed signals, the number of signal expansion coefficients  $p$ , the sparsity degree  $s$  and the variance parameters  $\sigma_a^2$  and  $\sigma^2$ . However, to lighten notation, we will not make this dependence explicit. Our goal is to develop a lower bound on  $\varepsilon$  using an information-theoretic method.

Having a lower bound for the minimax risk allows us to assess the performance of a given dictionary learning scheme.

In particular, if the MSE of a given algorithm is close to the minimax risk, or a lower bound to it, then there is little point to hope for finding improved techniques with substantially better performance.

## 2.2. Information Theory of Dictionary Learning

Our approach to bounding the minimax risk  $\varepsilon$  of (10) is to use the information-theoretic method put forward in [18, 14, 16]. However, the key challenge in applying this technique is the fact that the vector  $\mathbf{y}$  given by (1) does not follow a multivariate normal distribution. Indeed, due to the prior model for the coefficient vector  $\mathbf{x}$  (cf. (7)), the vector  $\mathbf{y}$  follows a Gaussian mixture model, with a mixture component associated with each specific value of the support index  $i$ .

In order to apply the information-theoretic technique, it is necessary to have a precise characterization of the mutual information  $I(\mathbf{Y}; l)$  between the observation  $\mathbf{Y}$  and a random index  $l$  which selects the generating dictionary  $\mathbf{D} = \mathbf{D}^{(l)}$  u.a.r. from a finite set  $\mathcal{D}_0 \subseteq \mathcal{D}$ . Obtaining a bound on  $I(\mathbf{Y}; l)$  typically involves the analysis of the Kullback Leibler (KL) divergence between the distributions of  $\mathbf{Y}$  implied by different dictionaries  $\mathbf{D} = \mathbf{D}^{(l)}$ . However, exact characterizations of the KL divergence between Gaussian mixture models is in general not possible and one has to resort to approximations or bounds [19].

A main conceptual contribution of this work is a strategy to avoid evaluating KL divergences between Gaussian mixture models. Instead, we rely on the following decomposition, which follows from the chain rule for mutual information,

$$\begin{aligned} I(\mathbf{Y}; l) &= I(\mathbf{Y}, \mathbf{i}; l) - I(l; \mathbf{i} | \mathbf{Y}) \\ &= I(\mathbf{Y}; l | \mathbf{i}) + \underbrace{I(l; \mathbf{i})}_{=0} - I(l; \mathbf{i} | \mathbf{Y}) \\ &= I(\mathbf{Y}; l | \mathbf{i}) - I(l; \mathbf{i} | \mathbf{Y}). \end{aligned} \quad (11)$$

Here,  $I(\mathbf{Y}; l | \mathbf{i})$  denotes the conditional mutual information between the observation  $\mathbf{Y}$  and the random index  $l$ , given the support indices  $\mathbf{i} = (i_1, \dots, i_N)$ . The components of the decomposition in (11) have particular interpretations. The term  $I(\mathbf{Y}; l | \mathbf{i})$  characterizes the difficulty of detecting the (index of the) generating dictionary  $\mathbf{D} = \mathbf{D}^{(l)}$ , if we had access to the indices  $i_k$  selecting the active coefficients of  $\mathbf{x}_k$ . The second term, i.e.,  $I(l; \mathbf{i} | \mathbf{Y})$  quantifies the dependence between the support of the sparse coefficient vector  $\mathbf{x}$  and the (index  $l$  of the) generating dictionary  $\mathbf{D} = \mathbf{D}^{(l)}$ , after observing  $\mathbf{Y}$ .

Since  $I(l; \mathbf{i} | \mathbf{Y}) \geq 0$  [13, Ch. 2], we can upper bound  $I(\mathbf{Y}; l)$  by upper bounding  $I(\mathbf{Y}; l | \mathbf{i})$ . Note that, conditioned on the support index  $i$ , the data vector  $\mathbf{y}$  in (1) follows a normal distribution with covariance matrix  $\Sigma_{y|i}$ , which renders the problem of upper bounding  $I(\mathbf{Y}; l | \mathbf{i})$  tractable. We detail this proof technique in Section 4.

## 3. A LOWER BOUND ON THE MINIMAX RISK

A typical requirement for sparse (compressed sensing) recovery to work well, even when the dictionary  $\mathbf{D}$  in (1) is known, is the validity of [20, 21]

$$m \geq c_0 s \log(p/s), \quad (12)$$

with some absolute constant  $c_0$ . Since we consider the more difficult problem of dictionary learning, i.e., we treat the dictionary as an unknown parameter, we expect (12) to be a necessary requirement for the existence of accurate dictionary learning schemes.

Our main result is the following lower bound on the minimax risk for a given dictionary learning problem.

**Theorem 3.1.** *Consider a dictionary learning problem based on  $N$  i.i.d. observations following the model (1) and the true dictionary satisfying (5) with  $r \leq 1/\sqrt{p}$ . Then, if*

$$p > 64, \text{ and } m \geq 192s(9 + 2 \log(p/s)), \quad (13)$$

*the minimax risk  $\varepsilon$  is lower bounded as*

$$\varepsilon \geq \min \left\{ r^2/16, \frac{\text{SNR}^{-1} p^2}{5120Ns} \right\}. \quad (14)$$

We highlight the fact that Theorem 3.1 does not place any assumptions (like incoherence or restricted isometry properties) on the underlying generating dictionary.

For sufficiently large sample-size  $N$ , such that  $N \gg p^2/s$ , the second bound in (14) will be in force. This bound shows a dependence on the sample-size via  $1/N$  which clearly makes sense. Indeed, by averaging the outcomes of a learning scheme over blocks of independent observations the MSE is expected to scale inversely proportional to the sample size  $N$ . This dependence of the MSE on the sample-size is also observed in the empirical results of simulation studies for specific learning schemes in [3, 1]. Moreover, the theoretic results presented in [3, 22] suggest that the estimation error of certain learning schemes, measured by the squared Frobenius norm, scales inversely proportional to  $N$ .

For the case of constant sparsity, i.e., when  $s \leq C_0$  for some constant (independent of  $p$ ) our lower bound scales as  $\Theta(p^2/N)$ , suggesting a sample-complexity of  $\Theta(p^2)$ . This scaling is considerably smaller than the sample complexity  $\mathcal{O}(p^3 m)$ , which [3] proved to be sufficient in the noisy and over-complete case, such that the estimator based on minimizing (2) performs well.

## 4. PROOF OF THE MAIN RESULT

The proof of Theorem 3.1 is based on reducing the minimax estimation problem (10) to a specific multiple hypothesis testing problem. In particular, we assume that the generating dictionary  $\mathbf{D}$  in (1) is taken from a finite subset  $\mathcal{D}_0 \triangleq \{\mathbf{D}^{(l)}\}_{l \in [L]} \subseteq \mathcal{X}(\mathbf{D}_0, r)$  for some  $L \in \mathbb{N}$ . This subset  $\mathcal{D}_0$  is constructed such that (i) any two distinct dictionaries  $\mathbf{D}^{(l)}, \mathbf{D}^{(l')} \in \mathcal{D}_0$  are separated by at least  $\sqrt{8\varepsilon}$ , i.e.,

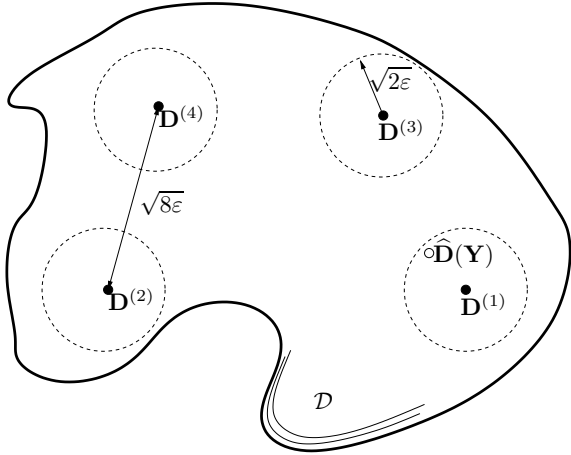


Fig. 1. Finite ensemble  $\mathcal{D}_0$  of size  $L = 4$ .

$\|\mathbf{D}^{(l)} - \mathbf{D}^{(l')}\|_F \geq \sqrt{8\varepsilon}$  and (ii) it is hard to detect the generating dictionary  $\mathbf{D}$  if it is drawn u.a.r. from  $\mathcal{D}_0$ . However, we do not specify a deterministic scheme to construct such a set  $\mathcal{D}_0$ . We merely use a probabilistic method to show that there must exist at least one such set  $\mathcal{D}_0$ . The existence of  $\mathcal{D}_0$  then yields, via Lemma 4.1, a relation between the sample-size  $N$  and the remaining model parameters  $m, p, s, \sigma_a, \sigma$  which has to be satisfied such that an estimator with worst-case MSE not exceeding  $\varepsilon$  may exist.

In Fig. 1, we sketch the idea of this method for the particular case of a subset  $\mathcal{D}_0 := \{\mathbf{D}^{(1)}, \dots, \mathbf{D}^{(4)}\}$  containing four dictionaries  $\mathbf{D}^{(l)} \in \mathcal{X}(\mathbf{D}_0, r)$ . We also show a realization of the estimator  $\hat{\mathbf{D}}(\mathbf{Y})$ . Two different dictionaries in  $\mathcal{D}_0$  are separated by at least  $\sqrt{8\varepsilon}$ . In particular, if  $\hat{\mathbf{D}}(\mathbf{Y})$  is a learning scheme achieving the minimax risk in (10), then the minimum distance detector

$$\operatorname{argmin}_{\mathbf{D}' \in \mathcal{D}_0} \|\hat{\mathbf{D}}(\mathbf{Y}) - \mathbf{D}'\|_F$$

recovers the correct dictionary  $\mathbf{D} \in \mathcal{D}_0$  if  $\hat{\mathbf{D}}(\mathbf{Y})$  belongs to the ball  $\mathcal{B}(\mathbf{D}, \sqrt{2\varepsilon})$  (indicated by a dashed circle in Fig. 1) centered at  $\mathbf{D}$  and with radius  $\sqrt{2\varepsilon}$ . The information-theoretic method [15, 14, 18] of lower bounding the minimax risk  $\varepsilon$  consists then in relating the probability  $\mathbb{P}\{\hat{\mathbf{D}}(\mathbf{Y}) \notin \mathcal{B}(\mathbf{D}, \sqrt{2\varepsilon})\}$  to the mutual information between the observation  $\mathbf{Y}$  and the dictionary  $\mathbf{D}$  which is assumed to be drawn u.a.r. from  $\mathcal{D}_0$ .

In particular, our analysis is based on the construction of a finite set  $\mathcal{D}_0 \triangleq \{\mathbf{D}^{(1)}, \dots, \mathbf{D}^{(L)}\} \subseteq \mathcal{X}(\mathbf{D}_0, r)$  of  $L$  distinct dictionaries belonging to  $\mathcal{D}$  (cf. (4)) having the following desiderata:

- For any two dictionaries  $\mathbf{D}^{(l)}, \mathbf{D}^{(l')} \in \mathcal{D}_0$ ,

$$\|\mathbf{D}^{(l)} - \mathbf{D}^{(l')}\|_F^2 \geq \bar{\delta}_{l,l'} 8\varepsilon. \quad (15)$$

- If the generating dictionary in (1) is chosen as  $\mathbf{D} = \mathbf{D}^{(l)} \in \mathcal{D}_0$ , where  $l$  is selected u.a.r. from  $[L]$ , then the conditional mutual information between  $\mathbf{Y}$  and  $l$ , given  $\mathbf{i}$ , is bounded as

$$I(\mathbf{Y}; l | \mathbf{i}) \leq \eta \quad (16)$$

with some given small  $\eta$ .

The following result gives precise conditions on the cardinality  $L$  and threshold  $\eta$  such that at least one subset  $\mathcal{D}_0 \subseteq \mathcal{X}(\mathbf{D}_0, r)$  of size  $L$  satisfying (15) as well as (16) is guaranteed to exist.

**Lemma 4.1.** *Consider a dictionary learning problem based on (5) with some  $r \leq 1/\sqrt{p}$ . Then, for any  $\varepsilon$  such that*

$$\varepsilon < r^2/16,$$

*there exists a set  $\mathcal{D}_0 \subseteq \mathcal{X}(\mathbf{D}_0, r)$  of cardinality  $L = e^{p/32}$  such that (15) and (16) are satisfied with  $\eta = 32\varepsilon N s \text{SNR}/p$ .*

The next result, which is the central argument of the information-theoretic method for lower bounding minimax risk, relates the cardinality  $L$  of a subset  $\mathcal{D}_0 \subseteq \mathcal{D}$  to the conditional mutual information  $I(\mathbf{Y}; l | \mathbf{i})$  between the observation  $\mathbf{Y}$  and a random index  $l$  selecting the generating dictionary u.a.r. from  $\mathcal{D}_0$ .

**Lemma 4.2.** *Consider the dictionary learning problem (1) with minimax risk  $\varepsilon$  ((10)) and a finite set  $\mathcal{D}_0 \subseteq \mathcal{X}(\mathbf{D}_0, r)$  consisting of  $L$  distinct dictionaries  $\mathbf{D}^{(l)} \in \mathbb{R}^{m \times p}$  such that*

$$\|\mathbf{D}^{(l)} - \mathbf{D}^{(l')}\|_F^2 \geq 8\bar{\delta}_{l,l'}\varepsilon.$$

*Then, it holds  $I(\mathbf{Y}; l | \mathbf{i}) \geq (1/2) \log_2(L) - 1$ .*

The proofs of Lemma 4.1 and 4.2 are omitted due to space limitations.

*Proof of Theorem 3.1:* According to Lemma 4.1, for any  $\varepsilon < r^2/8$ , with  $r \leq 1/\sqrt{p}$ , there exists a set  $\mathcal{D}_0 \subseteq \mathcal{X}(\mathbf{D}_0, r)$  of cardinality  $L = e^{p/32}$  satisfying (15) and (16) with  $\eta = 32N s \text{SNR}^2 \varepsilon/p$ . Applying Lemma 4.2 to the set  $\mathcal{D}_0$  yields, in turn,

$$32N s \text{SNR} \varepsilon/p \geq I(\mathbf{Y}; l | \mathbf{i}) \geq (1/2) \log_2(L) - 1$$

implying

$$\varepsilon \geq \frac{\text{SNR}^{-1}}{32N s} p((1/2) \log_2(L) - 1).$$

Since

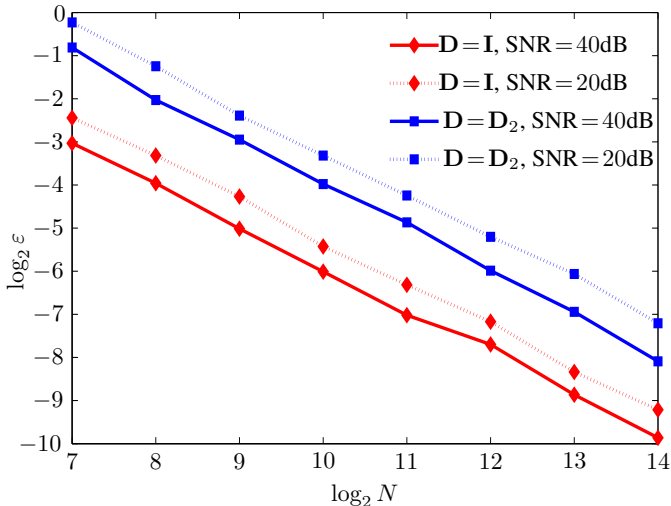
$$(1/2) \log_2(L) - 1 \geq 0.7p/32 - 1 \stackrel{(13)}{\geq} 0.2p/32,$$

we arrive at (14).

## 5. NUMERICAL EXPERIMENTS

One of the uses of the lower bound on the minimax risk stated in Theorem 3.1 is that it allows for an assessment of the performance of practical learning schemes. In this section we compare the lower bound (14) with the actual MSE of an (locally) efficient learning scheme  $\hat{\mathbf{D}}_{\text{ITKM}}(\mathbf{Y})$ , termed *iterative thresholding and K-means* (ITKM) algorithm, which has been proposed recently [22]. We applied the ITKM algorithm with sparsity parameter  $\tilde{s} = 1$ , using oracle initialization and signal normalization<sup>1</sup>, to a data matrix  $\mathbf{Y} \in \mathbb{R}^{m \times N}$ ,

<sup>1</sup>For background and notation, we refer to [22].



**Fig. 2.** MSE curves of the ITKM learning scheme for  $m = 8$ , with  $m = 8$ , whose columns are independent realizations of  $\mathbf{y}$  according to (1) with  $s = 2$ . For the underlying generating dictionary  $\mathbf{D}$  we choose the identity matrix  $\mathbf{I}$  and, in a second experiment, the concatenation of the identity matrix and the  $m \times m$  normalized Hadamard matrix  $\mathbf{F}_m^2$ , i.e.,  $\mathbf{D} = \mathbf{D}_2 := [\mathbf{I} \sqrt{1/m} \mathbf{F}_m]$ . For both choices for the generating dictionary we set  $m = 8$  and  $s = 2$ . In Fig. 2, we plot the actual MSE  $\varepsilon(\mathbf{D}, \hat{\mathbf{D}}_{\text{ITKM}}(\cdot))$  for varying sample-size  $N$  and different values of the SNR. The bound (14) correctly predicts the slope  $1/N$  of the curves. However, the absolute position of the lower bound (14) is significantly below that of the actual MSE curves. While this could mean that the performance of ITKM is far from optimum, there is also the possibility that the lower bound (14) can be tightened (made higher) considerably by taking also the term  $I(l; \mathbf{i} | \mathbf{Y})$  in (11) into account.

## 6. CONCLUSION

We derived a lower bound on the minimax risk for dictionary learning, which seems to be the first result of this kind. This lower bound yields, in turn, a characterization of the required sample-size, i.e., the sample-complexity, such that accurate learning schemes, regardless of computational complexity, may exist. Comparing our results with the sample-complexity of some popular learning schemes, which are mainly based on minimizing (2), reveals that there may be other algorithms requiring significantly fewer observations. Finally, we note that our lower bound complements the sufficient conditions on the sample-complexity for dictionary learning derived in [23].

## 7. ACKNOWLEDGMENT

The authors would like to thank Karin Schnass for sharing here expertise on practical dictionary learning schemes and

<sup>2</sup>For  $m$  being a power of 2, the Hadamard matrix  $\mathbf{F}_m$  is defined recursively by  $\mathbf{F}_1 = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$  and  $\mathbf{F}_m = \mathbf{F}_1 \otimes \mathbf{F}_{m/2}$ .

for providing some simulation results.

## 8. REFERENCES

- [1] K. Schnass, "On the identifiability of overcomplete dictionaries via the minimisation principle underlying k-svd," *Applied and Computational Harmonic Analysis*, 2014.
- [2] M. Aharon, M. Elad, and A. Bruckstein, "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Trans. Signal Processing*, vol. 54, no. 11, pp. 4311–4322, 2006.
- [3] R. Jenatton, R. Gribonval, and F. Bach, "Local stability and robustness of sparse dictionary learning in the presence of noise," *ArXiv e-prints*, Oct. 2012.
- [4] R. Gribonval and K. Schnass, "Dictionary identification - sparse matrix-factorization via  $\ell_1$ -minimization," *IEEE Trans. Inf. Theory*, vol. 56, no. 7, Jul. 2010.
- [5] M. Yaghoobi, T. Blumensath, and M.E. Davies, "Dictionary learning for sparse approximations with the majorization method," *IEEE Trans. Signal Processing*, vol. 57, no. 6, pp. 2178–2191, 2009.
- [6] J. T. Parker, P. Schniter, and V. Cevher, "Bilinear generalized approximate message passing," *ArXiv e-prints*, Oct. 2013.
- [7] David L. Donoho, Arian Maleki, and Andrea Montanari, "Message-passing algorithms for compressed sensing," *PNAS*, vol. 106, no. 45, pp. 18914–18919, 2009.
- [8] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, "Online dictionary learning for sparse coding," in *Proceedings of the 26th Annual International Conference on Machine Learning*, Montreal, Canada, 2009, ICML '09, pp. 689–696.
- [9] D. Spielman, H. Wang, and J. Wright, "Exact recovery of sparsely-used dictionaries," in *Conference on Learning Theory (arXiv:1206.5882)*, 2012.
- [10] S. Arora, R. Ge, and A. Moitra, "New algorithms for learning incoherent and overcomplete dictionaries," *arXiv:1308.6273*, 2013.
- [11] A. Agarwal, A. Anandkumar, and P. Netrapalli, "Exact recovery of sparsely used overcomplete dictionaries," *arXiv:1309.1952*, 2013.
- [12] A. Agarwal, A. Anandkumar, P. Jain, P. Netrapalli, and R. Tandon, "Learning sparsely used overcomplete dictionaries via alternating minimization," *arXiv:1310.7991*, 2013.
- [13] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, Wiley, New Jersey, 2 edition, 2006.
- [14] W. Wang, M. J. Wainwright, and K. Ramchandran, "Information-theoretic bounds on model selection for Gaussian Markov random field," in *Proc. IEEE ISIT-2010*, Austin, TX, Jun. 2010, pp. 1373–1377.
- [15] M. J. Wainwright, "Information-theoretic Limits on Sparsity Recovery in the High-Dimensional and Noisy Setting," *IEEE Trans. Inf. Theory*, vol. 55, no. 12, pp. 5728–5741, Dec. 2009.
- [16] E. J. Candès and M. A. Davenport, "How well can we estimate a sparse vector?," *Applied and Computational Harmonic Analysis*, vol. 34, no. 2, pp. 317–323, 2013.
- [17] T. T. Cai and H. H. Zhou, "Optimal rates of convergence for sparse covariance matrix estimation," *Ann. Stat.*, vol. 40, no. 5, pp. 2359–2763, 2012.
- [18] B. Yu, "Assouad, Fano, and Le Cam," in *Festschrift for Lucien Le Cam*, D. Pollard, E. Torgersen, and G. L. Yang, Eds., pp. 423–435. Springer New York, 1997.
- [19] J. R. Hershey and P. A. Olsen, "Approximating the Kullback Leibler Divergence Between Gaussian Mixture Models," in *Proc. IEEE ICASSP-2007*, 2007, pp. 317–320.
- [20] S. Foucart and H. Rauhut, *A Mathematical Introduction to Compressive Sensing*, Springer, New York, 2012.
- [21] Y. C. Eldar and G. Kutyniok, *Compressed Sensing: Theory and Applications*. Cambridge Univ. Press, Cambridge, UK, 2012.
- [22] K. Schnass, "Local identification of overcomplete dictionaries," *ArXiv e-prints*, 2014.
- [23] D. Vainsencher, S. Mannor, and A. M. Bruckstein, "The sample complexity of dictionary learning," *Journal Mach. Learn. Research*, vol. 12, pp. 3259–3281, 2011.