

# Nonparametric Bayesian Dictionary Learning for Analysis of Noisy and Incomplete Images

<sup>1</sup>Mingyuan Zhou, <sup>1</sup>Haojun Chen, <sup>1</sup>John Paisley, <sup>1</sup>Lu Ren, <sup>1</sup>Lingbo Li,

<sup>1</sup>Zhengming Xing, <sup>2</sup>David Dunson, <sup>3</sup>Guillermo Sapiro and <sup>1</sup>Lawrence Carin

<sup>1</sup>Department of Electrical and Computer Engineering and <sup>2</sup>Statistics Department  
Duke University, Durham, NC 27708-0291, USA

<sup>3</sup>Department of Electrical and Computer Engineering  
University of Minnesota, Minneapolis, MN 55455, USA

## Abstract

Nonparametric Bayesian methods are considered for recovery of imagery based upon compressive, incomplete and/or noisy measurements. A truncated beta-Bernoulli process is employed to infer an appropriate dictionary for the data under test, and also for image recovery. In the context of compressive sensing, significant improvements in image recovery are manifested using learned dictionaries, relative to using standard orthonormal image expansions. The compressive-measurement projections are also optimized for the learned dictionary. Additionally, we consider simpler (incomplete) measurements, defined by measuring a subset of image pixels, selected uniformly at random. Spatial inter-relationships within imagery are exploited through use of the Dirichlet and probit stick-breaking processes. Several example results are presented, with comparisons to other methods in the literature.

## I. INTRODUCTION

### A. Sparseness and dictionary learning

There has been significant recent interest in sparse image representations, in the context of denoising and interpolation [1], [13], [24]–[26], [28], [29], [32], compressive sensing (CS) [5], [12], and classification [41]. All of these applications exploit the fact that images may be sparsely represented in an appropriate dictionary. Most of the denoising, interpolation, and CS literature assumes “off-the-shelf” wavelet and DCT bases/dictionaries [21], but recent research has demonstrated the significant utility of learning an

Work partially supported by DOE, NSF, ONR, NGA, and ARO.

often over-complete dictionary matched to the signals of interest (*e.g.*, images) [1], [4], [12], [13], [24]–[26], [28], [29], [31], [32], [42].

Many of the existing methods for learning dictionaries are based on solving an optimization problem [1], [13], [24]–[26], [28], [29], in which one seeks to match the dictionary to the imagery of interest, while simultaneously encouraging a sparse representation. These methods have demonstrated state-of-the-art performance for denoising, super-resolution, interpolation, and inpainting. However, many existing algorithms for implementing such ideas also have some restrictions. For example, one must often assume access to the noise/residual variance, the size of the dictionary is set *a priori* or fixed via cross-validation type techniques, and a single (“point”) estimate is learned.

To mitigate the aforementioned limitations, dictionary learning has recently been cast as a factor-analysis problem, with the factor loadings corresponding to the dictionary elements (atoms). Utilizing nonparametric Bayesian methods like the beta process (BP) [30], [38], [43] and the Indian buffet process (IBP) [18], [22], one may for example infer the number of factors (dictionary elements) needed to fit the data. Further, one may place a prior on the noise or residual variance, with this inferred from the data [30], [43]. An approximation to the full posterior may be manifested via Gibbs sampling, yielding an ensemble of dictionary representations. Recent research has demonstrated that an ensemble of representations can be better than a single expansion [14], with such an ensemble naturally manifested by statistical models of the type described here.

### B. Exploiting structure and compressive measurements

In image analysis there is often additional information that may be exploited when learning dictionaries, with this well suited for Bayesian priors. For example, most natural images may be segmented, and it is probable that dictionary usage will be similar for regions within a particular segment class. To address this idea, we extend the model by employing a probit stick-breaking process (PSBP), with this a generalization of the Dirichlet process (DP) stick-breaking representation [35]. Related clustering techniques have proven successful in image processing [27]. The model clusters the image patches, with each cluster corresponding to a segment type; the PSBP encourages proximate and similar patches to be included within the same segment type, thereby performing image segmentation and dictionary learning simultaneously.

The principal focus of this paper is on applying hierarchical Bayesian algorithms to new compressive measurement techniques that have been developed recently. Specifically, we consider dictionary learning in the context of compressive sensing (CS) [5], [9], in which the measurements correspond to projections of typical image pixels. We consider dictionary learning performed “offline” based on representative

(training) images, with the learned dictionary applied within CS image recovery. We also consider the case for which the underlying dictionary is learned simultaneously with inversion (reconstruction), with this related to “blind” CS [17]. Finally, we design the CS projection matrix to be matched to the learned dictionary (when this is done offline), and demonstrate as in [12] that in practice this yields performance gains relative to conventional random CS projection matrices.

While CS is of interest for its potential to reduce the number of required measurements, it has the disadvantage of requiring the development of new classes of cameras. Such cameras are revolutionary and interesting [11], [36], but there have been decades of previous research performed on development of pixel-based cameras, and it would be desirable if such cameras could be modified simply to perform compressive measurements. We demonstrate that one may perform compressive measurements of natural images by simply sampling pixels measured by a conventional camera, with the pixels selected uniformly at random. This is closely related to recent research on matrix completion [6], [23], [34], but here we move beyond simple low-rank assumptions associated with such previous research.

### C. Contributions

This paper develops several hierarchical Bayesian models for learning dictionaries for analysis of imagery, with applications in denoising, interpolation and compressive sensing (CS). The inference is performed based on a Gibbs sampler, as is increasingly common in modern image processing [16]. Here we demonstrate how generalizations of the beta-Bernoulli process allow one to infer the dictionary elements directly based on the underlying degraded image, without any *a priori* training data, while simultaneously inferring the noise statistics and tolerating significant missingness in the imagery. This is achieved by exploiting the low-dimensional structure of most natural images, which implies that image patches may be represented in terms of a low-dimensional set of learned dictionary elements. Excellent results are realized in these applications, including for hyperspectral imagery, which has not been widely considered in these settings previously. We also demonstrate how the learned dictionary may be easily employed to define CS projection matrices that yield markedly improved image recovery, as compared to the much more typical random construction of such matrices.

The basic hierarchical Bayesian architecture developed here serves as a foundation that may be flexibly extended to incorporate additional prior information. As examples, we here show how dictionary learning may be readily coupled with clustering, through use of a Dirichlet process [15]. We also incorporate spatial information within the image via a probit stick-breaking process [33], with these extensions yielding significant advantages for the CS application. The basic modeling architecture may also exploit additional information manifested in the form of general covariates, as considered in a recent paper [44].

The remainder of the paper is organized as follows. In Section II we review the classes of problems being considered. The beta-Bernoulli process is discussed in Section III, with relationships made with previous work in this area, including those based on the Indian buffet process. The Dirichlet and probit stick-breaking processes are discussed in Section IV, and several example results are presented in Section V. Conclusions and a discussion of future work are provided in Section VI, and details of the inference equations are summarized in the Appendix.

## II. PROBLEMS UNDER STUDY

We consider data samples that may be expressed in the form

$$\mathbf{x}_i = \mathbf{D}\mathbf{w}_i + \boldsymbol{\epsilon}_i \quad (1)$$

where  $\mathbf{x}_i \in \mathbb{R}^P$ ,  $\boldsymbol{\epsilon}_i \in \mathbb{R}^P$ , and  $\mathbf{w}_i \in \mathbb{R}^K$ . The columns of the matrix  $\mathbf{D} \in \mathbb{R}^{P \times K}$  represent the  $K$  components of a dictionary with which  $\mathbf{x}_i$  is expanded. For our problem, the  $\mathbf{x}_i$  will correspond to  $B \times B$  (overlapped) pixel patches in an image [1], [13], [24], [25], [28], [43]. The set of vectors  $\{\mathbf{x}_i\}_{i=1,N}$  may be extracted from an image(s) of interest.

For the denoising problem, the vectors  $\boldsymbol{\epsilon}_i$  may represent sensor noise, in addition to (ideally small) residual from representation of the underlying signal as  $\mathbf{D}\mathbf{w}_i$ . To perform denoising, we place restrictions on the vectors  $\mathbf{w}_i$ , such that  $\mathbf{D}\mathbf{w}_i$  by itself does not exactly represent  $\mathbf{x}_i$ . A popular such restriction is that  $\mathbf{w}_i$  should be sparse, motivated by the idea that any particular  $\mathbf{x}_i$  may often be represented in terms of a small subset of representative dictionary elements, from the full dictionary defined by the columns of  $\mathbf{D}$ . There are several methods that have been developed recently to impose such a sparse representation, including  $\ell_1$ -based relaxation algorithms [24], [25], iterative algorithms [1], [13], and Bayesian methods [43]. One advantage of a Bayesian approach is that the noise/residual statistics may be nonstationary (with unknown noise statistics). Specifically, in addition to placing a sparseness-promoting prior on  $\mathbf{w}_i$ , we may also impose a prior on the components of  $\boldsymbol{\epsilon}_i$ . From the estimated posterior density function on model parameters, each component of  $\boldsymbol{\epsilon}_i$ , corresponding to the  $i$ th  $B \times B$  image patch, has its own variance. Given  $\{\mathbf{x}_i\}_{i=1,N}$ , our goal may be to simultaneously infer  $\mathbf{D}$  and  $\{\mathbf{w}_i\}_{i=1,N}$  (and implicitly  $\boldsymbol{\epsilon}_i$ ), and then the denoised version of  $\mathbf{x}_i$  is represented as  $\mathbf{D}\mathbf{w}_i$ .

In many applications the total number of pixels  $N \cdot P$  may be large. However, it is well known that compression algorithms may be used on  $\{\mathbf{x}_i\}_{i=1,N}$  after the measurements have been performed, to significantly reduce the quantity of data that need be stored or communicated. This compression indicates that while the data dimensionality  $N \cdot P$  may be large, the underlying information content may be relatively low. This has motivated the field of compressive sensing [5], [9], [11], [36], in which the

total number of measurements performed may be much less than  $N \cdot P$ . Toward this end, researchers have proposed *projection* measurements of the form

$$\mathbf{y}_i = \Sigma \mathbf{x}_i \quad (2)$$

where  $\Sigma \in \mathbb{R}^{n \times P}$  and  $\mathbf{y}_i \in \mathbb{R}^n$ , ideally with  $n \ll P$ . The projection matrix  $\Sigma$  has traditionally been constituted randomly [5], [9], with a binary or real alphabet (and  $\Sigma$  may also be a function of the specific patch, and generalized as  $\Sigma_i$ ). It is desirable that matrices  $\Sigma$  and  $\mathbf{D}$  be as incoherent as possible.

The recovery of  $\mathbf{x}_i$  from  $\mathbf{y}_i$  is an ill-posed problem unless restrictions are placed on  $\mathbf{x}_i$ . We may exploit the same class of restrictions used in the denoising problem; specifically, the observed data satisfy  $\mathbf{y}_i = \Phi \mathbf{w}_i + \nu_i$ , with  $\Phi = \Sigma \mathbf{D}$  and  $\nu_i = \Sigma \epsilon_i$ , and with sparse  $\mathbf{w}_i$ . Note that the sparseness constraint implies that  $\{\mathbf{w}_i\}_{i=1,N}$  (and hence  $\{\mathbf{x}_i\}_{i=1,N}$ ) occupy distinct subspaces of  $\mathbb{R}^P$ , selected from the overall linear subspace defined by the columns of  $\Phi$ .

In most applications of compressive sensing  $\mathbf{D}$  is assumed known, corresponding to an orthonormal basis (*e.g.*, wavelets or a DCT) [5], [9], [21]. However, such bases are not necessarily well matched to natural imagery, and it is desirable to consider design of dictionaries  $\mathbf{D}$  for this purpose [12]. One may even consider recovering  $\{\mathbf{x}_i\}_{i=1,N}$  from  $\{\mathbf{y}_i\}_{i=1,N}$  while simultaneously inferring  $\mathbf{D}$ . Thus, we again have a dictionary-learning problem, which may be coupled with optimization of the CS matrix  $\Sigma$ , such that it is matched to  $\mathbf{D}$  (defined by a low coherence between the rows of  $\Sigma$  and columns of  $\mathbf{D}$  [5], [9], [12], [21]).

### III. SPARSE FACTOR ANALYSIS WITH THE BETA-BERNOULLI PROCESS

When presenting example results, we will consider three problems. For *denoising*, it is assumed we measure  $\mathbf{x}_i = \mathbf{D} \mathbf{w}_i + \epsilon_i$ , where  $\epsilon_i$  represents measurement noise and model error; for the *compressive-sensing* application we observe  $\mathbf{y}_i = \Sigma(\mathbf{D} \mathbf{w}_i + \epsilon_i) = \Phi \mathbf{w}_i + \nu_i$ , with  $\Phi = \Sigma \mathbf{D}$  and  $\nu_i = \Sigma \epsilon_i$ ; and, finally, for the *interpolation* problem we observe  $Q_\phi(\mathbf{D} \mathbf{w}_i + \epsilon_i)$ , where  $Q_\phi(\mathbf{x}_i)$  is a vector of elements from  $\mathbf{x}_i$  contained within the set  $\phi$ . For all three problems our objective is to infer the underlying signal  $\mathbf{D} \mathbf{w}_i$ , with  $\mathbf{w}_i$  assumed sparse; we generally wish to simultaneously infer  $\mathbf{D}$  and  $\{\mathbf{w}_i\}_{i=1,N}$ . To address each of these problems, we consider a statistical model for  $\mathbf{x}_i = \mathbf{D} \mathbf{w}_i + \epsilon_i$ , placing Bayesian priors on  $\mathbf{D}$ ,  $\mathbf{w}_i$  and  $\epsilon_i$ ; the way the model is used is modified slightly for each specific application. For example, when considering interpolation, only the observed  $Q_\phi(\mathbf{x}_i)$  are used within the model likelihood function.

#### A. Beta-Bernoulli process for active-set selection

Let the binary vector  $\mathbf{z}_i \in \{0, 1\}^K$  denote which of the  $K$  columns of  $\mathbf{D}$  are used for representation of  $\mathbf{x}_i$  (active set); if a particular component of  $\mathbf{z}_i$  is equal to one, then the corresponding column of  $\mathbf{D}$

is used in the representation of  $\mathbf{x}_i$ . Hence, for the data  $\{\mathbf{x}_i\}_{i=1,N}$  there is an associated set of latent binary vectors  $\{\mathbf{z}_i\}_{i=1,N}$ , and the beta-Bernoulli process provides a convenient prior for these vectors [30], [38], [43]. Specifically, consider the model

$$\mathbf{z}_i \sim \prod_{k=1}^K \text{Bernoulli}(\pi_k), \quad \boldsymbol{\pi} \sim \prod_{k=1}^K \text{Beta}(a/K, b(K-1)/K) \quad (3)$$

where  $\pi_k$  is the  $k$ th component of  $\boldsymbol{\pi}$ , and  $a$  and  $b$  are model parameters; the impact of these parameters on the model are discussed below. Note that the use of product notation in (3) is meant to denote that each component of  $\mathbf{z}_i$  and  $\boldsymbol{\pi}$  are drawn independently from distributions of the same form.

Considering the limit  $K \rightarrow \infty$ , and after integrating out  $\boldsymbol{\pi}$ , the draws of  $\{\mathbf{z}_i\}_{i=1,N}$  may be constituted as follows. For each  $\mathbf{z}_i$ , draw  $c_i \sim \text{Poisson}(\frac{a}{b+i-1})$  and define  $C_i = \sum_{j=1}^i c_j$ , with  $C_0 = 0$ . Let  $z_{ik}$  represent the  $k$ th component of  $\mathbf{z}_i$ , and  $z_{ik} = 0$  for  $k > C_i$ . For  $k = 1, \dots, C_{i-1}$ ,  $z_{ik} \sim \text{Bernoulli}(\frac{n_{ik}}{b+i-1})$ , where  $n_{ik} = \sum_{j=1}^{i-1} z_{jk}$  ( $n_{ik}$  represents the total number of times the  $k$ th component of  $\{\mathbf{z}_j\}_{j=1,i-1}$  is one). For  $k = C_{i-1} + 1, \dots, C_i$ , we set  $z_{ik} = 1$ . Note that as  $a/(b+i-1)$  becomes small, with increasing  $i$ , it is probable that  $c_i$  will be small. Hence, with increasing  $i$ , the number of new non-zero components of  $\mathbf{z}_i$  diminishes. Further, as a consequence of  $\text{Bernoulli}(\frac{n_{ik}}{b+i-1})$ , when a particular component of the vectors  $\{\mathbf{z}_j\}_{j=1,i-1}$  is frequently one, it is more probable that it will be one for subsequent  $\mathbf{z}_j$ ,  $j \geq i$ . When  $b = 1$  this construction for  $\{\mathbf{z}_i\}_{i=1,N}$  corresponds to the Indian buffet process [18].

Since  $\mathbf{z}_i$  defines which columns of  $\mathbf{D}$  are used to represent  $\mathbf{x}_i$ , (3) imposes that it is probable that some columns of  $\mathbf{D}$  are used repeatedly among the set  $\{\mathbf{x}_i\}_{i=1,N}$ , while other columns of  $\mathbf{D}$  may be more specialized to particular  $\mathbf{x}_i$ . As demonstrated below, this has been found to be a good model when  $\{\mathbf{x}_i\}_{i=1,N}$  are patches of pixels extracted from natural images.

### B. Full hierarchical model

The hierarchical form of the model may now be expressed as

$$\begin{aligned} \mathbf{x}_i &= \mathbf{D}\mathbf{w}_i + \boldsymbol{\epsilon}_i \\ \mathbf{w}_i &= \mathbf{z}_i \odot \mathbf{s}_i \\ \mathbf{d}_k &\sim \mathcal{N}(0, P^{-1}\mathbf{I}_P) \\ \mathbf{s}_i &\sim \mathcal{N}(0, \gamma_s^{-1}\mathbf{I}_K) \\ \boldsymbol{\epsilon}_i &\sim \mathcal{N}(0, \gamma_\epsilon^{-1}\mathbf{I}_P) \end{aligned} \quad (4)$$

where  $\mathbf{d}_k$  represents the  $k$ th component (atom) of  $\mathbf{D}$ ,  $\odot$  represents the elementwise or Hadamard vector product,  $\mathbf{I}_P$  ( $\mathbf{I}_K$ ) represents a  $P \times P$  ( $K \times K$ ) identity matrix, and  $\{\mathbf{z}_i\}_{i=1,N}$  are drawn as in (3).

Conjugate hyperpriors  $\gamma_s \sim \text{Gamma}(c, d)$  and  $\gamma_\epsilon \sim \text{Gamma}(e, f)$  are also imposed. The construction in (4), and with the prior in (3) for  $\{\mathbf{z}_i\}_{i=1,N}$ , is henceforth referred to as the beta process factor analysis (BPFA) model. This model was first developed in [22], with a focus on general factor analysis; here we apply and extend this construction for image-processing applications.

Note that we impose independent Gaussian *priors* for  $\mathbf{d}_k$ ,  $\mathbf{s}_i$  and  $\epsilon_i$  for modeling convenience (conjugacy of consecutive terms in the hierarchical model). However, the inferred *posterior* for these terms is generally *not* independent or Gaussian. The independent priors essentially impose prior information about the *marginals* of the posterior of each component, while the inferred posterior accounts for statistical dependence as reflected in the data.

To make connections of this model to more-typical optimization-based approaches [24], [25], note that the negative logarithm of the posterior density function is

$$\begin{aligned} -\log p(\Theta|\mathcal{D}, \mathcal{H}) &= \frac{\gamma_\epsilon}{2} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{D}(\mathbf{s}_i \odot \mathbf{z}_i)\|_2^2 + \frac{P}{2} \sum_{k=1}^K \|\mathbf{d}_k\|_2^2 + \frac{\gamma_s}{2} \sum_{i=1}^N \|\mathbf{s}_i\|_2^2 \\ &\quad - \log f_{\text{Beta-Bern}}(\{\mathbf{z}_i\}_{i=1}^N; \mathcal{H}) - \log \text{Gamma}(\gamma_\epsilon|\mathcal{H}) - \log \text{Gamma}(\gamma_s|\mathcal{H}) + \text{Const.} \end{aligned} \quad (5)$$

where  $\Theta$  represents all unknown model parameters,  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1,N}$ ,  $f_{\text{Beta-Bern}}(\{\mathbf{z}_i\}_{i=1}^N; \mathcal{H})$  represents the beta-Bernoulli process prior in (3), and  $\mathcal{H}$  represents model hyper-parameters (*i.e.*,  $a, b, c, d, e$  and  $f$ ). Therefore, the typical  $\ell_2$  constraints [24], [25] on the dictionary elements  $\mathbf{d}_k$  and on the non-zero weights  $\mathbf{s}_i$  correspond here to the Gaussian priors employed in (4). However, rather than employing an  $\ell_1$  (Laplacian prior) constraint [24], [25] to impose sparseness on  $\mathbf{w}_i$ , we employ the beta-Bernoulli process and  $\mathbf{w}_i = \mathbf{s}_i \odot \mathbf{z}_i$ . The beta-Bernoulli process imposes that the binary  $\mathbf{z}_i$  should be sparse, *and* that there should be a relatively consistent (re)use of dictionary elements across the image, thereby also imposing self-similarity. Further, and perhaps most importantly, we do *not* constitute a point estimate, as one would do if a single  $\Theta$  was sought to minimize (5). We rather estimate the full posterior density  $p(\Theta|\mathcal{D}, \mathcal{H})$ , implemented via Gibbs sampling. A significant advantage of the hierarchical construction in (4) is that each Gibbs update equation is analytic, with detailed update equations provided in the Appendix. Note that consistent use of atoms is encouraged because the active sets are defined by the binary vectors  $\{\mathbf{z}_i\}_{i=1,N}$ , and these are all drawn from a shared probability vector  $\boldsymbol{\pi}$ ; this is distinct from drawing the active sets i.i.d. from a Laplacian prior. Further, the beta-Bernoulli prior imposes that many components of  $\mathbf{w}_i$  are exactly zero, while with a Laplacian prior many components are small but not exactly zero (hence the former is analogous to  $\ell_0$  regularization, with the latter closer to  $\ell_1$  regularization).

## IV. PATCH CLUSTERING VIA DIRICHLET AND PROBIT STICK-BREAKING PROCESSES

### A. Motivation

In the model discussed above each patch  $\mathbf{x}_i$  had a unique usage of dictionary atoms, defined by the binary vector  $\mathbf{z}_i$ , which selects columns of  $\mathbf{D}$ . One may wish to place further constraints on the model, thereby imposing a greater degree of statistical structure. For example, one may employ that the  $\mathbf{x}_i$  cluster, and that within each cluster each of the associated  $\mathbf{x}_i$  employ the same columns of  $\mathbf{D}$ . This is motivated by the idea that a natural image may be clustered into different types of textures or general image structure. However, rather than imposing that all  $\mathbf{x}_i$  within a given cluster use *exactly* the same columns of  $\mathbf{D}$ , one may want to impose that all  $\mathbf{x}_i$  within such a cluster share the same probability of dictionary usage, *i.e.*, that all  $\mathbf{x}_i$  within cluster  $c$  share the same probability of using columns of  $\mathbf{D}$ , defined by  $\boldsymbol{\pi}_c$ , rather than sharing a single  $\boldsymbol{\pi}$  for all  $\mathbf{x}_i$  (as in the original model above). Again, such clustering is motivated by the idea that natural images tend to segment into different textural or color forms. Below, we perform clustering in terms of the vectors  $\boldsymbol{\pi}_c$ , rather than explicit clustering of dictionary usage, which would entail cluster-dependent  $\mathbf{z}_c$ ; the “softer” nature of the former clustering structure is employed to retain model flexibility, while still encouraging sharing of parameters within clusters.

A question when performing such clustering concerns the *number* of clusters needed, this motivating the use of nonparametric methods, like those considered in the next subsections. Additionally, since the aforementioned clustering is motivated by the segmentations characteristic of natural images, it is desirable to explicitly utilize the spatial location of each image patch, encouraging that the patches  $\mathbf{x}_i$  in a particular segment/cluster are spatially contiguous. This latter goal motivates use of a probit stick-breaking process, as also detailed below.

### B. Dirichlet process

The Dirichlet process (DP) [15] constitutes a popular means of performing nonparametric clustering. A random draw from a DP,  $G \sim \text{DP}(\alpha G_0)$ , with precision  $\alpha \in \mathbb{R}^+$  and “base” measure  $G_0$ , may be constituted via the stick-breaking construction [35]

$$G = \sum_{l=1}^{\infty} \beta_l \delta_{\theta_l^*}, \quad \theta_l^* \sim G_0 \quad (6)$$

where  $\beta_l = V_l \prod_{h=1}^{l-1} (1 - V_h)$  and  $V_h \sim \text{Beta}(1, \alpha)$ . The  $\beta_l$  may be viewed as a sequence of fractional breaks from a “stick” of original length one, where the fraction of stick broken off on break  $l$  is  $V_l$ . The  $\theta_l^*$  are model parameters, associated with the  $l$ th data cluster. For our problem it has proven effective to

set  $G_0 = \prod_{k=1}^K \text{Beta}(a/K, b(K-1)/K)$  analogous to (3), and hence  $G = \sum_{l=1}^{\infty} \beta_l \delta_{\pi_l^*}$ . The  $\pi_l^*$ , drawn from  $G_0$ , correspond to distinct probability vectors for using the  $K$  dictionary elements (columns of  $\mathbf{D}$ ). For sample  $i$  we draw  $\pi_i \sim G$ , and a separate sparse binary vector  $\mathbf{z}_i$  is drawn for each sample  $x_i$ , as  $\mathbf{z}_i \sim \prod_{k=1}^K \text{Bernoulli}(\pi_{ik})$ , with  $\pi_{ik}$  the  $k$ th component of  $\pi_i$ . In practice we truncate the infinite sum for  $G$  to  $N_L$  elements, and impose  $V_{N_L} = 1$ , such that  $\sum_{l=1}^{N_L} \beta_l = 1$ . A (conjugate) gamma prior is placed on the DP parameter  $\alpha$ .

We may view this DP construction as an “Indian buffet franchise,” generalizing the Indian buffet analogy [18]. Specifically, there are  $N_L$  Indian buffet restaurants; each restaurant is composed of the same “menu” (columns of  $\mathbf{D}$ ), and is distinguished by different probabilities for selecting menu items. The “customers”  $\{x_i\}_{i=1,N}$  cluster based upon which restaurant they go to. The  $\{\pi_l^*\}_{l=1,N_L}$  represent the probability of using each column of  $\mathbf{D}$  in the respective  $N_L$  different buffets. The  $\{x_i\}_{i=1,N}$  cluster themselves among the different restaurants in a manner that is consistent with the characteristics of the data, with the model also simultaneously learning the dictionary/menu  $\mathbf{D}$ . Note that we typically make the truncation  $N_L$  large, and the posterior distribution infers the number of clusters actually needed to support the data, as represented by how many  $\beta_l$  are of significant value. The model in (4), with the above DP construction for  $\{\mathbf{z}_i\}_{i=1,N}$ , is henceforth referred to as DP-BPFA.

### C. Probit stick-breaking process

The DP yields a clustering of  $\{x_i\}_{i=1,N}$ , but it does not account for our knowledge of the location of each patch within the image. It is natural to expect that if  $x_i$  and  $x_{i'}$  are proximate then they are likely to be constituted in terms of similar columns of  $\mathbf{D}$ <sup>1</sup>. To impose this information, we employ the probit stick-breaking process (PSBP). A *logistic* stick-breaking process is discussed in detail in [33]. We employ the closely related probit version here because it may be easily implemented in a Gibbs sampler. We note that while the method in [33] is related to that discussed below, in [33] the concepts of learned dictionaries and beta-Bernoulli priors were not considered. Another related model, that employs a probit link function, is discussed in [7].

We augment the data as  $\{x_i, \mathbf{r}_i\}_{i=1,N}$ , where  $x_i$  again represents pixel values from the  $i$ th image patch, and  $\mathbf{r}_i \in \mathbb{R}^2$  represents the two-dimensional location of each patch. We wish to impose that proximate patches are more likely to be composed of the same or similar columns of  $\mathbf{D}$ . In the PSBP construction, all aspects of (4) are retained, except for the manner in which  $\mathbf{z}_i$  are constituted. Rather than drawing a

<sup>1</sup>Proximity can be modeled as in “spatial proximity,” as here developed in detail, or “feature proximity” as in non-local means and related approaches, see [27] and references therein.

single  $K$ -dimensional vector of probabilities  $\pi$  as in (3), we draw a *library* of such vectors:

$$\pi_l^* \sim \prod_{k=1}^K \text{Beta}(a/K, b(K-1)/K) , \quad l = 1, \dots, N_L \quad (7)$$

and each  $\pi_l^*$  is associated with a particular segment in the image. One  $\pi_i$  is associated with location  $r_i$ , and drawn

$$\pi_i \sim \sum_{l=1}^{N_L} \beta_l(r_i) \delta_{\pi_l^*} \quad (8)$$

with  $\sum_{l=1}^{N_L} \beta_l(r_i) = 1$  for all  $r_i$ , and  $\delta_{\pi_l^*}$  represents a point measure concentrated at  $\pi_l^*$ . Once  $\pi_i$  is associated with a particular  $x_i$ , the corresponding binary vector  $z_i$  is drawn as in the first line of (3). Note that the distinction between DP and PSBP is that in the former the mixture weights  $\{\beta_l\}_{l=1, N_L}$  are independent of spatial position  $r$ , while the latter explicitly utilizes  $r$  within  $\{\beta_l(r)\}_{l=1, N_L}$  (and below we impose that  $\beta_l(r)$  changes smoothly with  $r$ ).

The space-dependent weights are constructed as  $\beta_l(r) = V_l(r) \prod_{h=l}^{l-1} [1 - V_h(r)]$  where  $0 < V_l(r) < 1$  constitute space-dependent probabilities. We set  $V_{N_L} = 1$ , and for  $l \leq N_L - 1$  the  $V_l$  are space-dependent probit functions:

$$V_l(r) = \int_{-\infty}^{g_l(r)} dx \mathcal{N}(x|0, 1), \quad g_l(r) = \zeta_{l0} + \sum_{i=1}^N \zeta_{li} \mathcal{K}(r, r_i; \psi_l) \quad (9)$$

where  $\mathcal{K}(r, r_i; \psi_l)$  is a kernel characterized by parameter  $\psi_l$  and  $\{\zeta_{li}\}_{i=0, N}$  are a *sparse* set of real numbers. To implement the sparseness on  $\{\zeta_{li}\}_{i=0, N}$ , within the prior  $\zeta_{li} \sim \mathcal{N}(0, \alpha_{li}^{-1})$ , and (conjugate)  $\alpha_{li} \sim \text{Gamma}(a_0, b_0)$ , with  $(a_0, b_0)$  set to favor most  $\alpha_{li}$  being large (if  $\alpha_{li}$  is large, a draw  $\mathcal{N}(0, \alpha_{li}^{-1})$  is likely to be near zero, such that most  $\{\zeta_{li}\}_{i=0, N}$  are near zero). This sparseness-promoting construction is the same as that employed in the relevance vector machine (RVM) [40]. We here utilize a radial basis function (RBF) kernel  $\mathcal{K}(r, r_i; \psi_l) = \exp[-\|r_i - r\|_2/\psi_l]$ .

Each  $g_l(r)$  is encouraged to only be defined by a small set of localized kernel functions, and via the probit link function  $\int_{-\infty}^{g_l(r)} dx \mathcal{N}(x|0, 1)$  the probability  $V_l(r)$  is characterized by localized segments over which the probability  $V_l(r)$  is contiguous and smoothly varying. The  $V_l(r)$  constitute a space-dependent stick-breaking process. Since  $V_{N_L} = 1$ ,  $\sum_{l=1}^{N_L} \beta_l(r) = 1$  for all  $r$ .

The PSBP model is relatively simple to implement within a Gibbs sampler. For example, as indicated above, sparseness on  $\zeta_{li}$  is imposed as in the RVM, and the probit link function is simply implemented within a Gibbs sampler (which is why it was selected, rather than a logistic link function). Finally, we define a finite set of possible kernel parameters  $\{\psi_j\}_{j=1, N_p}$ , and a multinomial prior is placed on these parameters, with the multinomial probability vector drawn from a Dirichlet distribution [33] (each of the

$g_l(\mathbf{r})$  draws a kernel parameter from  $\{\psi_j\}_{j=1,N_p}$ . The model in (4), with the PSBP construction for  $\{\mathbf{z}_i\}_{i=1,N}$ , is henceforth referred to as PSBP-BPFA.

#### D. Discussion of proposed sparseness-imposing priors

The basic BPFA model is summarized in (4), and three related priors have been developed for the sparse binary vectors  $\{\mathbf{z}_i\}_{i=1,N}$ : (i) the basic truncated beta-Bernoulli process in (3), (ii) a DP-based clustering of the underlying  $\{\pi_i\}_{i=1,N}$ , and (iii) a PSBP clustering of  $\{\pi_i\}_{i=1,N}$  that exploits knowledge of the location of the image patches. For (ii) and (iii), the  $\mathbf{x}_i$  within a particular cluster have *similar*  $\mathbf{z}_i$ , rather than exactly the same binary vector; we also considered the latter, but this worked less well in practice. As discussed further when presenting results, for denoising and interpolation, all three methods yield comparable performance. However, for CS, (ii) and (iii) yield marked improvements in image-recovery accuracy relative to (i). In anticipation of these results, we provide a further discussion of the three priors on  $\{\mathbf{z}_i\}_{i=1,N}$  and on the three image-processing problems under consideration.

For the denoising and interpolation problems, we are provided with the data  $\{\mathbf{x}_i\}_{i=1,N}$ , albeit in the presence of noise and potentially with substantial missing pixels. However, for this problem  $N$  may be made quite large, since we may consider all possible (overlapping)  $B \times B$  patches. A given pixel (apart from near the edges of the image) is present in  $B^2$  different patches. Perhaps because we have such a large quantity of partially overlapping data, for denoising and interpolation we have found that beta-Bernoulli process in (3) is sufficient for inferring the underlying relationships between the different data  $\{\mathbf{x}_i\}_{i=1,N}$ , and processing these data collaboratively. However, the beta-Bernoulli construction does not explicitly segment the image, and therefore an advantage of the PSBP-BPFA construction is that it yields comparable denoising and interpolation performance as (3), while also simultaneously yielding an effective image segmentation.

For the CS problem, we measure  $\mathbf{y}_i = \Sigma \mathbf{x}_i$ , and therefore each of the  $n$  measurements associated with each image patch ( $\Sigma \in \mathbb{R}^{n \times P}$ ) loses the original pixels in  $\mathbf{x}_i$  (the projection matrix  $\Sigma$  may also change with each patch, denoted  $\Sigma_i$ ). Therefore, for CS one cannot consider all possible shifts of the patches, as the patches are predefined and fixed in the CS measurement (in the denoising and interpolation problems the patches are defined in the subsequent analysis). Therefore, for CS imposition of the clustering behavior via DP or PSBP provides important information, yielding state-of-the-art CS-recovery results.

#### E. Possible extensions

The “Indian buffet franchise” and probit stick-breaking process constructions considered above and in the results below draw the  $K$ -dimensional probability vectors  $\pi_l^*$  independently. This implies that

within the prior we impose no statistical correlation between the components of vectors  $\pi_l^*$  and  $\pi_{l'}^*$ , for  $l \neq l'$ . It may be desirable to impose such structure, imposing that there is a “global” probability of using particular dictionary elements, and the different mixture components within the DP/PSBP constructions correspond to specific draws from global statistics of dictionary usage. This will encourage the idea that there may be some “popular” dictionary elements that are shared across different mixture components (*i.e.*, popular across different buffets in the franchise). One can impose this additional structure via a *hierarchical* BP construction (HBP) [38], related to the *hierarchical* DP (HDP) [37]. Briefly, in an HBP construction one may draw the  $\pi_l^*$  via the hierarchical construction

$$\pi_l^* \sim \prod_{k=1}^K \text{Beta}(c\eta_k, c(1 - \eta_k)) , \quad \boldsymbol{\eta} \sim \prod_{k=1}^K \text{Beta}(a/K, b(K-1)/K) \quad (10)$$

where  $\eta_k$  is the  $k$ th component of  $\boldsymbol{\eta}$ . The vector  $\boldsymbol{\eta}$  constitutes “global” probabilities of using each of the  $K$  dictionary elements (across the “franchise”), and  $\pi_l^*$  defines the probability of dictionary usage for the  $l$ th buffet. This construction imposes statistical dependencies among the vectors  $\{\pi_l^*\}$ .

To simplify the presentation, in the below example results we do *not* consider the HBP construction, as good results have already been achieved with the (truncated) DP and PSBP models discussed above. We note that in these analyses the truncated versions of these models may actually help inference of statistical correlations among  $\{\pi_l^*\}$  within the *posterior* (since the set  $\{\pi_l^*\}$  is finite, and each vector  $\pi_l^*$  is of finite length). If we actually considered the infinite limit on  $K$  and on the number of mixture components, inference of such statistical relationships within the posterior may be undermined, because specialized dictionary elements may be constituted across the different franchises, rather than encouraging *sharing* of highly similar dictionary elements.

While we do not focus on the HBP construction here, a recent paper has employed the HBP construction in related dictionary learning for image-processing applications, yielding very encouraging results [44]. We therefore emphasize that the basic hierarchical Bayesian construction employed here is very flexible, and may be extended in many ways to impose additional structure.

## V. EXAMPLE RESULTS

### A. Reproducible research

The test results and the Matlab code to reproduce them can be downloaded from <http://www.ee.duke.edu/~mz1/Results/BPFAImage/>.

### B. Parameter settings

For all BPFA, DP-BPFA and PSBP-BPFA computations, the dictionary truncation level was set at  $K = 256$  or  $K = 512$  based on the size of the image. Not all  $K$  dictionary elements are used in the model; the truncated beta-Bernoulli process infers the subset of dictionary elements employed to represent the data  $\{\mathbf{x}_i\}_{i=1,N}$ . The larger the image, the more distinct types of structure are anticipated, and therefore the more dictionary elements are likely to be employed; however, very similar results are obtained with  $K = 512$  in all examples, just with more dictionary elements not employed for smaller images (therefore, to save computational resources, we set  $K = 256$  for the smaller images). The number of DP and PSBP sticks was set at  $N_L = 20$ . The library of PSBP parameters is defined as in [33]; the PSBP kernel locations,  $\{\mathbf{r}_i\}_{i=1,N}$ , were situated on a uniformly sampled grid in each image dimension, situated at every fourth pixel in each direction (the results are insensitive to many related definitions of  $\{\mathbf{r}_i\}_{i=1,N}$ ). The hyperparameters within the gamma distributions were set as  $c = d = e = f = 10^{-6}$ , as is typically done in models of this type [40] (the same settings were used for the gamma prior for the DP precision parameter  $\alpha$ ). The beta-distribution parameters are set as  $a = K$  and  $b = 1$  if random initialization is used or  $a = K$  and  $b = N/8$  if a singular value decomposition (SVD) based initialization is used. None of these parameters have been optimized or tuned. When performing inference, all parameters are initialized randomly (as a draw from the associated prior) or based on the SVD of the image under test. The Gibbs samplers for the BPFA, DP-BPFA and PSBP-BPFA have been found to mix and converge quickly, producing satisfactory results with as few as 20 iterations. The inferred images represent the average from the collection samples. All software was written in non-optimized Matlab. On a Dell Precision T3500 computer with a 2.4 GHz CPU, for  $N = 148,836$  patches of size  $8 \times 8 \times 3$  with 20% of the RGB pixels observed at random, the BPFA required about 2 minutes per Gibbs iteration (the DP version was comparable), and PSBP-BPFA required about 3 minutes per iteration. For the 106-band hyperspectral imagery, which employed  $N = 428,578$  patches of size  $4 \times 4 \times 106$  with 2% of the voxels observed uniformly at random, each Gibbs iteration required about 15 minutes.

### C. Denoising

The BPFA denoising algorithm is compared with the original KSVD [13], for both grey-scale and color images. Newer denoising algorithms include block matching with 3D filtering (BM3D) [8], the multiscale KSVD [29], and KSVD with the non-local mean constraints [26]. These algorithms assume the noise variance is known, while the proposed model automatically infers, as part of the same model, the noise variance from the image under test. There are existing methods for estimation of the noise variance, as

a preprocessing step, *e.g.*, via wavelet shrinkage [10]. However, it was shown in [43] that the denoising accuracy of methods like that in [13] can be sensitive to small errors in the estimated variance, which are likely to occur in practice. Additionally, when doing *joint* denoising and image interpolation, as we consider below, methods like that in [10] may not be applied directly to estimate the noise variance, as there is a large fraction of missing data. Moreover, the BPFA, DP-BPFA and PSBP-BPFA models infer a potentially non-stationary noise variance, with a broad prior on the variance imposed by the gamma distribution.

In the denoising examples we consider the BPFA model in (3); similar results are obtained via the DP-BPFA and PSBP-BPFA models discussed in Section IV. To infer the final signal value at each pixel, we average the associated pixel contribution from each of the overlapping patches in which it is included (this is true for all results presented below in which overlapping patches were used).

In Table I we consider images from [13]. The proposed BPFA performs very similarly to KSVD. As one representative example of the model's ability to infer the noise variance, we consider the Lena image from Table I. The mean inferred noise standard deviations are 5.83, 10.59, 15.53, 20.48, 25.44, 50.46 and 100.54 for images contaminated by noise with respective standard deviations of 5, 10, 15, 20, 25, 50 and 100. Each of these noise variances were automatically inferred using exactly the same model, with no changes to the gamma hyperparameters (while for the KSVD results it was assumed that the noise variance was known exactly *a priori*).

In Table II we present similar results, for denoising RGB images; the KSVD comparisons come from [28]. An example denoising result is shown in Figure 1. As another example of the BPFA's ability to infer the underlying noise variance, for the castle image, the mean (automatically) inferred variances are 5.15, 10.18, 15.22 and 25.23 for images with additive noise with true respective standard deviations 5, 10, 15 and 25. The sensitivity of the KSVD algorithm to a mismatch between the assumed and true noise variances is shown in Figure 1 in [43], and the insensitivity of BPFA to changes in the noise variance and to requiring knowledge of the noise variance is deemed an important advantage.

It is also important to note that the grey-scale KSVD results in Table I were initialized using an over-complete DCT dictionary, while the RGB KSVD results in Table II employed an extensive set of training imagery to learn a dictionary  $\mathbf{D}$  that was used to initialize the denoising computations. All BPFA, DP-BPFA and PSBP-BPFA results employ no training data, with the dictionary initialized at random using draws from the prior or with the SVD of the data under test.



Fig. 1. From left to right: the original horses image, the noisy horses image with the noise standard deviation of 25, the denoised image and the inferred dictionary with its elements ordered in the probability to be used (from top-left). The low-probability dictionary elements are never used to represent  $\{x_i\}_{i=1,N}$ , and are draws from the prior, showing the ability of the model to learn the number of dictionary elements needed for the data.

TABLE I

GREY-SCALE IMAGE DENOISING PSNR RESULTS, COMPARING KSVD [13] AND BPFA, USING PATCH SIZE  $8 \times 8$ . THE TOP AND BOTTOM PARTS OF EACH CELL ARE RESULTS OF KSVD AND BPFA, RESPECTIVELY.

$\sigma$	C.man	House	Peppers	Lena	Barbara	Boats	F.print	Couple	Hill
5	37.87	39.37	37.78	38.60	38.08	37.22	36.65	37.31	37.02
	37.32	39.18	37.24	38.20	37.94	36.43	36.29	36.77	36.24
10	33.73	35.98	34.28	35.47	34.42	33.64	32.39	33.52	33.37
	33.40	36.29	34.31	35.62	34.63	33.70	32.42	33.63	33.31
15	31.42	34.32	32.22	33.70	32.37	31.73	30.06	31.45	31.47
	31.34	34.52	32.46	33.93	32.61	31.97	30.23	31.73	31.64
20	29.91	33.20	30.82	32.38	30.83	30.36	28.47	30.00	30.18
	30.03	33.25	31.10	32.65	31.10	30.70	28.72	30.34	30.47
25	28.85	32.15	29.73	31.32	29.60	29.28	27.26	28.90	29.18
	28.99	32.24	30.00	31.63	29.88	29.70	27.58	29.28	29.57
50	25.73	27.95	26.13	27.79	25.47	25.95	23.24	25.32	26.27
	25.67	28.49	26.46	28.29	26.03	26.50	24.14	25.94	26.81
100	21.69	23.71	21.75	24.46	21.89	22.81	18.30	22.60	23.98
	21.93	24.37	22.73	24.95	22.13	23.32	20.44	23.01	24.22

#### D. Image interpolation

For the initial interpolation examples, we consider standard RGB images, with 80% of the RGB pixels missing uniformly at random (the data under test are shown in Figure 2). Results are first presented for

TABLE II

RGB IMAGE DENOISING PSNR RESULTS COMPARING KSVD [28] AND BPFA, BOTH USING A PATCH SIZE OF  $7 \times 7$ . THE TOP AND BOTTOM PARTS OF EACH CELL SHOW THE RESULTS OF KSVD AND BPFA, RESPECTIVELY.

$\sigma$	Castle	Mushroom	Train	Horses	Kangaroo
5	40.37	39.93	39.76	40.09	39.00
	40.34	39.73	39.38	39.96	39.00
10	36.24	35.60	34.72	35.43	34.06
	36.28	35.70	34.48	35.48	34.21
15	33.98	33.18	31.70	32.76	31.30
	34.04	33.41	31.63	32.98	31.68
25	31.19	30.26	28.16	29.81	28.39
	31.24	30.62	28.28	30.11	28.86

the Castle and Mushroom images, with comparisons between the BPFA model in (3) and the PSBP-BPFA model discussed in Section IV. The difference between the two is that the former is a “bag-of-patches” model, while the latter accounts for the spatial locations of the patches. Further, the PSBP-BPFA simultaneously performs image recovery and segmentation. The results are shown in Figure 3, presenting the mean reconstructed images and inferred segmentations. Each color in the inferred segmentation represents one PSBP mixture component, and the figure shows the last Gibbs iteration (to avoid issues with label switching between Gibbs iterations). While the BPFA does not directly yield a segmentation, its PSNR results are comparable to those inferred by PSBP-BPFA, as summarized in Table III.

TABLE III

COMPARISON OF INTERPOLATION OF THE CASTLE AND MUSHROOM IMAGES, BASED UPON OBSERVING 20% OF THE PIXELS, SELECTED UNIFORMLY AT RANDOM. RESULTS ARE SHOWN USING BPFA AND PSBP-BPFA, AND THE ANALYSIS IS SEPARATELY PERFORMED USING  $8 \times 8 \times 3$  AND  $5 \times 5 \times 3$  IMAGE PATCHES.

	Castle $8 \times 8 \times 3$	Castle $5 \times 5 \times 3$	Mushroom $8 \times 8 \times 3$	Mushroom $5 \times 5 \times 3$
BPFA	29.32	28.48	31.63	31.17
PSBP-BPFA	29.54	28.46	32.03	31.27

An important additional advantage of Bayesian models like BPFA, DP-BPFA and PSBP-BPFA is that they provide a measure of confidence in the accuracy of the inferred image. In Figure 4 we plot the



Fig. 2. Images with 80% of the RGB pixels missing at random. Although only 20% of the actual pixels are observed, in these figures the missing pixels are estimated based upon averaging all observed neighboring pixels within a  $5 \times 5$  spatial extent. Left: castle image (PSNR 22.58 dB), right: mushroom image (24.85 dB).

variance of the inferred error  $\{\epsilon_i\}_{i=1,N}$ , computed via the Gibbs collection samples.

To provide a more-thorough examination of model performance, in Table IV we present results for several well-studied grey-scale and RGB images, as a function of the fraction of pixels missing. All of these results are based upon BPFA, with DP-BPFA and PSBP-BPFA yielding similar results. Finally, in Table V we perform interpolation and denoising simultaneously, again with no training data and without prior knowledge of the noise level (again, as discussed above, it would be difficult to estimate the noise variance as a preprocessing step in this case using methods like that in [10], as in this case there are also a large number of missing pixels). An example result is shown in Figure 5. To our knowledge, this is the first time denoising and interpolation have been performed jointly, while simultaneously inferring the noise statistics.

For all of the examples considered above, for both grey-scale and RGB images, we also attempted a direct application of matrix completion based on the incomplete matrix  $\mathbf{X} \in \mathbb{R}^{P \times N}$ , with columns defined by the image patches (*i.e.*, for  $N$  patches, with  $P$  pixels in each, the incomplete matrix is of size  $P \times N$ , with the  $i$ th column defined by the pixels in  $\mathbf{x}_i$ ). We considered the algorithm in [20], using software from Prof. Candès' website. For most of the examples considered above, even after very careful



Fig. 3. PSBP-BPFA analysis with 80% of the RGB pixels missing uniformly at random (see Figure 2). The analysis is based on  $8 \times 8 \times 3$  image patches, considering all possible (overlapping) patches. For a given pixel, the results are the average based upon all patches in which it is contained. For each example, recovered image based on an average of Gibbs collection samples (left), and each color representing one of the PSBP mixture components (right).

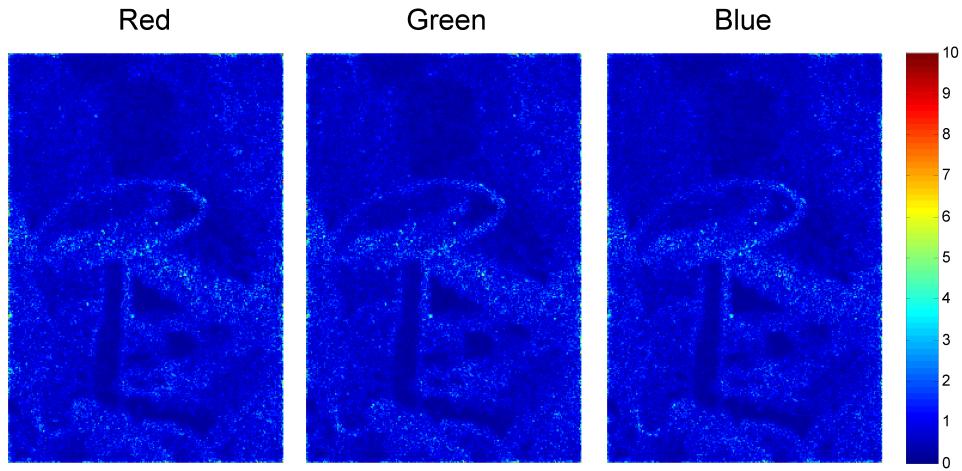


Fig. 4. Expected variance of each pixel for the (Mushroom) data considered in Figure 3.

tuning of the parameters, the algorithm diverged, suggesting that the low-rank assumptions were violated. For examples for which the algorithm did work, the PSNR values were typically 4 to 5 dB worse than those reported here for our model.

#### E. Interpolation of hyperspectral imagery

The basic BPFA, DP-BPFA and PSBP-BPFA technology may also be applied to hyperspectral imagery, and it is here where these methods may have significant practical utility. Specifically, the amount of data that need be measured and read off a hyperspectral camera is often enormous. By selecting a small

TABLE IV

TOP: BPFA GRAY-SCALE IMAGE INTERPOLATION PSNR RESULTS, USING PATCH SIZE  $8 \times 8$ . BOTTOM: BPFA RGB IMAGE INTERPOLATION PSNR RESULTS, USING PATCH SIZE  $7 \times 7$ .

data ratio	C.man	House	Peppers	Lena	Barbara	Boats	F.print	Man	Couple	Hill
20%	24.11	30.12	25.92	31.00	24.80	27.81	26.03	28.24	27.72	29.33
30%	25.71	33.14	28.19	33.31	27.52	30.00	09.01	30.06	30.00	31.21
50%	28.90	38.02	32.58	36.94	33.17	33.78	33.53	33.29	35.56	34.23
80%	34.70	43.03	37.73	41.27	40.76	39.50	40.17	39.11	38.71	38.75

data ratio	Castle	Mushroom	Train	Horses	Kangroo
20%	29.12	31.56	24.59	29.99	29.59
30%	32.02	34.63	27.00	32.52	32.21
50%	36.45	38.88	32.00	37.27	37.34
80%	41.51	42.56	40.73	41.97	42.74

TABLE V

SIMULTANEOUS IMAGE DENOISING AND INTERPOLATION PSNR RESULTS FOR BPFA, CONSIDERING THE BARBARA256 IMAGE AND USING PATCH SIZE  $8 \times 8$ .

$\sigma$	10%	20%	30%	50%	100%
0	23.47	26.87	29.83	35.60	42.94
5	23.34	26.73	29.27	33.61	37.70
10	23.16	26.07	28.17	31.17	34.31
15	22.66	25.17	26.82	29.31	32.14
20	22.17	24.27	25.62	27.90	30.55
25	21.68	23.49	24.72	26.79	29.30

fraction of voxels for measurement and read-out, selected uniformly at random, the quantity of data that need be handled is reduced substantially. Further, one may simply modify existing hyperspectral cameras. We consider hyperspectral data with 106 spectral bands, from HyMAP scene A.P. Hill, VA with permission from U.S. Army TEC. Because of the significant statistical correlation across the multiple spectral bands, the fraction of data that need be read is further reduced, relative to grey-scale or RGB imagery. In this example we considered 2% of the voxels, selected uniformly at random, and used image patches of size  $4 \times 4 \times 106$ . Other than the increased data dimensionality, nothing in the model was



Fig. 5. From left to right: the original barbara256 image, the noisy and incomplete barbara256 image with the noise standard deviation of 15 and 70% of its pixels missing at random (displayed based upon imputing missing pixels as the average of all observed pixels in a  $5 \times 5$  neighborhood), the restored image and the inferred dictionary with its elements ordered in the probability to be used (from top-left).

changed.

In Figure 6 we show example (mean) inferred images, at two (arbitrarily selected) spectral bands, as computed via BPFA. All 106 spectral bands are analyzed simultaneously. The average PSNR for the data cube (size  $845 \times 512 \times 106$ ) is 30.96 dB. While the PSNR value is of interest, for data of this type the more important question concerns the ability to classify different materials based upon the hyperspectral data. In a separate forthcoming paper we consider classification based on the full datacube, and based upon the BPFA-inferred datacube using 2% of the voxels, with encouraging results reported. We also tried the low-rank matrix completion algorithm from [20] for the hyperspectral data, and even after extensive parameter tuning, the algorithm diverged for all hyperspectral data considered.

In Table VI we summarize algorithm performance on another hyperspectral data set, composed of 210 spectral bands. We show the PSNR values as a function of percentage of observed data, and as a function of the size of the image patch. Note that the  $1 \times 1$  patches only exploit spectral information, while the other patch sizes exploit both spatial and spectral information.

#### F. Compressive sensing

We consider a CS example in which the image is divided into  $8 \times 8$  patches, with these constituting the underlying data  $\{\mathbf{x}_i\}_{i=1,N}$  to be inferred. For each of the  $N$  blocks, a vector of CS measurements  $\mathbf{y}_i = \Sigma \mathbf{x}_i$  is measured, where the number of projections per patch is  $n$ , and the total number of CS projections is  $n \cdot N$ . In our first examples the elements of  $\Sigma$  are constructed randomly, as draws from  $\mathcal{N}(0, 1)$ ; many other random projection classes may be considered [3] (and below we also consider optimized projections  $\Sigma$ , matched to the dictionary  $\mathbf{D}$ ). Each  $\mathbf{x}_i$  is assumed represented in terms of a dictionary  $\mathbf{x}_i = \mathbf{D}\mathbf{w}_i + \epsilon_i$ , and three constructions for  $\mathbf{D}$  were considered: (i) a DCT expansion; (ii) learning of  $\mathbf{D}$  using BPFA,

TABLE VI

BPFA HYPERSPECTRAL IMAGE INTERPOLATION PSNR RESULTS. FOR THIS EXAMPLE THE TEST IMAGE IS A  $150 \times 150$  URBAN IMAGE WITH 210 SPECTRAL BANDS. RESULTS ARE SHOWN AS A FUNCTION OF THE PERCENTAGE OF OBSERVED VOXELS, FOR DIFFERENT SIZED PATCHES (*e.g.*, THE  $4 \times 4$  CASE CORRESPONDS TO  $4 \times 4 \times 210$  “PATCHES”).

Observed data (%)	$1 \times 1$	$2 \times 2$	$3 \times 3$	$4 \times 4$
2	15.34	21.09	22.72	23.46
5	17.98	23.58	25.30	25.88
10	20.41	25.27	26.36	26.68
20	22.22	26.50	27.02	27.16

using training images; (*iii*) using the BPFA to perform *joint* CS inversion and learning of  $\mathbf{D}$ . For (*ii*), the training data consisted of 4000  $8 \times 8$  patches chosen at random from 100 images selected from the Microsoft database (<http://research.microsoft.com/en-us/projects/objectclassrecognition>). The dictionary was set to  $K = 256$ , and the offline beta process inferred a dictionary of size  $M = 237$ .

Representative CS reconstruction results are shown in Figure 7 (left) based upon a DCT dictionary, for a grey-scale version of the “castle” image. The results in Figure 7 (right) are based on a learned dictionary; except for the “online BP” results (where  $\mathbf{D}$  and  $\{\mathbf{w}_i\}_{i=1,N}$  are learned jointly), all of these results employ the same dictionary  $\mathbf{D}$  learned off-line as mentioned above, and the algorithms are distinguished by different ways of estimating  $\{\mathbf{w}_i\}_{i=1,N}$ . A range of CS-inversion algorithms are considered from the literature, and several BPFA-based constructions are considered as well for CS inversion. The online BPFA results (with no training data) are quite competitive with those based on a dictionary learned off-line.

Note that results based on a learned dictionary are markedly better than those based on the DCT; similar results were achieved when the DCT was replaced by a wavelet representation. For the DCT-based results, note that the DP-BPFA and PSBP-BPFA CS inversion results are significantly better than those of all other CS inversion algorithms. The results reported here are consistent with tests we performed using over 100 images from the aforementioned Microsoft database, not reported here in detail for brevity.

In all previous results the projection matrix  $\Sigma$  was constituted randomly. We now consider a simple means of matching  $\Sigma$  to a  $\mathbf{D}$  learned offline, based upon representative training images. Assume a learned  $\mathbf{D} \in \mathbb{R}^{P \times K}$ , with  $K > P$ , which may be represented via SVD as  $\mathbf{D} = \mathbf{U}\Lambda\mathbf{V}^T$ ;  $\mathbf{U} \in \mathbb{R}^{P \times P}$  and  $\mathbf{V} \in \mathbb{R}^{K \times P}$  are each composed of orthonormal columns, and  $\Lambda$  is a  $P \times P$  diagonal matrix. The columns of  $\mathbf{U}$  span the linear subspace of  $\mathbb{R}^P$  in which the columns of  $\mathbf{D}$  reside. Further, since the

columns of  $\mathbf{D}$  are generally *not* orthonormal, each column of  $\mathbf{D}$  is “spread out” when expanded in the columns of  $\mathbf{U}$ . Therefore, one expects that  $\mathbf{U}$  and  $\mathbf{D}$  are incoherent. Hence, a simple means of matching CS projections to the data is to define the rows of  $\Sigma$  in terms of randomly selected columns of  $\mathbf{U}$ . This was done in Figure 8 for the grey-scale “castle” image, using the same learned dictionary as considered in Figure 7. It is observed that this procedure yields a marked improvement in CS recovery accuracy, for all CS inversion algorithms considered.

Concerning computational costs, all CS inversions were run efficiently on PCs, with the specifics computational times dictated by the detailed Matlab implementation and the machine run on. A rough ranking of the computational speeds, from fastest to slowest, is as follows: StOMP-CFAR, Fast BCS, OMP, BPFA, LARS/Lasso, Online BPFA, DP-BPFA, PSBP-BPFA, VB BCS, Basis Pursuit; in this list, algorithms BPFA through Basis Pursuits have approximately the same computational costs.

The improved performance of the CS inversion based upon the learned dictionaries is manifested as a consequence of the structure that is imposed on the underlying image while performing inversion. The early CS inversion algorithms, of the type considered in the above comparisons, imposed that the underlying image is sparsely represented in an appropriate basis (we showed results here based upon a DCT expansion of the  $8 \times 8$  blocks over which CS inversion was performed, and similar results were manifested using a wavelet expansion). The imposition of such sparseness does not take into account the additional structure between the wavelet coefficients associated with natural images. Recently researchers have utilized such structure to move beyond sparseness, and achieve even better CS-inversion quality [2], [19]; in this work structural relationships and correlations *between* basis-function coefficients are accounted for. Additionally, there has been recent statistical research that has moved beyond sparsity, and that are of interest in the context of CS inversion [39]. In tests we omit for brevity, the algorithms in [2], [19] yield performance that is comparable to the best results to the right in Figure 7. Consequently, the imposition of structure in the form of learned dictionaries can be achieved using conventional basis expansions but with more-sophisticated inversion techniques, that account for structure in imagery. Therefore, the main advantage of the learned dictionary in the context of CS is that it provides a very convenient means of defining projection matrices that are matched to natural imagery, as done in Figure 8. Once those projection matrices are specified, the CS inversion may be employed based upon dictionaries as discussed here, or based upon more-traditional expansions (DCT or wavelets) and newer CS inversion methods [2], [19].

## VI. CONCLUSIONS

The truncated beta-Bernoulli process has been employed to learn dictionaries matched to image patches  $\{\mathbf{x}_i\}_{i=1,N}$ . The basic nonparametric Bayesian model is termed a beta process factor analysis (BPFA) framework, and extensions have also been considered. Specifically, the Dirichlet process (DP) has been employed to cluster the  $\{\mathbf{x}_i\}_{i=1,N}$ , encouraging similar dictionary-element usage within respective clusters. Further, the probit stick-breaking process (PSBP) has been used to impose that proximate patches are more likely to be clustered similarly (imposing that they are more probable to employ similar dictionary elements). All inference has been performed by a Gibbs sampler, with analytic update equations. The PBFA, DP-BPFA and PSBP-BPFA have been applied to three problems in image processing: (i) denoising, (ii) image interpolation based upon a subset of pixels selected uniformly at random, and (iii) learning dictionaries for compressive sensing and also compressive sensing inversion. We have also considered jointly performing (i) and (ii). Important advantages of the proposed methods are: (i) a full posterior on model parameters are inferred, and therefore “error bars” may be placed on the inverted images; (ii) the noise variance need not be known, it is inferred within the analysis and may be nonstationary, and it may be inferred in the presence of significant missing pixels; (iii) while training data may be used to initialize the dictionary learning, this is not needed, and the BPFA results are highly competitive even based upon random initializations. In the context of compressive sensing, the DP-BPFA and PSBP-BPFA results are state of the art, significantly better than existing published methods. Finally, based upon the learned dictionary, a simple method has been constituted for optimizing the CS projections.

The interpolation problem is related to CS, in that we exploit the fact that  $\{\mathbf{x}_i\}_{i=1,N}$  reside on a low-dimensional subspace of  $\mathbb{R}^P$ , such that the total number of measurements is small relative to  $N \cdot P$  (recall  $\mathbf{x}_i \in \mathbb{R}^P$ ). However, in CS one employs projection measurements  $\Sigma \mathbf{x}_i$ , where  $\Sigma \in \mathbb{R}^{n \times P}$ , ideally with  $n \ll P$ . The interpolation problem corresponds to the special case in which the rows of  $\Sigma$  are randomly selected rows of the  $P \times P$  identity matrix. This problem is closely related to the problem of matrix completion [6], [23], [34], where the incomplete matrix  $\mathbf{X} \in \mathbb{R}^{P \times N}$  has columns defined by  $\{\mathbf{x}_i\}_{i=1,N}$ .

While the PSBP-BPFA successfully segmented the image while performing denoising and interpolation of missing pixels, we found that the PSNR performance of direct BPFA analysis performed very close to that of PSBP-BPFA in those applications. The use of PSBP-BPFA utilizes the spatial location of the image patches employed in the analysis, and therefore it removes the exchangeability assumption associated with the simple BPFA (the location of the patches may be interchanged within the BPFA, without affecting the inference). However, since in the denoising and interpolation problems we have many *overlapping*

patches, the extra information provided by PSBP-BPFA does not appear to be significant. By contrast, in the CS inversion problem we do not have overlapping patches, and PSBP-BPFA provided significant performance gains relative to BPFA alone.

#### APPENDIX: GIBBS SAMPLING INFERENCE

The Gibbs sampling update equations are given below; we provide the update equations for the BPFA, and the DP and PSBP versions are relatively simple extensions. Below,  $\Sigma_i$  represents the projection matrix on the data, for image patch  $x_i$ . For the CS problem,  $\Sigma_i$  is typically fully populated, while for the interpolation problem each row of  $\Sigma_i$  is all zeros except for a single one, corresponding to the specific pixel that is measured. The update equations are the conditional probability of each parameter, conditioned on all other parameters in the model.

#### Sample $d_k$

$$p(d_k|-) \propto \prod_{i=1}^N \mathcal{N}(y_i; \Sigma_i D(s_i \odot z_i), \gamma_\epsilon^{-1} I_{\|\Sigma_i\|_0}) \mathcal{N}(d_k; 0, P^{-1} I_P)$$

It can be shown that  $d_k$  can be drawn from a normal distribution

$$p(d_k|-) \sim \mathcal{N}(\mu_{d_k}, \Sigma_{d_k})$$

with the covariance  $\Sigma_{d_k}$  and mean  $\mu_{d_k}$  expressed as

$$\begin{aligned} \Sigma_{d_k} &= \left( P I + \gamma_\epsilon \sum_{i=1}^N z_{ik}^2 s_{ik}^2 \Sigma_i^T \Sigma_i \right)^{-1} \\ \mu_{d_k} &= \gamma_\epsilon \Sigma_{d_k} \sum_{i=1}^N z_{ik} s_{ik} \tilde{x}_i^{-k} \end{aligned}$$

where

$$\tilde{x}_i^{-k} = \Sigma_i^T y_i - \Sigma_i^T \Sigma_i D(s_i \odot z_i) + \Sigma_i^T \Sigma_i d_k (s_{ik} \odot z_{ik}).$$

**Sample  $z_k$ :**  $[z_{1k}, z_{2k}, \dots, z_{Nk}]$

$$p(z_{ik}|-) \propto \mathcal{N}(y_i; \Sigma_i D(s_i \odot z_i), \gamma_\epsilon^{-1} I_{\|\Sigma_i\|_0}) \text{Bernoulli}(z_{ik}; \pi_k)$$

The posterior probability that  $z_{ik} = 1$  is proportional to

$$p_1 = \pi_k \exp \left[ -\frac{\gamma_\epsilon}{2} (s_{ik}^2 d_k^T \Sigma_i^T \Sigma_i d_k - 2 s_{ik} d_k^T \tilde{x}_i^{-k}) \right]$$

and the posterior probability that  $z_{ik} = 0$  is proportional to

$$p_0 = 1 - \pi_k$$

so  $z_{ik}$  can be drawn from a Bernoulli distribution as

$$z_{ik} \sim \text{Bernoulli}\left(\frac{p_1}{p_0 + p_1}\right). \quad (11)$$

**Sample  $s_k:$**   $[s_{1k}, s_{2k}, \dots, s_{Nk}]$

$$p(s_{ik}|-) \propto \mathcal{N}(\mathbf{y}_i; \boldsymbol{\Sigma}_i \mathbf{D}(s_i \odot \mathbf{z}_i), \gamma_\epsilon^{-1} \mathbf{I}_{\|\boldsymbol{\Sigma}_i\|_0}) \mathcal{N}(\mathbf{s}_i; 0, \gamma_s^{-1} \mathbf{I}_K)$$

It can be shown that  $s_{ik}$  can be drawn from a normal distribution

$$p(s_{ik}|-) \sim \mathcal{N}(\mu_{s_{ik}}, \Sigma_{s_{ik}}) \quad (12)$$

with the variance  $\Sigma_{s_{ik}}$  and mean  $\mu_{s_{ik}}$  expressed as

$$\Sigma_{s_{ik}} = (\gamma_s + \gamma_\epsilon z_{ik}^2 \mathbf{d}_k^T \boldsymbol{\Sigma}_i^T \boldsymbol{\Sigma}_i \mathbf{d}_k)^{-1}$$

$$\mu_{s_{ik}} = \gamma_\epsilon \sum_{s_{ik}} z_{ik} \mathbf{d}_k^T \boldsymbol{\Sigma}_i^T \boldsymbol{\Sigma}_i \tilde{\mathbf{x}}_i^{-k}.$$

Note  $z_{ik}$  is equal to either 1 or 0,  $\Sigma_{s_{ik}}$  and  $\mu_{s_{ik}}$  can be further expressed as

$$\begin{aligned} \Sigma_{s_{ik}} &= \begin{cases} (\gamma_s + \gamma_\epsilon \mathbf{d}_k^T \boldsymbol{\Sigma}_i^T \boldsymbol{\Sigma}_i \mathbf{d}_k)^{-1} & \text{if } z_{ik} = 1 \\ \gamma_s^{-1} & \text{if } z_{ik} = 0 \end{cases} \\ \mu_{s_{ik}} &= \begin{cases} \gamma_\epsilon \sum_{s_{ik}} z_{ik} \mathbf{d}_k^T \boldsymbol{\Sigma}_i^T \boldsymbol{\Sigma}_i \tilde{\mathbf{x}}_i^{-k} & \text{if } z_{ik} = 1 \\ 0 & \text{if } z_{ik} = 0 \end{cases}. \end{aligned}$$

**Sample  $\pi_k$**

$$p(\pi_k|-) \propto \text{Beta}(\pi_k; a, b) \prod_{i=1}^N \text{Bernoulli}(z_{ik}; \pi_k)$$

It can be shown that  $\pi_k$  can be drawn from a Beta distribution as

$$p(\pi_k|-) \sim \text{Beta}\left(\frac{a}{K} + \sum_{i=1}^N z_{ik}, \frac{b_0(K-1)}{K} + N - \sum_{i=1}^N z_{ik}\right)$$

**Sample  $\gamma_s$**

$$p(\gamma_s|-) \propto \Gamma(\gamma_s; c_0, d_0) \prod_{i=1}^N \mathcal{N}(\mathbf{s}_i; 0, \gamma_s^{-1} \mathbf{I}_K)$$

It can be shown that  $\gamma_s$  can be drawn from a Gamma distribution as

$$p(\gamma_s|-) \sim \Gamma \left( c_0 + \frac{1}{2}KN, d_0 + \frac{1}{2} \sum_{i=1}^N \mathbf{s}_i^T \mathbf{s}_i \right)$$

**Sample**  $\gamma_\epsilon$

$$p(\gamma_\epsilon|-) \propto \Gamma(\gamma_\epsilon; e_0, f_0) \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i; \boldsymbol{\Sigma}_i \mathbf{D}(\mathbf{s}_i \odot \mathbf{z}_i), \gamma_\epsilon^{-1} \mathbf{I}_{\|\boldsymbol{\Sigma}_i\|_0}) \quad (13)$$

It can be shown that  $\gamma_\epsilon$  can be drawn from a Gamma distribution as

$$p(\gamma_\epsilon|-) \sim \Gamma \left( e_0 + \frac{1}{2} \sum_{i=1}^N \|\boldsymbol{\Sigma}_i\|_0, f_0 + \frac{1}{2} \sum_{i=1}^N \|\boldsymbol{\Sigma}_i^T \mathbf{y}_i - \boldsymbol{\Sigma}_i^T \boldsymbol{\Sigma}_i \mathbf{D}(\mathbf{s}_i \odot \mathbf{z}_i)\|_{\ell_2}^2 \right). \quad (14)$$

Note that  $\boldsymbol{\Sigma}_i^T \boldsymbol{\Sigma}_i$  is a sparse identity matrix,  $\boldsymbol{\Sigma}_{d_k}$  is a diagonal matrix, and  $\mathbf{Z}$  is a sparse matrix, it is easy to find that only basic arithmetical operations are needed and many unnecessary calculations can be avoided, leading to fast computation and low memory requirement.

## REFERENCES

- [1] M. Aharon, M. Elad, and A. M. Bruckstein. K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Trans. Signal Processing*, 54:4311–4322, 2006.
- [2] R. G. Baraniuk, V. Cevher, M. F. Duarte, and C. Hegde. Model-based compressive sensing. *IEEE Trans. Information Theory*, pages 1982–2001, 2010.
- [3] R.G. Baraniuk. Compressive sensing. *IEEE Signal Processing Magazine*, 24:118–124, 2007.
- [4] A.M. Bruckstein, D.L. Donoho, and M. Elad. From sparse solutions of systems of equations to sparse modeling of signals and images. *SIAM review*, 51:34–81, 2007.
- [5] E. Candès and T. Tao. Near-optimal signal recovery from random projections: universal encoding strategies? *IEEE Trans. Information Theory*, 52:5406–5425, 2006.
- [6] E. J. Candès and T. Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE Trans. Inform. Theory*, 2010.
- [7] Y. Chung and D.B. Dunson. Nonparametric bayes conditional distribution modeling with variable selection. *Journal of the American Statistical Association*, 104:1646–1660, 2009.
- [8] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian. Image denoising by sparse 3d transform-domain collaborative filtering. *IEEE Trans. Image Processing*, 16:2007, 2007.
- [9] D.L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52:1289–1306, 2006.
- [10] D.L. Donoho, I.M. Johnstone, G. Kerkyacharian, and D. Picard. Wavelet shrinkage: asymptopia. *Journal of the Royal Statistical Society, Ser. B*, pages 371–394, 1995.
- [11] M.F. Duarte, M.A. Davenport, D. Takhar, J.N. Laska, T. Sun, K.F. Kelly, and R.G. Baraniuk. Single-pixel imaging via compressive sampling. *IEEE Signal Processing Magazine*, 2008.
- [12] J.M. Duarte-Carvajalino and G. Sapiro. Learning to sense sparse signals: simultaneous sensing matrix and sparsifying dictionary optimization. *IEEE Transactions on Image Processing*, pages 1395–1408, 2009.

- [13] M. Elad and M. Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Trans. Image Processing*, 15:3736–3745, 2006.
- [14] M. Elad and I. Yavneh. A weighted average of sparse representations is better than the sparsest one alone. *Preprint*, 2010.
- [15] T. Ferguson. A Bayesian analysis of some nonparametric problems. *Annals of Statistics*, 1:209–230, 1973.
- [16] S. Geman and D. Geman. Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Trans. Pattern Anal. Mach. Intell.*, pages 721–741, 1984.
- [17] S. Gleichman and Y.C. Eldar. Blind compressed sensing. *Preprint (on Arxiv.org)*.
- [18] T.L. Griffiths and Z. Ghahramani. Infinite latent feature models and the Indian buffet process. In *Proc. Advances in Neural Information Processing Systems*, pages 475–482, 2005.
- [19] L. He and L. Carin. Exploiting structure in wavelet-based bayesian compressive sensing. *IEEE Trans. Signal Processing*, pages 3488–3497, 2009.
- [20] E.J. Candès J.-F. Cai and Z. Shen. A singular value thresholding algorithm for matrix completion. *SIAM J. on Optimization*, pages 1956–1982, 2008.
- [21] S. Ji, Y. Xue, and L. Carin. Bayesian compressive sensing. *IEEE Trans. Signal Processing*, 56:2346–2356, 2008.
- [22] D. Knowles and Z. Ghahramani. Infinite sparse factor analysis and infinite independent components analysis. In *Proc. International Conference on Independent Component Analysis and Signal Separation*, 2007.
- [23] N.D. Lawrence and R. Urtasun. Non-linear matrix factorization with gaussian processes. In *Proc. International Conference on Machine Learning*, pages 601–608, 2009.
- [24] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online dictionary learning for sparse coding. In *Proc. International Conference on Machine Learning*, 2009.
- [25] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Supervised dictionary learning. In *Proc. Neural Information Processing Systems*, 2008.
- [26] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Non-local sparse models for image restoration. In *Proc. International Conference on Computer Vision*, 2009.
- [27] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Non-local sparse models for image restoration. In *International Conference on Computer Vision*, 2009.
- [28] J. Mairal, M. Elad, and G. Sapiro. Sparse representation for color image restoration. *IEEE Trans. Image Processing*, 17:53–69, 2008.
- [29] J. Mairal, G. Sapiro, , and M. Elad. Learning multiscale sparse representations for image and video restoration. *SIAM Multiscale Modeling and Simulation*, 7:214 – 241, 2008.
- [30] J. Paisley and L. Carin. Nonparametric factor analysis with beta process priors. In *Proc. International Conference on Machine Learning*, 2009.
- [31] R. Raina, A. Battle, H. Lee, B. Packer, and A.Y. Ng. Self-taught learning: transfer learning from unlabeled data. In *Proc. International Conference on Machine Learning*, 2007.
- [32] M. Ranzato, C. Poultney, S. Chopra, and Y. Lecun. Efficient learning of sparse representations with an energy-based model. In *Proc. Neural Information Processing Systems*, 2006.
- [33] L. Ren, L. Du, D. Dunson, and L. Carin. The logistic stick breaking process. *J. Machine Learning Research*, preprint.
- [34] R. Salakhutdinov and A. Mnih. Bayesian probabilistic matrix factorization using markov chain monte carlo. In *Proc. International Conference on Machine Learning*, pages 880–887, 2008.
- [35] J. Sethuraman. A constructive definition of Dirichlet priors. *Statistica Sinica*, 4:639–650, 1994.
- [36] M. Shankar, N.P. Pitsianis, and D.J. Brady. Compressive video sensors using multichannel imagers. *Appl. Opt.*, 49, 2010.

- [37] Y. W. Teh, M.I. Jordan, M.J. Beal, and D.M. Blei. Hierarchical dirichlet processes. *Journal of the American Statistical Association*, 101, 2004.
- [38] R. Thibaux and M. I. Jordan. Hierarchical beta processes and the Indian buffet process. In *Proc. International Conference on Artificial Intelligence and Statistics*, 2007.
- [39] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused lasso. *J. Royal Stat. Soc. Ser. B.*, pages 91–108, 2005.
- [40] M. Tipping. Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, 1:211–244, June 2001.
- [41] J. Wright, A.Y. Yang, A. Ganesh, S.S. Sastry, and Y. Ma. Robust face recognition via sparse representation. *IEEE Trans. Pattern Analysis Machine Intelligence*, 31:210–227, 2009.
- [42] J. Yang, J. Wright, T. Huang, and Y. Ma. Image super-resolution via sparse representation. *IEEE Transactions on Image Processing*, 2009.
- [43] M. Zhou, H. Chen, J. Paisley, L. Ren, G. Sapiro, and L. Carin. Non-parametric bayesian dictionary learning for sparse image representations. In *Proc. Neural Information Processing Systems*, 2009.
- [44] M. Zhou, H. Yang, G. Sapiro, D. Dunson, and L. Carin. Dependent hierarchical beta process for image interpolation and denoising. In *Proc. Artificial Intelligence and Statistics (AISTATS)*, 2011.

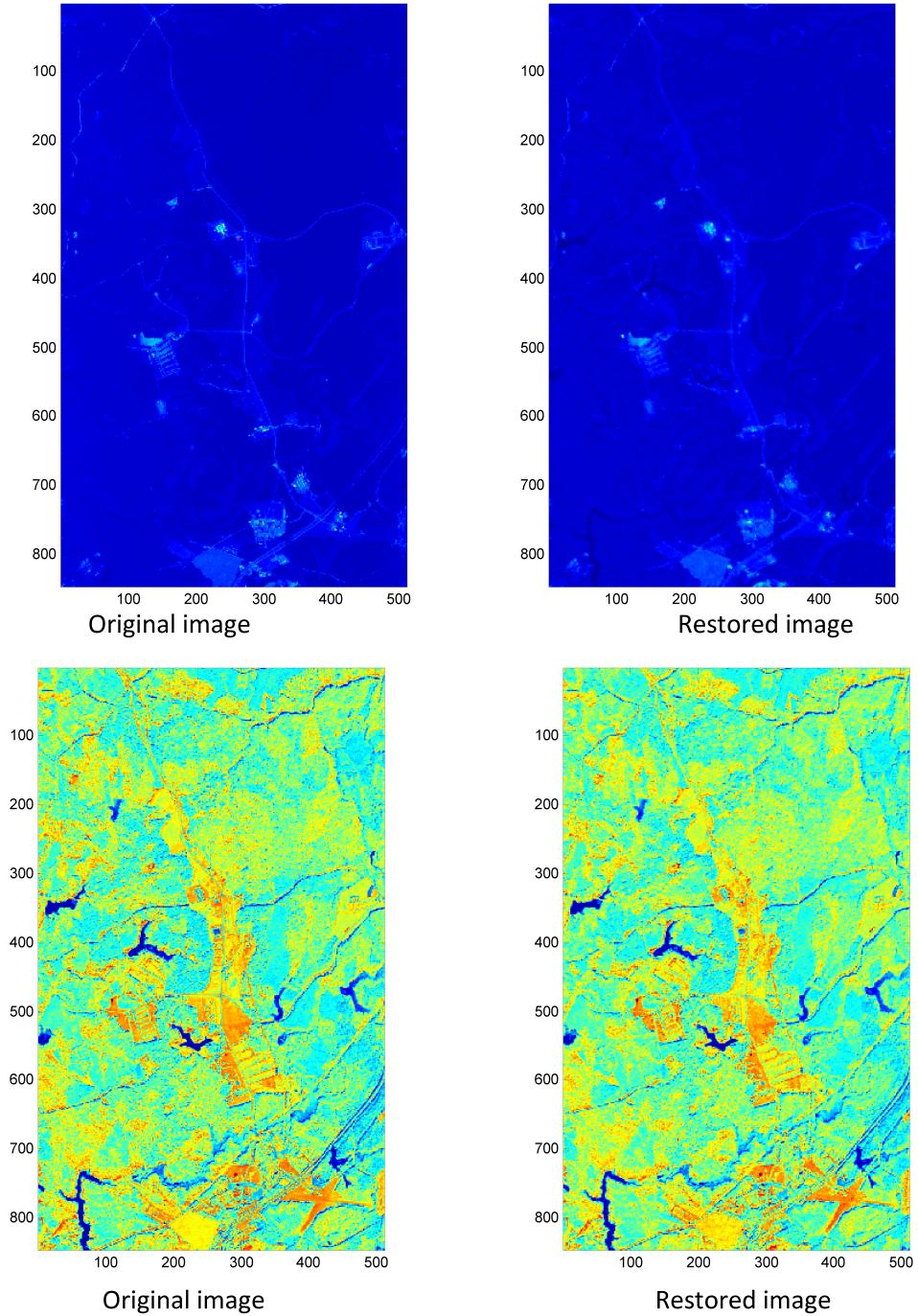


Fig. 6. Comparison of recovered band (average from Gibbs collection iterations) for hyperspectral imagery with 106 spectral bands. The interpolation is performed using 2% of the hyperspectral datacube, selected uniformly at random. The analysis employs  $4 \times 4 \times 106$  patches. All spectral bands are analyzed at once, and here the data (recovered and original) are shown (arbitrarily) for bands 1 (top) and 50 (bottom). Results are computed using the BPFA model.

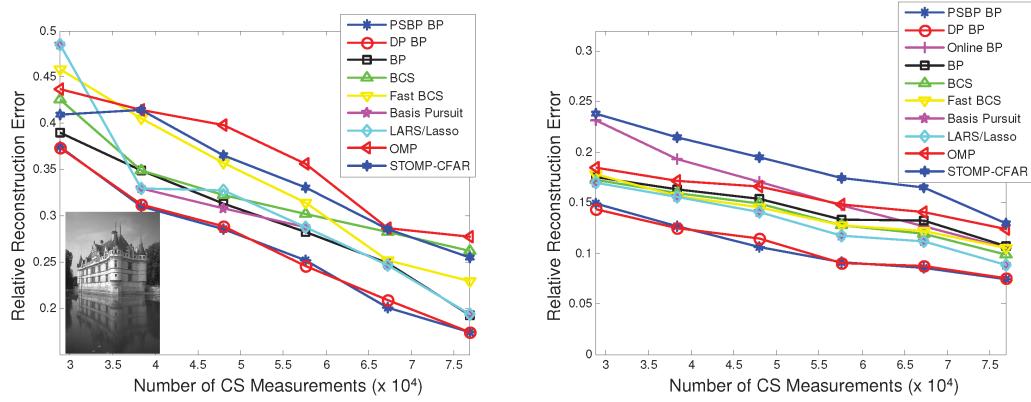


Fig. 7. Left: Compressive sensing (CS) results on grey-scale Castle image, based on a DCT dictionary  $\mathbf{D}$ . The CS projection matrix  $\Sigma$  is constituted randomly, with elements drawn iid from  $\mathcal{N}(0, 1)$ . Results are shown using the DP-BPFA and PSBP-BPFA models in Section IV. Comparisons are also made with several CS inversion algorithms from the literature. Right: Same as on the left but based on a learned dictionary  $\mathbf{D}$  instead of DCT. The online BP results employ BPFA to learn  $\mathbf{D}$  and do CS inversion jointly. All other results are based upon a learned  $\mathbf{D}$  with learning performed offline using distinct training images.

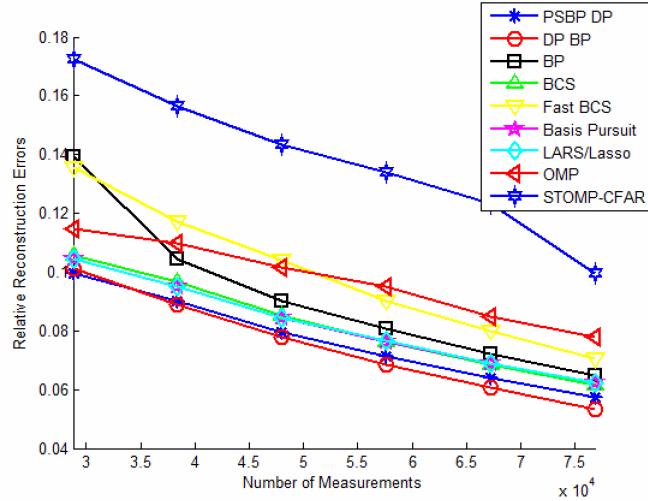


Fig. 8. Compressive sensing (CS) results on grey-scale Castle image, based on a learned dictionary  $\mathbf{D}$  (learning performed offline, using distinct training data). The projection matrix  $\Sigma$  is matched to  $\mathbf{D}$ , based upon an SVD of  $\mathbf{D}$ .