

---

# Towards improving discriminative reconstruction via simultaneous dense and sparse coding

---

Anonymous Author  
Anonymous Institution

## Abstract

Discriminative features extracted from the sparse coding model have been shown to perform well for classification and reconstruction. Recent deep learning architectures have further improved reconstruction in inverse problems by considering new dense priors learned from data. We propose a novel dense and sparse coding model that integrates both representation capability and discriminative features. The model considers the problem of recovering a dense vector  $\mathbf{x}$  and a sparse vector  $\mathbf{u}$  given measurements of the form  $\mathbf{y} = \mathbf{Ax} + \mathbf{Bu}$ . Our first analysis proposes a natural geometric condition based on the minimal angle between spanning subspaces corresponding to the measurement matrices  $\mathbf{A}$  and  $\mathbf{B}$  to establish the uniqueness of solutions to the linear system. The second analysis shows that, under mild assumptions, a convex program recovers the dense and sparse components. We validate the effectiveness of the proposed model on simulated data and propose a dense and sparse autoencoder (DenSaE) tailored to learning the dictionaries from the dense and sparse model. We demonstrate that a) DenSaE denoises natural images better than architectures derived from the sparse coding model ( $\mathbf{Bu}$ ), b) in the presence of noise, training the biases in the latter amounts to implicitly learning the  $\mathbf{Ax} + \mathbf{Bu}$  model, c)  $\mathbf{A}$  and  $\mathbf{B}$  capture low- and high-frequency contents, respectively, and d) compared to the sparse coding model, DenSaE offers a balance between discriminative power and representation.

## 1 Introduction

Given a data set, learning a dictionary in which each example admits a sparse representation is tremendously useful in a number of tasks (Aharon et al., 2006; Mairal et al., 2011). This problem, known as sparse coding (Olshausen and Field, 1997) or dictionary learning (Garcia-Cardona and Wohlberg, 2018), has been the subject of significant investigation in recent years in the signal processing community. A growing body of work has mapped the sparse coding problem into encoders for sparse recovery (Gregor and Lecun, 2010), and into autoencoders for purely classification (Rolfe and LeCun, 2013) or denoising (Simon and Elad, 2019; Tolooshams et al., 2020) purposes.

Autoencoders are widely used for unsupervised learning. Their integration with supervised tasks and classifiers has become popular for their regularization power and reduction of the generalization gap (Vincent et al., 2010; Epstein et al., 2018; Epstein and Meir, 2019). Rolfe and LeCun (2013) have shown benefits of autoencoders and sparse features in discriminative tasks.

For data reconstruction, recent work has highlighted some limitations of convolutional sparse coding (CSC) autoencoders (Simon and Elad, 2019) and its multi-layer and deep generalizations (Sulam et al., 2019; Zazo et al., 2019). Simon and Elad (2019) argue that the sparsity levels that CSC allows can only accommodate very sparse vectors, making it unsuitable to capture all features of signals such as natural images, and propose to compute the minimum mean-squared error solution under the CSC model, which is a dense vector capturing a richer set of features.

To address the aforementioned limitations of classical sparse coding, we propose a dense and sparse coding model that represents a signal as the sum of two components: one that admits a dense representation  $\mathbf{x}$  in a dictionary  $\mathbf{A}$  that is useful for reconstruction, and another whose representation  $\mathbf{u}$  is discriminative and sparse in a second dictionary  $\mathbf{B}$ . Based on empirical evidence, the authors in (Zazo et al., 2019) argue that a multi-layer extension of this model can, in principle,

have arbitrary depth. However, to our knowledge, the dense and sparse coding model has not been yet fully analyzed. Our contributions are

**Conditions for identifiability and recovery by convex optimization:** We derive conditions under which the dense and sparse representation is unique. We then propose a convex program for recovery that minimizes  $\|\mathbf{Ax}\|_2^2 + \|\mathbf{u}\|_1$ , subject to linear constraints.

**Phase-transition curves:** We demonstrate through simulations that the convex program can successfully solve the dense and sparse coding problem.

**Discriminative reconstruction:** We propose a dense and sparse autoencoder (DenSaE) that has competitive discriminative power and improves the representation capability compared to sparse networks.

**Organization:** Section 2 discusses the theoretical analysis of the dense and sparse coding problem. Phase transition, classification, and denoising experiments appear in Section 3. We conclude in Section 4.

### 1.1 Related work

We comment on the most closely related models. Given the measurements  $\mathbf{y}$ , the problem of recovering  $\mathbf{x}$  and  $\mathbf{u}$  is similar in flavor to sparse recovery in the union of dictionaries (Donoho and Huo, 2001; Elad and Bruckstein, 2002; Donoho and Elad, 2003; Soltani and Hegde, 2017; Studer et al., 2011; Studer and Baraniuk, 2014). Most results in this literature take the form of an uncertainty principle that relates the sum of the sparsity of  $\mathbf{x}$  and  $\mathbf{u}$  to the mutual coherence between  $\mathbf{A}$  and  $\mathbf{B}$ , and which guarantees that the representation is unique and identifiable by  $\ell_1$  minimization. To our knowledge, the analysis of this program is novel and in sharp contrast to classical settings in sparse approximation, in which the objective consists of a single sparsifying norm, rather than the combination of different norms. Robust PCA (Candès et al., 2011), which decomposes a matrix as the sum of low-rank and sparse matrices, uses the combination of the  $\ell_1$  and nuclear norms, giving it a flavor similar to our problem.

Our model resembles weighted LASSO (Lian et al., 2018; Mansour and Saab, 2017). Compared to weighted LASSO, we can directly map the weighted LASSO objective  $\|\mathbf{W}\boldsymbol{\alpha}\|_1$  to  $\|\mathbf{u}\|_1$  by letting  $\boldsymbol{\alpha} = [\mathbf{x} \ \mathbf{u}]^T$  and choosing appropriately the entries of a diagonal matrix  $\mathbf{W}$ , with  $W_{ij} \in \{0, 1\}$ ; however, in the weighted LASSO formulation, constraints can only be enforced on the sparse component  $\mathbf{u}$ . Our work differs in that a significant part of our analysis is the directed Euclidean norm constraint on  $\mathbf{x}$ , which recovers a unique solution  $\mathbf{x}^* \in \text{Ker}(\mathbf{A})^\perp$ . Our model can also be interpreted as a special case of Morphological Component Analysis

(MCA) (Elad et al., 2005) for  $K = 2$ ,  $\mathbf{s} = \sum_{k=1}^K \Phi_k \boldsymbol{\alpha}_k$ , with, however, some distinct differences: i) MCA encodes different morphological structures via the dictionaries  $\Phi_k$ . We encode a smooth morphological component via the whole product  $\mathbf{Ax}$ , which is conceptually different, and ii) we make no assumption of sparsity on the dense component  $\mathbf{x}$ . This leads to an optimization objective that is the combination of  $\ell_1$  and  $\ell_2$  norms, unlike that of MCA. Finally, a bare application of noisy sparse coding would treat  $\mathbf{e} = \mathbf{Ax}$  as arbitrary noise, hence i) recovers  $\mathbf{u}$  approximately and ii) cannot recover  $\mathbf{x}$ . However, in our analysis, the term  $\mathbf{Ax}$  is not just undesired noise but represents a sought-out feature. We can recover both  $\mathbf{x}$  and  $\mathbf{u}$  exactly. See Appendix C for a comparison of our model to noisy compressive sensing. We note that the full dense and sparse coding model is  $\mathbf{y} = \mathbf{Ax} + \mathbf{Bu} + \mathbf{e}$  where  $\mathbf{e}$  is Gaussian noise.

**Notation.** Lowercase and uppercase boldface letters denote column vectors and matrices, respectively. Given a vector  $\mathbf{x} \in \mathbb{R}^n$  and a support set  $S \subset \{1, \dots, n\}$ ,  $\mathbf{x}_S$  denotes the restriction of  $\mathbf{x}$  to indices in  $S$ . For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times p}$ ,  $\mathbf{A}_S$  is a submatrix of size  $m \times |S|$  with column indices in  $S$ . The column space of a matrix  $\mathbf{A}$  (the span of the columns of  $\mathbf{A}$ ) is designated by  $\text{Col}(\mathbf{A})$ , its null space by  $\text{Ker}(\mathbf{A})$ . We denote the Euclidean,  $\ell_1$  and  $\ell_\infty$  norms of a vector, respectively as  $\|\mathbf{x}\|_2$ ,  $\|\mathbf{x}\|_1$ , and  $\|\mathbf{x}\|_\infty$ . The operator and infinity norm of a matrix  $\mathbf{A}$  are respectively denoted as  $\|\mathbf{A}\|$  and  $\|\mathbf{A}\|_\infty$ . The sign function, applied componentwise to a vector  $\mathbf{x}$ , is denoted by  $\text{sgn}(\mathbf{x})$ . The indicator function is denoted by  $\mathbb{1}$ . The column vector  $\mathbf{e}_i$  denotes the vector of zeros except a 1 at the  $i$ -th location. The orthogonal complement of a subspace  $\mathbf{W}$  denoted by  $\mathbf{W}^\perp$ . The operator  $\mathcal{P}_{\mathbf{W}}$  denotes the orthogonal projection operator onto the subspace  $\mathbf{W}$ .

## 2 Theoretical Analysis

The dense and sparse coding problem studies the solutions of the linear system  $\mathbf{y} = \mathbf{Ax} + \mathbf{Bu}$ . Given matrices  $\mathbf{A} \in \mathbb{R}^{m \times p}$  and  $\mathbf{B} \in \mathbb{R}^{m \times n}$  and a vector  $\mathbf{y} \in \mathbb{R}^m$ , the goal is to provide conditions under which there is a unique solution  $(\mathbf{x}^*, \mathbf{u}^*)$ , where  $\mathbf{u}^*$  is  $s$ -sparse, and an algorithm for recovering it.

### 2.1 Uniqueness results for the feasibility problem

In this subsection, we study the uniqueness of solutions to the linear system accounting for the different structures the measurement matrices  $\mathbf{A}$  and  $\mathbf{B}$  can have. For more details of all the different cases we consider, we refer the reader to Appendix A. The main result of this subsection is Theorem 3 which, under a natural geometric condition based on the minimum principal

angle between the column space of  $\mathbf{A}$  and the span of  $s$  columns in  $\mathbf{B}$ , establishes a uniqueness result for the dense and sparse coding problem. Since the vector  $\mathbf{u}$  in the proposed model is sparse, we consider the classical setting of an overcomplete measurement matrix  $\mathbf{B}$  with  $n \gg m$ . The next theorem provides a uniqueness result assuming a certain direct sum representation of the space  $\mathbb{R}^m$ .

**Theorem 1** Assume that there exists at least one solution to  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$ , namely the pair  $(\mathbf{x}^*, \mathbf{u}^*)$ . Let  $S$ , with  $|S| = s$ , denote the support of  $\mathbf{u}^*$ . If  $\mathbf{B}_S$  has full column rank and  $\mathbb{R}^m = \text{Col}(\mathbf{A}) \oplus \text{Col}(\mathbf{B}_S)$ , the only unique solution to the linear system, with the condition that any feasible  $s$ -sparse vector  $\mathbf{u}$  is supported on  $S$  and any feasible  $\mathbf{x}$  is in  $\text{Ker}(\mathbf{A})^\perp$ , is  $(\mathbf{x}^*, \mathbf{u}^*)$ .

**Proof 1** Let  $(\mathbf{x}, \mathbf{u})$ , with  $\mathbf{u}$  supported on  $S$  and  $\mathbf{x}^* \in \text{Ker}(\mathbf{A})^\perp$ , be another solution pair. It follows that  $\mathbf{A}\delta_1 + \mathbf{B}_S(\delta_2)_S = \mathbf{0}$  where  $\delta_1 = \mathbf{x} - \mathbf{x}^*$  and  $\delta_2 = \mathbf{u}_S - \mathbf{u}_S^*$ . Let  $\mathbf{U} \in \mathbb{R}^{m \times r}$  and  $\mathbf{V} \in \mathbb{R}^{m \times q}$  be matrices whose columns are the orthonormal bases of  $\text{Col}(\mathbf{A})$  and  $\text{Col}(\mathbf{B}_S)$  respectively. The equation  $\mathbf{A}\delta_1 + \mathbf{B}_S(\delta_2)_S = \mathbf{0}$  can equivalently be written as  $\sum_{i=1}^r \langle \mathbf{A}\delta_1, \mathbf{U}_i \rangle \mathbf{U}_i + \sum_{i=1}^q \langle \mathbf{B}_S(\delta_2)_S, \mathbf{V}_i \rangle \mathbf{V}_i = \mathbf{0}$  with  $\mathbf{U}_i$  and  $\mathbf{V}_i$  denoting the  $i$ -th column of  $\mathbf{U}$  and  $\mathbf{V}$  respectively. More compactly, we have  $[\mathbf{U} \ \mathbf{V}] \begin{bmatrix} \{\langle \mathbf{A}\delta_1, \mathbf{U}_i \rangle\}_{i=1}^r \\ \{\langle \mathbf{B}_S(\delta_2)_S, \mathbf{V}_i \rangle\}_{i=1}^q \end{bmatrix} = \mathbf{0}$ .

Noting that the matrix  $[\mathbf{U} \ \mathbf{V}]$  has full column rank, the homogeneous problem admits the trivial solution implying that  $\mathbf{A}\delta_1 = \mathbf{0}$  and  $\mathbf{B}_S(\delta_2)_S = \mathbf{0}$ . Since  $\mathbf{B}_S$  has full column rank and  $\delta_1 \in \{\text{Ker}(\mathbf{A}) \cap \text{Ker}(\mathbf{A})^\perp\}$ , it follows that  $\delta_1 = \delta_2 = \mathbf{0}$ . Therefore,  $(\mathbf{x}^*, \mathbf{u}^*)$  is the unique solution.

The uniqueness result in the above theorem hinges on the representation of the space  $\mathbb{R}^m$  as the direct sum of the subspaces  $\text{Col}(\mathbf{A})$  and  $\text{Col}(\mathbf{B}_S)$ . We use the definition of the minimal principal angle between two subspaces, and its formulation in terms of singular values (Björck and Golub, 1973), to derive an explicit geometric condition for the uniqueness analysis of the linear system in the general case.

**Definition 2** Let  $\mathbf{U} \in \mathbb{R}^{m \times r}$  and  $\mathbf{V} \in \mathbb{R}^{m \times q}$  be matrices whose columns are the orthonormal basis of  $\text{Col}(\mathbf{A})$  and  $\text{Col}(\mathbf{B})$  respectively. The minimum principal angle between the subspaces  $\text{Col}(\mathbf{A})$  and  $\text{Col}(\mathbf{B})$  is defined as follows

$$\cos(\mu(\mathbf{U}, \mathbf{V})) = \max_{\mathbf{u} \in \text{Col}(\mathbf{U}), \mathbf{v} \in \text{Col}(\mathbf{V})} \frac{\mathbf{u}^T \mathbf{v}}{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2}, \quad (1)$$

The minimum angle  $\mu(\mathbf{U}, \mathbf{V})$  is also equal to the largest singular value of  $\mathbf{U}^T \mathbf{V}$ ,  $\cos(\mu(\mathbf{U}, \mathbf{V})) = \sigma_1(\mathbf{U}^T \mathbf{V})$ .

**Theorem 3** Assume that there exists at least one solution to  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$ , namely the pair  $(\mathbf{x}^*, \mathbf{u}^*)$ .

Let  $S$ , with  $|S| = s$ , denote the support of  $\mathbf{u}^*$ . Assume that  $\mathbf{B}_S$  has full column rank. Let  $\mathbf{U} \in \mathbb{R}^{m \times r}$  and  $\mathbf{V} \in \mathbb{R}^{m \times q}$  be matrices whose columns are the orthonormal bases of  $\text{Col}(\mathbf{A})$  and  $\text{Col}(\mathbf{B}_S)$  respectively. If  $\cos(\mu(\mathbf{U}, \mathbf{V})) = \sigma_1(\mathbf{U}^T \mathbf{V}) < 1$ , the only unique solution to the linear system, with the condition that any feasible  $s$ -sparse vector  $\mathbf{u}$  is supported on  $S$  and any feasible  $\mathbf{x}$  is in  $\text{Ker}(\mathbf{A})^\perp$ , is  $(\mathbf{x}^*, \mathbf{u}^*)$ .

**Proof 2** Consider any candidate solution pair  $(\mathbf{x}^* + \delta_1, \mathbf{u}^* + \delta_2)$ . We will prove uniqueness by showing that  $\mathbf{A}\delta_1 + \mathbf{B}_S(\delta_2)_S = \mathbf{0}$  if and only if  $\delta_1 = \mathbf{0}$  and  $(\delta_2)_S = \mathbf{0}$ . Using the orthonormal basis set  $\mathbf{U}$  and  $\mathbf{V}$ ,  $\mathbf{A}\delta_1 + \mathbf{B}_S\mathbf{u}_S$  can be represented as:  $\mathbf{A}\delta_1 + \mathbf{B}_S(\delta_2)_S = [\mathbf{U} \ \mathbf{V}] \begin{bmatrix} \mathbf{U}^T \mathbf{A}\delta_1 \\ \mathbf{V}^T \mathbf{B}_S(\delta_2)_S \end{bmatrix}$ . For simplicity of notation, let  $\mathbf{K}$  denote the block matrix:  $\mathbf{K} = [\mathbf{U} \ \mathbf{V}]$ . If we can show that the columns of  $\mathbf{K}$  are linearly independent, it follows that  $\mathbf{A}\delta_1 + \mathbf{B}_S(\delta_2)_S = \mathbf{0}$  if and only if  $\mathbf{A}\delta_1 = \mathbf{0}$  and  $\mathbf{B}_S(\delta_2)_S = \mathbf{0}$ . We now consider the matrix  $\mathbf{K}^T \mathbf{K}$  which has the following representation

$$\begin{aligned} \mathbf{K}^T \mathbf{K} &= \begin{bmatrix} [\mathbf{I}]_{r \times r} & [\mathbf{U}^T \mathbf{V}]_{r \times q} \\ [\mathbf{V}^T \mathbf{U}]_{q \times r} & [\mathbf{I}]_{q \times q} \end{bmatrix} \\ &= \begin{bmatrix} [\mathbf{I}]_{r \times r} & [\mathbf{0}]_{r \times q} \\ [\mathbf{0}]_{q \times r} & [\mathbf{I}]_{q \times q} \end{bmatrix} + \begin{bmatrix} [\mathbf{0}]_{r \times r} & [\mathbf{U}^T \mathbf{V}]_{r \times q} \\ [\mathbf{V}^T \mathbf{U}]_{q \times r} & [\mathbf{0}]_{q \times q} \end{bmatrix}. \end{aligned}$$

With the singular value decomposition of  $\mathbf{U}^T \mathbf{V}$  being  $\mathbf{U}^T \mathbf{V} = \mathbf{Q}\mathbf{\Sigma}\mathbf{R}^T$ , the last matrix in the above representation has the following equivalent form

$$\begin{bmatrix} \mathbf{0} & \mathbf{U}^T \mathbf{V} \\ \mathbf{V}^T \mathbf{U} & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{Q} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{\Sigma} \\ \mathbf{\Sigma} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{Q} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{bmatrix}^T. \quad \text{It}$$

now follows that  $\begin{bmatrix} \mathbf{0} & \mathbf{U}^T \mathbf{V} \\ \mathbf{V}^T \mathbf{U} & \mathbf{0} \end{bmatrix}$  is similar to the matrix

$$\begin{bmatrix} \mathbf{0} & \mathbf{\Sigma} \\ \mathbf{\Sigma} & \mathbf{0} \end{bmatrix}. \quad \text{Hence, the nonzero eigenvalues of } \mathbf{K}^T \mathbf{K}$$

are  $1 \pm \sigma_i$ ,  $1 \leq i \leq \min(p, q)$ , with  $\sigma_i$  denoting the  $i$ -th largest singular value of  $\mathbf{U}^T \mathbf{V}$ . Using the assumption  $\sigma_1 < 1$  results the bound  $\lambda_{\min}(\mathbf{K}^T \mathbf{K}) > 0$ . It follows that the columns of  $\mathbf{K}$  are linearly independent, and hence  $\mathbf{A}\delta_1 = \mathbf{0}$  and  $\mathbf{B}_S(\delta_2)_S = \mathbf{0}$ . Since  $\mathbf{B}_S$  is full column rank and  $\delta_1 \in \{\text{Ker}(\mathbf{A}) \cap \text{Ker}(\mathbf{A})^\perp\}$ , it follows that  $\delta_1 = \mathbf{0}$  and  $(\delta_2)_S = \mathbf{0}$ . This concludes the proof.

A restrictive assumption of the above theorem is that the support of the sought-after  $s$ -sparse solution  $\mathbf{u}^*$  is known. We can remove this assumption by considering  $\text{Col}(\mathbf{A})$  and  $\text{Col}(\mathbf{B}_T)$  where  $T$  is an arbitrary subset of  $\{1, 2, \dots, n\}$  with  $|T| = s$ . More precisely, we state the following corollary whose proof is similar to the proof of Theorem 3.

**Corollary 4** Assume that there exists at least one solution to  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$ , namely the pair  $(\mathbf{x}^*, \mathbf{u}^*)$ . Let  $S$ , with  $|S| = s$ , denote the support of  $\mathbf{u}^*$  and  $T$  be an arbitrary subset of  $\{1, 2, \dots, n\}$  with  $|T| \leq s$ . Assume

that any  $2s$  columns of  $\mathbf{B}$  are linearly independent. Let  $\mathbf{U} \in \mathbb{R}^{m \times p}$  and  $\mathbf{V} \in \mathbb{R}^{m \times q}$  be matrices whose columns are the orthonormal bases of  $\text{Col}(\mathbf{A})$  and  $\text{Col}(\mathbf{B}_{S \cup T})$  respectively. If  $\mu(\mathbf{U}, \mathbf{V}) = \sigma_1(\mathbf{U}^T \mathbf{V}) < 1$ , holds for all choices of  $T$ , the only unique solution to the linear system is  $(\mathbf{x}^*, \mathbf{u}^*)$  with the condition that any feasible  $\mathbf{u}$  is  $s$ -sparse and any feasible  $\mathbf{x}$  is in  $\text{Ker}(\mathbf{A})^\perp$ ,

Of interest is the identification of simple conditions such that  $\sigma_1(\mathbf{U}^T \mathbf{V}) < 1$ . The following theorem proposes one such condition to establish uniqueness of the dense and sparse coding problem.

**Theorem 5** Assume that there exists at least one solution to  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$ , namely the pair  $(\mathbf{x}^*, \mathbf{u}^*)$ . Let  $S$ , with  $|S| = s$ , denote the support of  $\mathbf{u}^*$ . Assume that  $\mathbf{B}_S$  has full column rank. Let  $\mathbf{U} \in \mathbb{R}^{m \times r}$  and  $\mathbf{V} \in \mathbb{R}^{m \times q}$  be matrices whose columns are the orthonormal bases of  $\text{Col}(\mathbf{A})$  and  $\text{Col}(\mathbf{B}_S)$  respectively. Let  $\max_{i,j} |(\mathbf{U}^T \mathbf{V})_{i,j}| = \mu$ . If  $s < \frac{1}{\sqrt{r}\mu}$ , the only unique solution to the linear system, with the condition that any feasible  $s$ -sparse vector  $\mathbf{u}$  is supported on  $S$  and any feasible  $\mathbf{x}$  is in  $\text{Ker}(\mathbf{A})^\perp$ , is  $(\mathbf{x}^*, \mathbf{u}^*)$ .

**Proof 3** It suffices to show that  $\sigma_1 < 1$ . Noting that  $\sigma_1 = \|\mathbf{U}^T \mathbf{V}\|_2$ , we use the following matrix norm inequality  $\|\mathbf{U}^T \mathbf{V}\|_2 \leq \sqrt{r} \|\mathbf{U}^T \mathbf{V}\|_\infty$  as follows:  $\sigma_1 \leq \sqrt{r} \|\mathbf{U}^T \mathbf{V}\|_\infty \leq \sqrt{r} \mu s < 1$ .

The constant  $\mu$  is the coherence of the matrix  $\mathbf{U}^T \mathbf{V}$  (Donoho et al., 2005; Tropp, 2004). The above result states that if the mutual coherence of  $\mathbf{U}^T \mathbf{V}$  is small, we can accommodate increased sparsity of the underlying signal component  $\mathbf{u}^*$ . We note that, up to a scaling factor,  $\sigma_1(\mathbf{U}^T \mathbf{V})$  is the block coherence of  $\mathbf{U}$  and  $\mathbf{V}$  (Eldar et al., 2010). However, unlike the condition in (Eldar et al., 2010), we don't restrict the dictionaries  $\mathbf{A}$  and  $\mathbf{B}$  to have linearly independent columns. In the next subsection, we propose a convex program to recover the dense and sparse vectors. Theorem 8 establishes uniqueness and complexity results for the proposed optimization program.

## 2.2 Dense and sparse recovery via convex optimization

Given that the dense and sparse coding problem seeks a dense vector  $\mathbf{x}^*$  and a sparse solution  $\mathbf{u}^*$ , with measurements given as  $\mathbf{y} = \mathbf{A}\mathbf{x}^* + \mathbf{B}\mathbf{u}^*$ , we propose the following convex optimization program

$$\min_{\mathbf{x}, \mathbf{u}} \|\mathbf{A}\mathbf{x}\|_2^2 + \|\mathbf{u}\|_1 \quad \text{s.t.} \quad \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}. \quad (2)$$

In this section, we show that, under certain conditions, the above minimization problem admits a unique solution. Our proof is a non-trivial adaptation of the

existing analysis in (Kueng and Gross, 2014) for the anisotropic compressive sensing problem. This analysis is based on a single measurement matrix and can not be directly applied to our scenario. Let  $\mathbf{a}_1, \dots, \mathbf{a}_m$  be a sequence of zero-mean i.i.d random vectors drawn from some distribution  $F$  on  $\mathbb{R}^p$  and let  $\mathbf{b}_1, \dots, \mathbf{b}_m$  be a sequence of zero-mean i.i.d random vectors drawn from some distribution  $G$  on  $\mathbb{R}^n$ . We can eliminate the dense component in the linear constraint by projecting the vector  $\mathbf{y}$  onto the orthogonal complement of  $\text{Col}(\mathbf{A})$  to obtain  $\mathcal{P}_{\text{Col}(\mathbf{A})^\perp}(\mathbf{y}) = \mathcal{P}_{\text{Col}(\mathbf{A})^\perp}(\mathbf{B}\mathbf{u})$ . With this, the matrix  $\mathcal{P}_{\text{Col}(\mathbf{A})^\perp}(\mathbf{B})$  is central in the analysis to follow. We define the matrix  $\mathbf{C} \mathbf{C}^T = \frac{1}{\sqrt{m}} \sum_{i=1}^m \mathbf{e}_i \mathbf{c}_i^T$  where  $\mathbf{c}_i = [\mathcal{P}_{\text{Col}(\mathbf{A})^\perp}(\mathbf{B})]^T \mathbf{e}_i$  denotes the  $i$ -th measurement vector corresponding to a row of this matrix. Further technical discussion on the matrix  $\mathbf{C}$  is deferred to **Appendix B**. We use the measurement matrix  $\mathbf{C}$  introduced above and adapt the anisotropic compressive sensing theory in (Kueng and Gross, 2014) to analyze uniqueness of the proposed program. Below, we give brief background to this theory highlighting important assumptions and results following the notation closely therein.

**Anisotropic compressive sensing:** Given a sequence of zero-mean i.i.d random vectors  $\mathbf{d}_1, \dots, \mathbf{d}_m$  drawn from some distribution  $F$  on  $\mathbb{R}^n$ , with measurements  $\mathbf{y} = \mathbf{D}\mathbf{u}^*$ , the anisotropic compressive sensing problem studies the following optimization program

$$\min_{\mathbf{u}} \|\mathbf{u}\|_1 \quad \text{s.t.} \quad \mathbf{y} = \mathbf{D}\mathbf{u}, \quad (3)$$

where  $\mathbf{D} = \frac{1}{\sqrt{m}} \sum_{i=1}^m \mathbf{e}_i \mathbf{d}_i^T$  and  $\mathbf{u}^*$  is the sought-out sparse solution. The analysis makes three important assumptions.

**Completeness:** The covariance matrix  $\mathbf{\Sigma}$  is invertible with condition number denoted by  $\kappa$ .

**Incoherence:** The incoherence parameter is the smallest number  $\nu$  such that

$$\max_{1 \leq i \leq n} |\langle \mathbf{d}, \mathbf{e}_i \rangle|^2 \leq \nu \quad \text{and} \quad \max_{1 \leq i \leq n} |\langle \mathbf{d}, E[\mathbf{c}\mathbf{c}^*]^{-1} \mathbf{e}_i \rangle|^2 \leq \nu \quad (4)$$

hold almost surely.

**Conditioning of the covariance matrix:** We start with the following definition of the  $s$ -sparse condition number restated from (Kueng and Gross, 2014).

**Definition 6** (Kueng and Gross, 2014) The largest and smallest  $s$ -sparse eigenvalue of a matrix  $\mathbf{X}$  are given by

$$\lambda_{\max}(s, \mathbf{X}) := \max_{\mathbf{v}, \|\mathbf{v}\|_0 \leq s} \frac{\|\mathbf{X}\mathbf{v}\|_2}{\|\mathbf{v}\|_2}$$

$$\lambda_{\min}(s, \mathbf{X}) := \min_{\mathbf{v}, \|\mathbf{v}\|_0 \leq s} \frac{\|\mathbf{X}\mathbf{v}\|_2}{\|\mathbf{v}\|_2}.$$

The  $s$ -sparse condition number of  $\mathbf{X}$  is  $\text{cond}(s, \mathbf{X}) = \frac{\lambda_{\max}(s, \mathbf{X})}{\lambda_{\min}(s, \mathbf{X})}$ .

Given these assumptions, the main result in (Kueng and Gross, 2014) reads

**Theorem 7** (Kueng and Gross, 2014) *With  $\kappa_s = \max\{\text{cond}(s, \mathbf{\Sigma}), \text{cond}(s, \mathbf{\Sigma}^{-1})\}$  let  $\mathbf{u} \in \mathbb{C}^n$  be an  $s$ -sparse vector and let  $\omega \geq 1$ . If the number of measurements fulfills  $m \geq C\kappa_s \nu \omega^2 s \log n$ , then the solution  $\mathbf{u}$  of the convex program (3) is unique and equal to  $\mathbf{u}^*$  with probability at least  $1 - e^{-\omega}$ .*

The proof of Theorem 7 is based on the dual certificate approach. The idea is to first propose a dual certificate vector  $\mathbf{v}$  with sufficient conditions that ensure uniqueness of the minimization problem. It then remains to construct the dual certificate satisfying the conditions. We seek a similar result for the uniqueness of the convex program corresponding to the dense and sparse coding model. However, the standard analysis can not be directly applied since it only considers a single measurement matrix. This requires us to analyze the matrix  $\mathbf{C}$  introduced earlier.

The anisotropic compressive sensing analysis in (Kueng and Gross, 2014) assumes the following conditions on the dual certificate  $\mathbf{v}$

$$\|\mathbf{v}_S - \text{sgn}(\mathbf{u}_S^*)\|_2 \leq \frac{1}{4} \quad \text{and} \quad \|\mathbf{v}_{S^\perp}\|_\infty \leq \frac{1}{4}. \quad (5)$$

The following condition follows from the assumptions in Theorem 7

$$\|\mathbf{\Delta}_S\|_2 \leq 2\|\mathbf{\Delta}_{S^\perp}\|_2, \quad (6)$$

where  $\mathbf{\Delta} \in \text{Ker}(\mathbf{D})$ . The conditions (5) and (6) will be used in the proof of our main result. The main part of the technical analysis in (Kueng and Gross, 2014) is using the assumptions in Theorem 7 and showing that the above conditions (5) and (6) hold with high probability.

**Main result:** Using the the background discussed above, we assume completeness, incoherence, and conditioning of the covariance matrix  $\mathbf{\Sigma}$ . Our main result is stated below.

**Theorem 8** *Assume that there exists at least one solution to  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$ , namely the pair  $(\mathbf{x}^*, \mathbf{u}^*)$ . Let  $\omega \geq 1$  and define  $\kappa_s = \max\{\text{cond}(s, \mathbf{\Sigma}), \text{cond}(s, \mathbf{\Sigma}^{-1})\}$ . Assume the two conditions*

$$\|\mathbf{B}_S^T \mathbf{A}\| \leq \frac{1}{32\|\mathbf{x}^*\|_2}, \quad \|\mathbf{B}_{S^\perp}^T \mathbf{A}\|_\infty \leq \frac{1}{32\|\mathbf{x}^*\|_\infty}. \quad (7)$$

*If the number of measurements fulfills  $m \geq C\kappa_s \nu \omega^2 s \log n$ , then the solution of the convex program (2) is unique and equal to  $(\mathbf{x}^*, \mathbf{u}^*)$  with probability at least  $1 - e^{-\omega}$ .*

**Proof sketch 1** *Consider a feasible solution pair  $(\mathbf{x}^* + \delta_1, \mathbf{u}^* + \delta_2)$  and let the function  $f(\mathbf{x}, \mathbf{u})$  denote the objective in the optimization program. The idea of the proof is to show that any feasible solution is not minimal in the objective value,  $f(\mathbf{x}^* + \delta_1, \mathbf{u}^* + \delta_2) > f(\mathbf{x}, \mathbf{u})$ . Using duality of the  $\ell_1$  norm and characterization of the subgradient  $\mathbf{\Lambda}$  of the  $\ell_1$  norm, we first show that  $f(\mathbf{x}^* + \delta_1, \mathbf{u}^* + \delta_2) > f(\mathbf{x}^*, \mathbf{u}^*) + \langle \text{sgn}(\mathbf{u}_S^*) + \mathbf{\Lambda} - \mathbf{v} - 2\mathbf{B}^T \mathbf{A} \mathbf{x}^*, \delta_2 \rangle$  where  $\mathbf{v} \in \text{Col}(\mathbf{C}^T)$ , with  $\mathbf{C} = \mathbf{B} - \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{B}$  denoting the dual certificate. It then remains to show that the term  $\langle \text{sgn}(\mathbf{u}_S^*) + \mathbf{\Lambda} - \mathbf{v} - 2\mathbf{B}^T \mathbf{A} \mathbf{x}^*, \delta_2 \rangle$  is positive. To show this, we further analyze this term and make use of the assumptions of the theorem, the dual certificate conditions (5), and the deviation inequality in (6) to arrive at the desired result. For a complete proof, see **Appendix B**.*

**Complexity compared to  $\ell_1$  minimization:** The sample complexity of solving the convex program corresponding to the dense and sparse coding problem is larger than that of  $\ell_1$  minimization for the compressive sensing problem. Essentially, the constants  $\kappa_s$  and  $\nu$  in our analysis are expected to scale with  $p + n$ , in contrast to the compressive sensing analysis where they scale with  $n$ .

## 3 Experiments

### 3.1 Phase transition curves

We generate *phase transition curves* and present how the success rate of the recovery, using the proposed model, changes under different scenarios. To generate the data, we fix the number of columns of  $\mathbf{B}$  to be  $n = 100$ . Then, we vary the sampling ratio  $\sigma = \frac{m}{n+p} \in [0.05, 0.95]$  and the sparsity ratio  $\rho = \frac{s}{m}$  in the same range. The sensing matrix in our model is  $[\mathbf{A} \ \mathbf{B}]$ , hence the apparent difference in the definition of  $\sigma$  compared to “traditional” compressive sensing. In the case where we revert to the compressive sensing scenario ( $p = 0$ ), the ratios coincide.

We generate random matrices  $\mathbf{A} \in \mathbb{R}^{m \times p}$  and  $\mathbf{B} \in \mathbb{R}^{m \times n}$  whose columns have expected unit norm. The vector  $\mathbf{u} \in \mathbb{R}^n$  has  $s$  randomly chosen indices, whose entries are drawn according to a standard normal distribution, and  $\mathbf{x} \in \mathbb{R}^p$  is generated as follows: we generate a random vector  $\gamma \in \mathbb{R}^m$ , and then construct  $\mathbf{x} = \mathbf{A}^T \gamma$ . The construction ensures that  $\mathbf{x}$  does not belong in the null space of  $\mathbf{A}$ , and hence ignores trivial solutions with respect to this dense component. We normalize both  $\mathbf{x}$  and  $\mathbf{u}$  to have unit norm, and generate the measurement vector  $\mathbf{y} \in \mathbb{R}^m$  as  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$ . We solve the convex optimization problem in (2) to obtain the numerical solution pair  $(\hat{\mathbf{x}}, \hat{\mathbf{u}})$  using CVXPY, and

register a successful recovery if both  $\frac{\|\hat{\mathbf{x}} - \mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \epsilon$  and  $\frac{\|\hat{\mathbf{u}} - \mathbf{u}\|_2}{\|\mathbf{u}\|_2} \leq \epsilon$ , with  $\epsilon = 10^{-3}$ . For each choice of  $\sigma$  and  $\rho$  we average 100 independent runs to estimate the success rate.

Figure 1 shows the phase transition curves for  $p \in \{0.1m, 0.5m\}$  to highlight different ratios between  $p$  and  $n$ . We observe that increasing  $p$  leads to a deterioration in performance. This is expected, as this creates a greater *overlap* on the spaces spanned by  $\mathbf{A}$  and  $\mathbf{B}$ . We can view our model as explicitly modeling the noise of the system. In such a case, the number of columns of  $\mathbf{A}$  explicitly encodes the complexity of the noise model: as  $p$  increases, so does the span of the noise space.

Extending the signal processing interpretation, note that we model the noise signal  $\mathbf{x}$  as a dense vector, which can be seen as encoding smooth areas of the signal that correspond to *low-frequency* components. On the contrary, the signal  $\mathbf{u}$  has, by construction, a sparse structure, containing *high-frequency* information, an interpretation that will be further validated in the next subsection. Further numerical experiments comparing the dense and sparse coding model to noisy compressive sensing can be found in **Appendix C**.

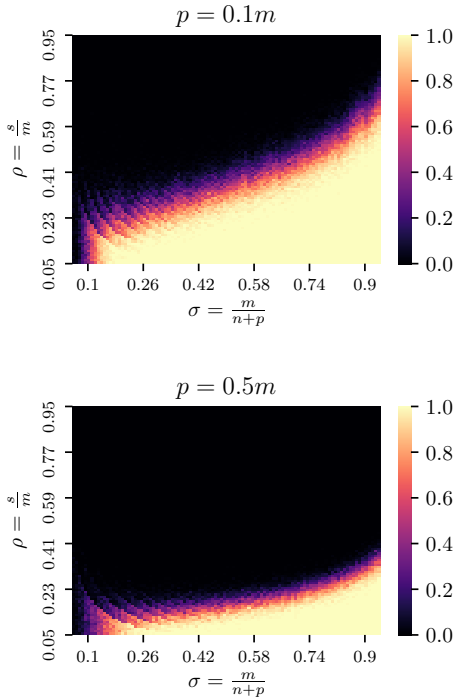


Figure 1: Phase transition curves for  $p = 0.1m$  (top) and  $p = 0.5m$  (bottom).

### 3.2 Classification and image denoising

We formulate the dense and sparse dictionary learning problem as minimizing the objective  $\min_{\mathbf{A}, \mathbf{B}, \{\mathbf{x}^j\}_{j=1}^J, \{\mathbf{u}^j\}_{j=1}^J} \sum_{j=1}^J \frac{1}{2} \|\mathbf{y}^j - \mathbf{A}\mathbf{x}^j - \mathbf{B}\mathbf{u}^j\|_2^2 + \frac{1}{2\lambda_x} \|\mathbf{A}\mathbf{x}^j\|_2^2 + \lambda_u \|\mathbf{u}^j\|_1$ , where  $J$  is the number of images,  $\lambda_x$  controls the smoothness of  $\mathbf{A}\mathbf{x}^j$  and  $\lambda_u$  controls the degree of sparsity. Based on the objective, we use deep unfolding to construct a unfolding neural network (Tolooshams et al., 2020; Gregor and Lecun, 2010), which we term the dense and sparse autoencoder (DenSaE), tailored to learning the dictionaries from the dense and sparse model. The encoder maps  $\mathbf{y}^j$  into a dense vector  $\mathbf{x}_T^j$  and a sparse one  $\mathbf{u}_T^j$  by unfolding  $T$  proximal gradient iterations. The decoder reconstructs the image. For classification, we use  $\mathbf{u}_T$  and  $\mathbf{x}_T$  as inputs to a linear classifier  $\mathbf{C}$  that maps them to the predicted class  $\hat{\mathbf{q}}$ . We learn the dictionaries  $\mathbf{A}$  and  $\mathbf{B}$ , as well as the classifier  $\mathbf{C}$ , by minimizing the weighted reconstruction (AE) and classification (Logistic) loss (i.e.,  $(1 - \beta)$  AE +  $\beta$  Logistic). Figure 2 shows the DenSaE architecture (for details see **Appendix D**).

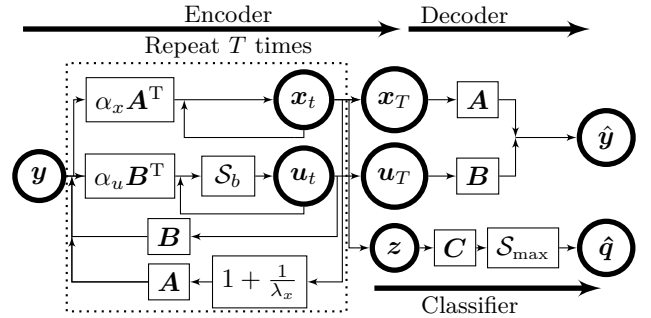


Figure 2: DenSaE. The vector  $\mathbf{z}$  is normalized stacked features with a column of 1 (i.e.,  $\mathbf{z}^T = [\mathbf{1}, \frac{[\mathbf{x}_T, \mathbf{u}_T]}{\|[\mathbf{x}_T, \mathbf{u}_T]\|}]$ ),  $\mathcal{S}_b$  is soft-thresholding, and  $\mathcal{S}_{\max}$  is softmax.

We examined the following questions

- How do the discriminative, reconstruction, and denoising capabilities change as we vary the number of filters in  $\mathbf{A}$  vs.  $\mathbf{B}$ ?
- What is the performance of DenSaE compared to sparse coding networks?
- What data characteristics does the model capture?

As baselines, we trained two variants,  $\text{CSCNet}_{\text{hyp}}^{\text{tied}}$  and  $\text{CSCNet}_{\text{LS}}^{\text{tied}}$ , of CSCNet (Simon and Elad, 2019), an architecture tailored to dictionary learning for the sparse coding problem. In  $\text{CSCNet}_{\text{hyp}}^{\text{tied}}$ , the bias is a shared hyper-parameter. In  $\text{CSCNet}_{\text{LS}}^{\text{tied}}$ , we learn a different bias for each filter by minimizing the reconstruction loss. When the dictionaries are non-convolutional, we call the network SCNet.

Table 1: DenSaE’s performance on MNIST test dataset from both disjoint (D) and joint ( $J_\beta$ ) training.

		SCNet <sup>LS</sup> <sub>hyp</sub>	SCNet <sup>tied</sup> <sub>hyp</sub>	$\frac{5\mathbf{A}}{395\mathbf{B}}$	$\frac{25\mathbf{A}}{375\mathbf{B}}$	$\frac{200\mathbf{A}}{200\mathbf{B}}$
$\frac{\mathbf{A}}{\mathbf{A+B}}$ model		-	-	1.25	6.25	50
D	Acc.	94.16	98.32	98.18	98.18	96.98
	Rec.	1.95	6.80	6.83	6.30	3.04
	$\frac{\mathbf{A}}{\mathbf{A+B}}$ class	-	-	0	0	0
	$\frac{\mathbf{A}}{\mathbf{A+B}}$ rec.	-	-	8	28	58
$J_{0.75}$	Acc.	96.91	98.18	98.19	98.23	97.64
	Rec.	2.17	1.24	0.75	1.11	<b>0.51</b>
	$\frac{\mathbf{A}}{\mathbf{A+B}}$ class	-	-	8	8	84
	$\frac{\mathbf{A}}{\mathbf{A+B}}$ rec.	-	-	8	36	8
$J_1$	Acc.	96.06	98.59	<b>98.61</b>	98.56	98.40
	Rec.	71.20	47.70	32.61	30.20	25.57
	$\frac{\mathbf{A}}{\mathbf{A+B}}$ class	-	-	16	46	42
	$\frac{\mathbf{A}}{\mathbf{A+B}}$ rec.	-	-	0	2	4

### 3.2.1 DenSaE strikes a balance between discriminative capability and reconstruction

We study the case when DenSaE is trained on the MNIST dataset for joint reconstruction and classification purposes. We show a) how the explicit imposition of sparse and dense representations in DenSaE helps to balance discriminative and representation power, and b) that DenSaE outperforms SCNet. We warm start the training of the classifier using dictionaries obtained first training the autoencoder, i.e., with  $\beta = 0$ .

**Characteristics of the representations  $\mathbf{x}_T$  and  $\mathbf{u}_T$ :** To evaluate the discriminative power of the representations learned by only training the autoencoder, we first trained the classifier *given* the representations (i.e., first train  $\mathbf{A}$  and  $\mathbf{B}$  with  $\beta = 0$ , then train  $\mathbf{C}$  with  $\beta = 1$ ). We call this disjoint training. The first four rows of section D from Table 1 show, respectively, the classification accuracy (Acc.),  $\ell_2$  reconstruction loss (Rec.), and the relative contributions, expressed as a percentage, of the dense or sparse representations to classification and reconstruction for disjoint training. Each col of  $[\mathbf{A} \ \mathbf{B}]$ , and of  $\mathbf{C}$ , corresponds to either a dense or a sparse feature. For reconstruction, we find the indices of the 50 most important columns and report the proportion of these that represent dense features. For each of the 10 classes (rows of  $\mathbf{C}$ ), we find the indices of the 5 most important columns (features) and compute the proportion of the total of 50 indices that represent dense features. The first row of Table 1 shows the proportion of rows of  $[\mathbf{A} \ \mathbf{B}]$  that represent dense features. Comparing this row, respectively to the third and fourth row of section D reveals the importance of  $\mathbf{x}$  for reconstruction, and of  $\mathbf{u}$  for classification. Indeed, the first two rows of section D show that, as the proportion of dense features increases, DenSaE

gains reconstruction capability but results in a lower classification accuracy. Moreover, in DenSaE, the most important features in classification are all from  $\mathbf{B}$ , and the contribution of  $\mathbf{A}$  in reconstruction is greater than its percentage in the model, which clearly demonstrates that dense and sparse coding autoencoders balance discriminative and representation power.

The table also shows that DenSaE outperforms SCNet<sup>LS</sup><sub>hyp</sub> in classification and SCNet<sup>tied</sup><sub>hyp</sub> in reconstruction. We observed that in the absence of noise, training SCNet<sup>LS</sup><sub>hyp</sub> results in dense features with negative biases, hence, making its performance close to DenSaE with large number of atoms in  $\mathbf{A}$ . We see that SCNet<sup>LS</sup><sub>hyp</sub> in absence of a supervised classification loss fails to learn discriminative features useful for classification. On the other hand, enforcing sparsity in SCNet<sup>tied</sup><sub>hyp</sub> suggests that sparse representations are useful for classification.

**How do roles of  $\mathbf{x}_T$  and  $\mathbf{u}_T$  change as we vary  $\beta$  in joint training?:** In joint training of the autoencoder and the classifier, it is natural to expect that the reconstruction loss should increase compared to disjoint training. This is indeed the case for SCNet<sup>LS</sup><sub>hyp</sub>; as we go from disjoint to joint training and as  $\beta$  increases (Table 1, sections labeled J), the reconstruction loss increases and classification accuracy has an overall increase. However, for  $\beta < 1$ , joint training of both networks that enforce some sparsity on their representations, SCNet<sup>tied</sup><sub>hyp</sub> and DenSaE, improves reconstruction and classification. Moreover, as we increase the importance of classification loss (i.e., increase  $\beta$ ), the contribution of dense representations decreases in reconstruction and increases in discrimination.

For purely discriminative training ( $\beta = 1$ ), DenSaE outperforms both SCNet<sup>LS</sup><sub>hyp</sub> and SCNet<sup>tied</sup><sub>hyp</sub> in classification accuracy and representation capability. We speculate that this likely results from the fact that, by construction the encoder from DenSaE seeks to produce two sets of representations: namely a dense one, mostly important for reconstruction and a sparse one, useful for classification. In some sense, the dense component acts as a prior that promotes good reconstruction. More detailed results can be found in **Appendix D**.

**Remarks:** As our network is non-convolutional, we do not compare it to the state-of-the-art, a convolutional network. We do not compare our results with the network in (Rolfe and LeCun, 2013) as that work does not report reconstruction loss and it involves a sparsity enforcing loss that change the learning behaviour.

### 3.2.2 Denoising

We trained DenSaE for supervised image denoising when  $\beta = 0$  using BSD432 and tested it on BSD68 (Martin et al., 2001) (see **Appendix D** for details). We



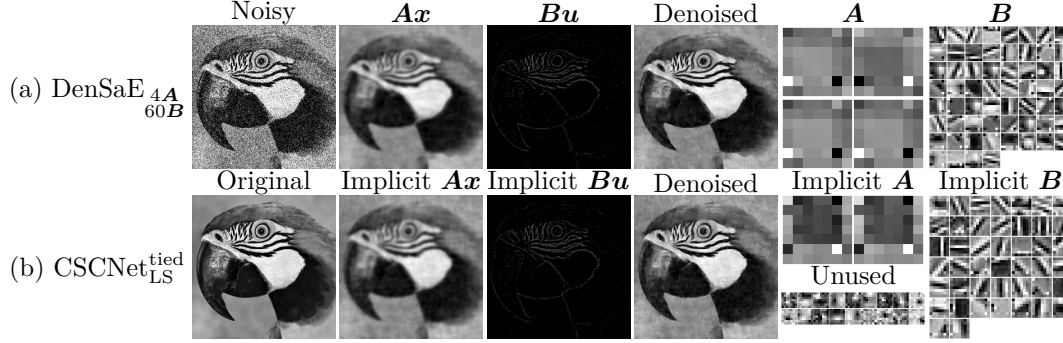


Figure 3: Visualization of a test image for  $\tau = 50$ . (a) DenSaE ( $4\mathbf{A}, 60\mathbf{B}$ ), (b) CSCNet<sub>LS</sub><sup>tied</sup>.

Table 2: DenSaE’s denoising performance on test BSD68 as the ratio of filters in  $\mathbf{A}$  and  $\mathbf{B}$  changes.

$\tau$	1A63B	4A60B	8A56B	16A48B	32A32B
15	<b>30.21</b>	30.18	30.18	30.14	29.89
25	<b>27.70</b>	<b>27.70</b>	27.65	27.56	27.26
50	<b>24.81</b>	<b>24.81</b>	24.43	24.44	23.68
75	23.31	<b>23.33</b>	23.09	22.09	20.09

Table 3: DenSaE vs. CSCNet on test BSD68.

$\tau$	DenSaE	CSCNet <sub>hyp</sub> <sup>tied</sup>	CSCNet <sub>LS</sub> <sup>tied</sup>
15	30.21	30.12	<b>30.34</b>
25	27.70	27.51	<b>27.75</b>
50	<b>24.81</b>	24.54	<b>24.81</b>
75	<b>23.33</b>	22.83	23.32

varied the ratio of number of filters in  $\mathbf{A}$  and  $\mathbf{B}$  as the overall number of filters was kept constant. We evaluate the model in the presence of Gaussian noise with standard deviation of  $\tau = \{15, 25, 50, 75\}$ .

**Ratio of number of filters in  $\mathbf{A}$  and  $\mathbf{B}$ :** Table 2 shows that the smaller the number of filters associated with  $\mathbf{A}$ , the better DenSaE can denoise images. We hypothesize that this is a direct consequence of our findings from **Section 2** that the smaller the number of columns of  $\mathbf{A}$ , the easier the recovery  $\mathbf{x}$  and  $\mathbf{u}$ .

**Dense and sparse coding vs. sparse coding:** Table 3 shows that DenSaE (best network from Table 2) denoises images better than CSCNet<sub>hyp</sub><sup>tied</sup>, suggesting that the dense and sparse coding model represents images better than sparse coding.

**Dictionary characteristics:** Figure 3(a) shows the decomposition of a noisy test image ( $\tau = 50$ ) by DenSaE. The figure demonstrates that  $\mathbf{Ax}$  captures low-frequency content while  $\mathbf{Bu}$  captures high-frequency details (edges). This is corroborated by the smoothness of the filters associated with  $\mathbf{A}$ , and the Gabor-like nature of those associated with  $\mathbf{B}$  (Mehrotra et al., 1992). We observed similar performance when we tuned  $\lambda_x$ , and found that, as  $\lambda_x$  decreases,  $\mathbf{Ax}$  captures a lower frequencies, and  $\mathbf{Bu}$  a broader range.

**CSCNet implicitly learns  $\mathbf{Ax} + \mathbf{Bu}$  model in the**

**presence of noise:** We observed that CSCNet<sub>LS</sub><sup>tied</sup> comprises three groups of filters: one with small bias, one with intermediate ones, and a third with large values (see **Appendix D** for bias visualizations). We found that the feature maps associated with the large bias values are all zero. Moreover, the majority of features are associated with intermediate bias values, and are sparse, in contrast to the small number of feature maps with small bias values, which are dense. These observations suggest that *autoencoder implementing the sparse coding model ( $\mathbf{y} = \mathbf{Bu}$ ), when learning the biases by minimizing reconstruction error, implicitly perform two functions*. First, they select the optimal number of filters. Second, they partition the filters into two groups: one that yields a dense representation of the input, and another that yields a sparse one. In other words, the architectures trained in this manner *implicitly learn the dense and sparse coding model ( $\mathbf{y} = \mathbf{Ax} + \mathbf{Bu}$ )*. Figure 3(b) shows the filters.

## 4 Conclusions

This paper proposed a novel dense and sparse coding model for a flexible representation of a signal as  $\mathbf{y} = \mathbf{Ax} + \mathbf{Bu}$ . Our first result gives a verifiable condition that guarantees uniqueness of the model. Our second result uses tools from RIPless compressed sensing to show that, with sufficiently many linear measurements, a convex program with  $\ell_1$  and  $\ell_2$  regularizations can recover the components  $\mathbf{x}$  and  $\mathbf{u}$  uniquely with high probability. Numerical experiments on synthetic data confirm our observations.

We proposed a dense and sparse autoencoder, DenSaE, tailored to dictionary learning for the  $\mathbf{Ax} + \mathbf{Bu}$  model. DenSaE, naturally decomposing signals into low- and high-frequency components, provides a balance between learning dense representations that are useful for reconstruction and discriminative sparse representations. We showed the superiority of DenSaE to sparse autoencoders for data reconstruction and its competitive performance in classification.



## References

- M. Aharon, M. Elad, and A. Bruckstein, “ $k$ -svd: An algorithm for designing overcomplete dictionaries for sparse representation,” *IEEE Transactions on signal processing*, vol. 54, no. 11, pp. 4311–22, 2006.
- J. Mairal, F. Bach, and J. Ponce, “Task-driven dictionary learning,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 4, pp. 791–804, 2011.
- B. A. Olshausen and D. J. Field, “Sparse coding with an overcomplete basis set: A strategy employed by v1?” *Vision research*, vol. 37, no. 23, pp. 3311–3325, 1997.
- C. Garcia-Cardona and B. Wohlberg, “Convolutional dictionary learning: A comparative review and new algorithms,” *IEEE Transactions on Computational Imaging*, vol. 4, no. 3, pp. 366–81, Sep. 2018.
- K. Gregor and Y. Lecun, “Learning fast approximations of sparse coding,” in *Proc. International Conference on Machine Learning (ICML)*, 2010, pp. 399–406.
- J. T. Rolfe and Y. LeCun, “Discriminative recurrent sparse auto-encoders,” *arXiv preprint arXiv:1301.3775*, 2013.
- D. Simon and M. Elad, “Rethinking the csc model for natural images,” in *Proc. Advances in Neural Information Processing Systems*, 2019, pp. 2271–2281.
- B. Tolooshams, S. Dey, and D. Ba, “Deep residual autoencoders for expectation maximization-inspired dictionary learning,” *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2020.
- P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, P.-A. Manzagol, and L. Bottou, “Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion,” *Journal of machine learning research*, vol. 11, no. 12, 2010.
- B. Epstein, R. Meir, and T. Michaeli, “Joint autoencoders: a flexible meta-learning framework,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2018, pp. 494–509.
- B. Epstein and R. Meir, “Generalization bounds for unsupervised and semi-supervised learning with autoencoders,” *arXiv preprint arXiv:1902.01449*, 2019.
- J. Sulam, A. Aberdam, A. Beck, and M. Elad, “On multi-layer basis pursuit, efficient algorithms and convolutional neural networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2019.
- J. Zazo, B. Tolooshams, and D. Ba, “Convolutional dictionary learning in hierarchical networks,” in *Proc. 2019 IEEE 8th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*. IEEE, 2019, pp. 131–135.
- D. L. Donoho and X. Huo, “Uncertainty principles and ideal atomic decomposition,” *IEEE transactions on information theory*, vol. 47, no. 7, pp. 2845–2862, 2001.
- M. Elad and A. M. Bruckstein, “A generalized uncertainty principle and sparse representation in pairs of bases,” *IEEE Transactions on Information Theory*, vol. 48, no. 9, pp. 2558–2567, 2002.
- D. L. Donoho and M. Elad, “Optimally sparse representation in general (nonorthogonal) dictionaries via  $\ell_1$  minimization,” *Proc. the National Academy of Sciences*, vol. 100, no. 5, pp. 2197–2202, 2003.
- M. Soltani and C. Hegde, “Fast algorithms for demixing sparse signals from nonlinear observations,” *IEEE Transactions on Signal Processing*, vol. 65, no. 16, pp. 4209–4222, 2017.
- C. Studer, P. Kuppinger, G. Pope, and H. Bolcskei, “Recovery of sparsely corrupted signals,” *IEEE Transactions on Information Theory*, vol. 58, no. 5, pp. 3115–3130, 2011.
- C. Studer and R. G. Baraniuk, “Stable restoration and separation of approximately sparse signals,” *Applied and Computational Harmonic Analysis*, vol. 37, no. 1, pp. 12–35, 2014.
- E. J. Candès, X. Li, Y. Ma, and J. Wright, “Robust principal component analysis?” *Journal of the ACM (JACM)*, vol. 58, no. 3, pp. 1–37, 2011.
- L. Lian, A. Liu, and V. K. Lau, “Weighted lasso for sparse recovery with statistical prior support information,” *IEEE Transactions on Signal Processing*, vol. 66, no. 6, pp. 1607–1618, 2018.
- H. Mansour and R. Saab, “Recovery analysis for weighted  $\ell_1$ -minimization using the null space property,” *Applied and Computational Harmonic Analysis*, vol. 43, no. 1, pp. 23–38, 2017.
- M. Elad, J.-L. Starck, P. Querre, and D. L. Donoho, “Simultaneous cartoon and texture image inpainting using morphological component analysis (mca),” *Applied and computational harmonic analysis*, vol. 19, no. 3, pp. 340–358, 2005.
- A. Björck and G. H. Golub, “Numerical methods for computing angles between linear subspaces,” *Mathematics of computation*, vol. 27, no. 123, pp. 579–594, 1973.
- D. L. Donoho, M. Elad, and V. N. Temlyakov, “Stable recovery of sparse overcomplete representations in the presence of noise,” *IEEE Transactions on information theory*, vol. 52, no. 1, pp. 6–18, 2005.
- J. A. Tropp, “Greed is good: Algorithmic results for sparse approximation,” *IEEE Transactions on Information Theory*, vol. 50, no. 10, pp. 2231–2242, 2004.
- Y. C. Eldar, P. Kuppinger, and H. Bolcskei, “Block-sparse signals: Uncertainty relations and efficient recovery,” *IEEE Transactions on Signal Processing*, vol. 58, no. 6, pp. 3042–3054, 2010.
- R. Kueng and D. Gross, “Ripless compressed sensing from anisotropic measurements,” *Linear Algebra and its Applications*, vol. 441, pp. 110–123, 2014.
- D. Martin, C. Fowlkes, D. Tal, and J. Malik, “A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics,” in *Proc. 8th Int’l Conf. Computer Vision*, vol. 2, July 2001, pp. 416–423.
- R. Mehrotra, K. Namuduri, and N. Ranganathan, “Gabor filter-based edge detection,” *Pattern Recognition*, vol. 25, no. 12, pp. 1479–94, 1992.