

# Dictionary learning of sound speed profiles

Michael Bianco<sup>a)</sup> and Peter Gerstoft

*Scripps Institution of Oceanography, University of California San Diego, La Jolla, California 92093-0238, USA*

(Received 21 August 2016; revised 6 January 2017; accepted 17 February 2017; published online 13 March 2017)

To provide constraints on the inversion of ocean sound speed profiles (SSPs), SSPs are often modeled using empirical orthogonal functions (EOFs). However, this regularization, which uses the leading order EOFs with a minimum-energy constraint on the coefficients, often yields low resolution SSP estimates. In this paper, it is shown that dictionary learning, a form of unsupervised machine learning, can improve SSP resolution by generating a dictionary of shape functions for sparse processing (e.g., compressive sensing) that optimally compress SSPs; both minimizing the reconstruction error and the number of coefficients. These learned dictionaries (LDs) are not constrained to be orthogonal and thus, fit the given signals such that each signal example is approximated using few LD entries. Here, LDs describing SSP observations from the HF-97 experiment and the South China Sea are generated using the K-SVD algorithm. These LDs better explain SSP variability and require fewer coefficients than EOFs, describing much of the variability with one coefficient. Thus, LDs improve the resolution of SSP estimates with negligible computational burden.

© 2017 Acoustical Society of America. [<http://dx.doi.org/10.1121/1.4977926>]

[ZHM]

Pages: 1749–1758

## I. INTRODUCTION

Inversion for ocean sound speed profiles (SSPs) using acoustic data is a non-linear and highly underdetermined problem.<sup>1</sup> To ensure physically realistic solutions while moderating the size of the parameter search, SSP inversion has often been regularized by modeling SSP as the sum of leading order empirical orthogonal functions (EOFs).<sup>2-7</sup> However, regularization using EOFs often yields low resolution estimates of ocean SSPs, which can be highly variable with fine scale fluctuations. In this paper, it is shown that the resolution of SSP estimates are improved using dictionary learning,<sup>8-13</sup> a form of unsupervised machine learning, to generate a dictionary of regularizing shape functions from SSP data for parsimonious representation of SSPs.

Many signals, including natural images,<sup>14,15</sup> audio,<sup>16</sup> and seismic profiles<sup>17</sup> are well approximated using sparse (few) coefficients, provided a dictionary of shape functions exist under which their representation is sparse. Given a  $K$ -dimensional signal, a dictionary is defined as a set of  $N$ ,  $\ell_2$ -normalized vectors which describe the signal using few coefficients. The sparse processor is then an  $\ell_2$ -norm cost function with an  $\ell_0$ -norm penalty on the number of non-zero coefficients. Signal sparsity is exploited for a number of purposes including signal compression and denoising.<sup>9</sup> Applications of compressive sensing,<sup>18</sup> one approximation to the  $\ell_0$ -norm sparse processor, have in ocean acoustics shown improvements in beamforming,<sup>19-22</sup> geoacoustic inversion,<sup>23</sup> and estimation of ocean SSPs.<sup>24</sup>

Dictionaries that approximate a given class of signals using few coefficients can be designed using dictionary learning.<sup>9</sup> Dictionaries can be generated *ad hoc* from common

shape functions such as wavelets or curvelets, however extensive analysis is required to find an optimal set of prescribed shape functions. Dictionary learning proposes a more direct approach: given enough signal examples for a given signal class, learn a dictionary of shape functions that approximate signals within the class using few coefficients. These learned dictionaries (LDs) have improved compression and denoising results for image and video data over *ad hoc* dictionaries.<sup>9,11</sup> Dictionary learning has been applied to denoising problems in seismics<sup>25</sup> and ocean acoustics,<sup>26,27</sup> as well as to structural acoustic health monitoring.<sup>28</sup>

The K-SVD algorithm,<sup>12</sup> a popular dictionary learning method, finds a dictionary of vectors that optimally partition the data from the training set such that the few dictionary vectors describe each data example. Relative to EOFs which are derived using principal component analysis (PCA),<sup>29,30</sup> these LDs are not constrained to be orthogonal. Thus, LD's provide potentially better signal compression because the vectors are on average, nearer to the signal examples (see Fig. 1).<sup>13</sup>

In this paper, LDs describing one dimensional (1D) ocean SSP data from the HF-97 experiment,<sup>31,32</sup> and from the South China Sea (SCS)<sup>33</sup> are generated using the K-SVD algorithm and the reconstruction performance is evaluated against EOF methods. In Sec. II, EOFs, sparse reconstruction methods, and compression are introduced. In Sec. III, the K-SVD dictionary learning algorithm is explained. In Sec. IV, SSP reconstruction results are given for LDs and EOFs. It is shown that each shape function within the resulting LDs explain more SSP variability than the leading order EOFs trained on the same data. Further, it is demonstrated that SSPs can be reconstructed up to acceptable error using as few as one non-zero coefficient. This compression can improve the resolution of ocean SSP estimates with negligible computational burden.

<sup>a)</sup>Electronic mail: mbianco@ucsd.edu

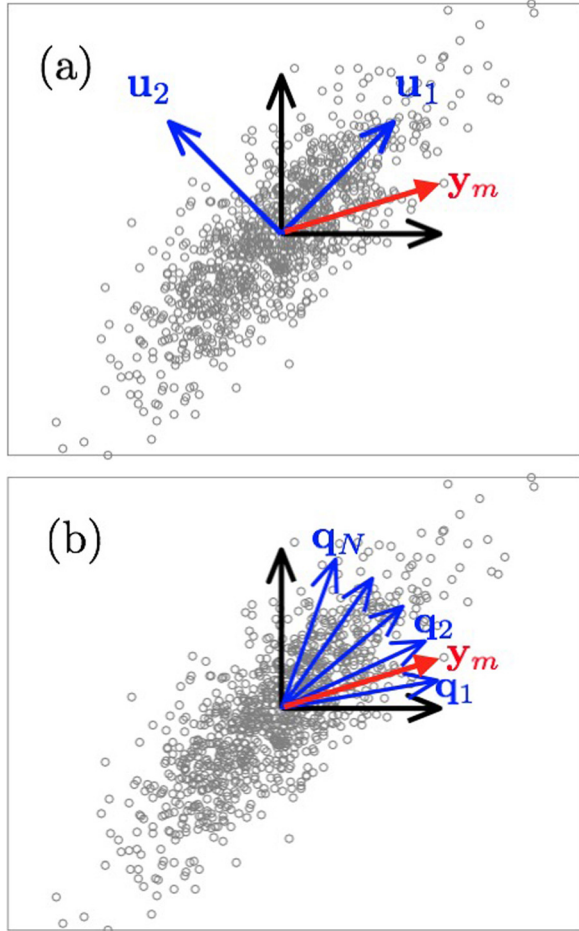


FIG. 1. (Color online) (a) EOF vectors  $[\mathbf{u}_1, \mathbf{u}_2]$  and (b) overcomplete LD vectors  $[\mathbf{q}_1, \dots, \mathbf{q}_N]$  for arbitrary 2D Gaussian distribution relative to arbitrary 2D data observation  $\mathbf{y}_m$ .

*Notation:* In the following, vectors are represented by bold lower-case letters and matrices by bold uppercase letters. The  $\ell_p$ -norm of the vector  $\mathbf{x} \in \mathbb{R}^N$  is defined as  $\|\mathbf{x}\|_p = (\sum_{n=1}^N |x_n|^p)^{1/p}$ . Using similar notation, the  $\ell_0$ -norm is defined as  $\|\mathbf{x}\|_0 = \sum_{n=1}^N |x_n|^0 = \sum_{n=1}^N 1_{|x_n|>0}$ . The  $\ell_p$ -norm of the matrix  $\mathbf{A} \in \mathbb{R}^{K \times M}$  is defined as  $\|\mathbf{A}\|_p = (\sum_{m=1}^M \sum_{k=1}^K |a_k^m|^p)^{1/p}$ . The Frobenius norm ( $\ell_2$ -norm) of the matrix  $\mathbf{A}$  is written as  $\|\mathbf{A}\|_{\mathcal{F}}$ . The hat symbol  $\hat{\cdot}$  appearing above vectors and matrices indicates approximations to the true signals or coefficients.

## II. EOFs AND COMPRESSION

### A. EOFs and PCA

Empirical orthogonal function (EOF) analysis seeks to reduce the dimension of continuously sampled space-time fields by finding spatial patterns which explain much of the variance of the process. These spatial patterns or EOFs correspond to the principal components, from principal component analysis (PCA), of the temporally varying field.<sup>29</sup> Here, the field is a collection of zero-mean ocean SSP anomaly vectors  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_M] \in \mathbb{R}^{K \times M}$ , which are sampled over  $K$  discrete points in depth and  $M$  instants in time. The mean

value of the  $M$  original observations is subtracted to obtain  $\mathbf{Y}$ . The variance of the SSP anomaly at each depth sample  $k$ ,  $\sigma_k^2$ , is defined as

$$\sigma_k^2 = \frac{1}{M} \sum_{m=1}^M (y_m^k)^2, \quad (1)$$

where  $[y_1^k, \dots, y_M^k]$  are the SSP anomaly values at depth sample  $k$  for  $M$  time samples.

The singular value decomposition (SVD)<sup>34</sup> finds the EOFs as the eigenvectors of  $\mathbf{Y}\mathbf{Y}^T$  by

$$\mathbf{Y}\mathbf{Y}^T = \mathbf{P}\mathbf{\Lambda}^2\mathbf{P}^T, \quad (2)$$

where  $\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_L] \in \mathbb{R}^{K \times L}$  are EOFs (eigenvectors) and  $\mathbf{\Lambda}^2 = \text{diag}([\lambda_1^2, \dots, \lambda_L^2]) \in \mathbb{R}^{L \times L}$  are the total variances of the data along the principal directions defined by the EOFs  $\mathbf{p}_l$  with

$$\sum_{k=1}^K \sigma_k^2 = \frac{1}{M} \text{tr}(\mathbf{\Lambda}^2). \quad (3)$$

The EOFs  $\mathbf{p}_l$  with  $\lambda_1^2 \geq \dots \geq \lambda_L^2$  are spatial features of the SSPs which explain the greatest variance of  $\mathbf{Y}$ . If the number of training vectors  $M \geq K$ ,  $L = K$  and  $[\mathbf{p}_1, \dots, \mathbf{p}_L]$  form a basis in  $\mathbb{R}^K$ .

### B. SSP reconstruction using EOFs

Since the leading-order EOFs often explain much of the variance in  $\mathbf{Y}$ , the representation of anomalies  $\mathbf{y}_m$  can be compressed by retaining only the leading order EOFs  $P < L$ ,

$$\hat{\mathbf{y}}_m = \mathbf{Q}_P \hat{\mathbf{x}}_{P,m}, \quad (4)$$

where  $\mathbf{Q}_P \in \mathbb{R}^{K \times P}$  is here the dictionary containing the  $P$  leading-order EOFs and  $\hat{\mathbf{x}}_{P,m} \in \mathbb{R}^P$  is the coefficient vector. Since the entries in  $\mathbf{Q}_P$  are orthonormal, the coefficients are solved by

$$\hat{\mathbf{x}}_{P,m} = \mathbf{Q}_P^T \mathbf{y}_m. \quad (5)$$

For ocean SSPs, usually no more than  $P=5$  EOF coefficients have been used to reconstruct ocean SSPs.<sup>4,7</sup>

### C. Sparse reconstruction

A signal  $\mathbf{y}_m$ , whose model is sparse in the dictionary  $\mathbf{Q}_N = [\mathbf{q}_1, \dots, \mathbf{q}_N] \in \mathbb{R}^{K \times N}$  ( $N$ -entry sparsifying dictionary for  $\mathbf{Y}$ ), is reconstructed to acceptable error using  $T \ll K$  vectors  $\mathbf{q}_n$ .<sup>9</sup> The problem of estimating few coefficients in  $\mathbf{x}_m$  for reconstruction of  $\mathbf{y}_m$  can be phrased using the canonical sparse processor

$$\hat{\mathbf{x}}_m = \arg \min_{\mathbf{x}_m \in \mathbb{R}^N} \|\mathbf{y}_m - \mathbf{Q}\mathbf{x}_m\|_2 \text{ subject to } \|\mathbf{x}_m\|_0 \leq T. \quad (6)$$

The  $\ell_0$ -norm penalizes the number of non-zero coefficients in the solution to a typical  $\ell_2$ -norm cost function. The  $\ell_0$ -norm constraint is non-convex and imposes combinatorial search for the exact solution to Eq. (6). Since exhaustive

search generally requires a prohibitive number of computations, approximate solution methods such as matching pursuit (MP) and basis pursuit (BP) are preferred.<sup>9</sup> In this paper, orthogonal matching pursuit (OMP)<sup>35</sup> is used as the sparse solver. For small  $T$ , OMP achieves similar reconstruction accuracy relative to BP methods, but with much greater speed.<sup>9</sup>

It has been shown that non-orthogonal, overcomplete dictionaries  $\mathbf{Q}_N$  with  $N > K$  (complete,  $N = K$ ) can be designed to minimize both error and number of non-zero coefficients  $T$ , and thus provide greater compression over orthogonal dictionaries.<sup>9,13,16</sup> While overcomplete dictionaries can be designed by concatenating ortho-bases of wavelets or Fourier shape functions, better compression is often achieved by adapting the dictionary to the data under analysis using dictionary learning techniques.<sup>12,13</sup> Since Eq. (6) promotes sparse solutions, it provides criteria for the design of dictionary  $\mathbf{Q}$  for adequate reconstruction of  $\mathbf{y}_m$  with a minimum number of non-zero coefficients. Rewriting Eq. (7) with

$$\min_{\mathbf{Q}} \left\{ \min_{\mathbf{X}} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_{\mathcal{F}}^2 \text{ subject to } \forall_m, \|\mathbf{x}_m\|_0 \leq T \right\}, \quad (7)$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_M]$  is the matrix of coefficient vectors corresponding to examples  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_M]$ , reconstruction error is minimized relative to the dictionary  $\mathbf{Q}$  as well as relative to the sparse coefficients.

In this paper, the K-SVD algorithm, a clustering based dictionary learning method, is used to solve Eq. (7). The K-SVD is an adaptation of the K-means algorithm for vector quantization (VQ) codebook design (a.k.a. the generalized Lloyd algorithm).<sup>16</sup> The LD vectors  $\mathbf{q}_n$  from this technique partition the feature space of the data rather than  $\mathbb{R}^K$ , increasing the likelihood that  $\mathbf{y}_m$  is as a linear combination of few vectors  $\mathbf{q}_n$  in the solution to Eq. (6) (see Fig. 1). By increasing the number of vectors  $N \geq K$  for overcomplete dictionaries, and thus the number of partitions in feature space, the sparsity of the solutions can be increased further.<sup>13</sup>

#### D. Vector quantization

VQ (Ref. 16) compresses a class of  $K$ -dimensional signals  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_M] \in \mathbb{R}^{K \times M}$  by optimally mapping  $\mathbf{y}_m$  to a set of code vectors  $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_N] \in \mathbb{R}^{K \times N}$  for  $N < M$ , called a codebook. The signals  $\mathbf{y}_m$  are then quantized or replaced by the best code vector choice from  $\mathbf{C}$ .<sup>16</sup> The mapping that minimizes mean squared error (MSE) in reconstruction

$$\text{MSE}(\mathbf{Y}, \hat{\mathbf{Y}}) = \frac{1}{N} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_{\mathcal{F}}^2, \quad (8)$$

where  $\hat{\mathbf{Y}} = [\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_M]$  is the vector quantized  $\mathbf{Y}$ , is the assignment of each vector  $\mathbf{y}_m$  to the code vectors  $\mathbf{c}_n$  based on minimum  $\ell_2$ -distance (nearest neighbor metric). Thus the  $\ell_2$ -distances from the code vectors  $\mathbf{c}_n$  define a set of partitions  $(R_1, \dots, R_N) \in \mathbb{R}^K$  (called Voronoi cells)

$$R_n = \{i | \forall_{l \neq n}, \|\mathbf{y}_i - \mathbf{c}_n\|_2 < \|\mathbf{y}_i - \mathbf{c}_l\|_2\}, \quad (9)$$

where if  $\mathbf{y}_i$  falls within the cell  $R_n$ ,  $\hat{\mathbf{y}}_i$  is  $\mathbf{c}_n$ . These cells are shown in Fig. 2(a). This is stated formally by defining a selector function  $S_n$  as

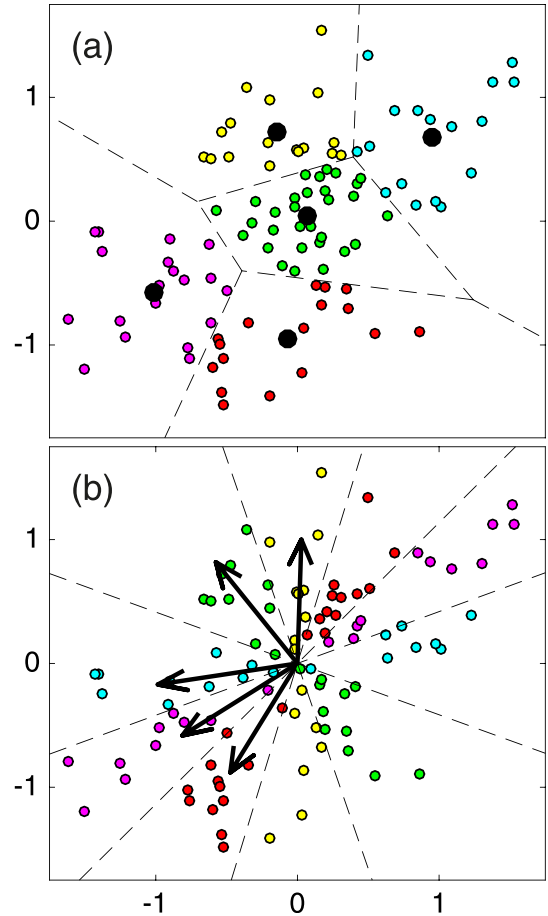


FIG. 2. (Color online) Partitioning of Gaussian random distribution ( $\sigma_1 = 0.75$ ,  $\sigma_2 = 0.5$ ) using (a) five codebook vectors (K-means, VQ) and with (b) five dictionary vectors from dictionary learning (K-SVD,  $T = 1$ ).

$$S_n(\mathbf{y}_m) = \begin{cases} 1 & \text{if } \mathbf{y}_m \in R_n \\ 0 & \text{otherwise.} \end{cases} \quad (10)$$

The vector quantization step is then

$$\hat{\mathbf{y}}_m = \sum_{n=1}^N S_n(\mathbf{y}_m) \mathbf{c}_n. \quad (11)$$

The operations in Eqs. (9) and (10) are analogous to solving the sparse minimization problem:

$$\hat{\mathbf{x}}_m = \arg \min_{\mathbf{x}_m \in \mathbb{R}^N} \|\mathbf{y}_m - \mathbf{C}\mathbf{x}_m\|_2 \text{ subject to } \|\mathbf{x}_m\|_0 = 1, \quad (12)$$

where the non-zero coefficients  $x_m^n = 1$ . In this problem, selection of the coefficient in  $\mathbf{x}_m$  corresponds to mapping the observation vector  $\mathbf{y}_m$  to  $\mathbf{c}_n$ , similar to the selector function  $S_n$ . The vector quantized  $\mathbf{y}_m$  is thus written, alternately from Eq. (11), as

$$\hat{\mathbf{y}}_m = \mathbf{C}\hat{\mathbf{x}}_m. \quad (13)$$

#### E. K-means

Given the MSE metric [Eq. (8)], VQ codebook vectors  $[\mathbf{c}_1, \dots, \mathbf{c}_N]$  which correspond to the centroids of the data

$\mathbf{Y}$  within  $(R_1, \dots, R_N)$  minimize the reconstruction error. The assignment of  $\mathbf{c}_n$  as the centroid of  $\mathbf{y}_j \in R_n$  is

$$\mathbf{c}_n = \frac{1}{|R_n|} \sum_{j \in R_n} \mathbf{y}_j, \quad (14)$$

where  $|R_n|$  is the number of vectors  $\mathbf{y}_j \in R_n$ .

The K-means algorithm shown in Table I, iteratively updates  $\mathbf{C}$  using the centroid condition Eq. (14) and the  $\ell_2$  nearest-neighbor criteria Eq. (9) to optimize the code vectors for VQ. The algorithm requires an initial codebook  $\mathbf{C}^0$ . For example,  $\mathbf{C}^0$  can be  $N$  random vectors in  $\mathbb{R}^K$  or selected observations from the training set  $\mathbf{Y}$ . The K-means algorithm is guaranteed to improve or leave unchanged the MSE distortion after each iteration and converges to a local minimum.<sup>12,16</sup>

### III. DICTIONARY LEARNING

Two popular algorithms for dictionary learning, the method of optimal directions (MOD)<sup>13</sup> and the K-SVD,<sup>12</sup> are inspired by the iterative K-means codebook updates for VQ (Table I). The  $N$  columns of the dictionary  $\mathbf{Q}$ , like the entries in codebook  $\mathbf{C}$ , correspond to partitions in  $\mathbb{R}^K$ . However, they are constrained to have unit  $\ell_2$ -norm and thus separate the magnitude (coefficients  $\mathbf{x}_m$ ) from the shapes (dictionary entries  $\mathbf{q}_n$ ) for the sparse processing objective Eq. (6). When  $T = 1$ , the  $\ell_2$ -norm in Eq. (6) is minimized by the dictionary entry  $\mathbf{q}_n$  that has the greatest inner product with example  $\mathbf{y}_m$ .<sup>9</sup> Thus for  $T = 1$ ,  $[\mathbf{q}_1, \dots, \mathbf{q}_N]$  define radial partitions of  $\mathbb{R}^K$ . These partitions are shown in Fig. 2(b) for a hypothetical 2D ( $K = 2$ ) random data set. This corresponds to a special case of VQ, called gain-shape VQ.<sup>16</sup> However, for sparse processing, only the shapes of the signals are quantized. The gains, which are the coefficients  $\mathbf{x}_m$ , are solved. For  $T > 1$ , the sparse solution is analogous to VQ, assigning examples  $\mathbf{y}_m$  to dictionary entries in  $\mathbf{Q}$  for up to  $T$  non-zero coefficients in  $\mathbf{x}_m$ .

Given these relationships between sparse processing with dictionaries and VQ, the MOD<sup>13</sup> and K-SVD<sup>12</sup> algorithms attempt to generalize the K-means algorithm to

optimization of dictionaries for sparse processing for  $T \geq 1$ . They are two-step algorithms which reflect the two update steps in the K-means codebook optimization: (1) partition data  $\mathbf{Y}$  into regions  $(R_1, \dots, R_N)$  corresponding to  $\mathbf{c}_n$  and (2) update  $\mathbf{c}_n$  to centroid of examples  $\mathbf{y}_m \in R_n$ . The K-means algorithm is generalized to the dictionary learning problem Eq. (7) as two steps:

- (1) Sparse coding: Given dictionary  $\mathbf{Q}$ , solve for up to  $T$  non-zero coefficients in  $\mathbf{x}_m$  corresponding to examples  $\mathbf{y}_m$  for  $m = [1, \dots, M]$ .
- (2) Dictionary update: Given coefficients  $\mathbf{X}$ , solve for  $\mathbf{Q}$  which minimizes reconstruction error for  $\mathbf{Y}$ .

The sparse coding step (1), which is the same for both MOD and K-SVD, is accomplished using any sparse solution method, including matching pursuit and basis pursuit. The algorithms differ in the dictionary update step.

#### A. The K-SVD algorithm

The K-SVD algorithm is here chosen for its computational efficiency, speed, and convergence to local minima (at least for  $T = 1$ ). The K-SVD algorithm sequentially optimizes the dictionary entries  $\mathbf{q}_n$  and coefficients  $\mathbf{x}_m$  for each update step using the SVD, and thus also avoids the matrix inverse. For  $T = 1$ , the sequential updates of the K-SVD provide optimal dictionary updates for gain-shape VQ.<sup>12,16</sup> Optimal updates to the gain-shape dictionary will, like K-means updates, either improve or leave unchanged the MSE and convergence to a local minimum is guaranteed. For  $T > 1$ , convergence of the K-SVD updates to a local minimum depends on the accuracy of the sparse-solver used in the sparse coding stage.<sup>12</sup>

In the K-SVD algorithm, each dictionary update step  $i$  sequentially improves both the entries  $\mathbf{q}_n \in \mathbf{Q}^i$  and the coefficients in  $\mathbf{x}_m \in \mathbf{X}^i$ , without change in support. Expressing the coefficients as row vectors  $\mathbf{x}_T^n \in \mathbb{R}^N$  and  $\mathbf{x}_T^j \in \mathbb{R}^N$ , which relate all examples  $\mathbf{Y}$  to  $\mathbf{q}_n$  and  $\mathbf{q}_j$ , respectively, the  $\ell_2$ -penalty from Eq. (7) is rewritten as

$$\begin{aligned} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_{\mathcal{F}}^2 &= \left\| \mathbf{Y} - \sum_{n=1}^N \mathbf{q}_n \mathbf{x}_T^n \right\|_{\mathcal{F}}^2 \\ &= \|\mathbf{E}_j - \mathbf{q}_j \mathbf{x}_T^j\|_{\mathcal{F}}^2, \end{aligned} \quad (15)$$

where

$$\mathbf{E}_j = \left( \mathbf{Y} - \sum_{n \neq j} \mathbf{q}_n \mathbf{x}_T^n \right). \quad (16)$$

Thus, in Eq. (15) the  $\ell_2$ -penalty is separated into an error term  $\mathbf{E}_j = [\mathbf{e}_{j,1}, \dots, \mathbf{e}_{j,M}] \in \mathbb{R}^{K \times M}$ , which is the error for all examples  $\mathbf{Y}$  if  $\mathbf{q}_j$  is excluded from their reconstruction, and the product of the excluded entry  $\mathbf{q}_j$  and coefficients  $\mathbf{x}_T^j \in \mathbb{R}^N$ .

An update to the dictionary entry  $\mathbf{q}_j$  and coefficients  $\mathbf{x}_T^j$  which minimizes Eq. (15) is found by taking the SVD of  $\mathbf{E}_j$ , which provides the best rank-1 approximation of  $\mathbf{E}_j$ . However, many of the entries in  $\mathbf{x}_T^j$  are zero (corresponding to examples which do not use  $\mathbf{q}_j$ ). To properly update  $\mathbf{q}_j$  and

TABLE I. The K-means algorithm (Ref. 16).

Given: training vectors $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_M] \in \mathbb{R}^{K \times M}$	
Initialize: index $i = 0$ , codebook $\mathbf{C}^0 = [\mathbf{c}_1^0, \dots, \mathbf{c}_N^0] \in \mathbb{R}^{K \times N}$ , MSE <sup>0</sup> solving Eq. (8)-(11)	
I:	Update codebook
	1. Partition $\mathbf{Y}$ into $N$ regions $(R_1, \dots, R_N)$ by
	$R_n = \{i   \forall i \neq n, \ \mathbf{y}_i - \mathbf{c}_n^i\ _2 < \ \mathbf{y}_i - \mathbf{c}_j^i\ _2\}$ [Eq. (9)]
	2. Make code vectors centroids of $\mathbf{y}_j$ in partitions $R_n$
	$\mathbf{c}_n^{i+1} = \frac{1}{ R_n^i } \sum_{j \in R_n^i} \mathbf{y}_j$
II:	Check error
	1. Calculate MSE <sup><math>i+1</math></sup> from updated codebook $\mathbf{C}^{i+1}$
	2. If $ \text{MSE}^{i+1} - \text{MSE}^i  < \eta$
	$i = i + 1$ , return to I
	else
	end

$\mathbf{x}_T^j$  with SVD, Eq. (15) must be restricted to examples  $\mathbf{y}_m$  which use  $\mathbf{q}_j$ ,

$$\|\mathbf{E}_j^R - \mathbf{q}_j \mathbf{x}_R^j\|_{\mathcal{F}}^2, \quad (17)$$

where  $\mathbf{E}_j^R$  and  $\mathbf{x}_R^j$  are entries in  $\mathbf{E}_j$  and  $\mathbf{x}_T^j$ , respectively, corresponding to examples  $\mathbf{y}_m$  which use  $\mathbf{q}_j$ , and are defined as

$$\mathbf{E}_j^R = \{\mathbf{e}_{j,l} | \forall l, x_l^j \neq 0\}, \quad \mathbf{x}_R^j = \{x_l^j | \forall l, x_l^j \neq 0\}. \quad (18)$$

Thus for each K-SVD iteration, the dictionary entries and coefficients are sequentially updated as the SVD of  $\mathbf{E}_j^R = \mathbf{USV}^T$ . The dictionary entry  $\mathbf{q}_j^i$  is updated with the first column in  $\mathbf{U}$  and the coefficient vector  $\mathbf{x}_R^j$  is updated as the product of the first singular value  $\mathbf{S}(1, 1)$  with the first column of  $\mathbf{V}$ . The K-SVD algorithm is given in Table II.

The dictionary  $\mathbf{Q}$  is initialized using  $N$  randomly selected,  $\ell_2$ -normalized examples from  $\mathbf{Y}$ .<sup>9,12</sup> During the iterations, one or more dictionary entries may become unused. If this occurs, the unused entries are replaced using the most poorly represented examples  $\mathbf{y}_m$  ( $\ell_2$ -normlized), determined by reconstruction error.

#### IV. EXPERIMENTAL RESULTS

To demonstrate the usefulness of the dictionary learning approach, we here analyze two data sets: (1) thermistor data from the HF-97 acoustics experiment,<sup>31,32</sup> conducted off the coast of Point Loma, CA and (2) conductivity, temperature, and depth (CTD) data collected across the Luzon Strait near the South China Sea (SCS).<sup>33</sup> Training data  $\mathbf{Y}$  were derived from the data sets by converting raw thermistor and CTD data to SSPs and subtracting the mean. The HF-97 thermistor data were recorded every 15 s, over a 48 h period, from 14 to 70 m depth, with 4 m spacing (15 points). The full 11 488 profile

data set was down-sampled to  $M = 1000$  profiles for the training set, and SSPs were interpolated to  $K = 30$  points using a shape-preserving cubic spline. The SCS CTD data were recorded at about 1 m resolution from 116 to 496 m depth (384 points). From the SCS data set,  $M = 755$  profiles were used as the training set, and the profiles were uniformly down-sampled to  $K = 50$  points. The SSP data sets are shown in Fig. 3. Both data sets have small and large spatiotemporal variations.

EOFs were calculated from the SVD [Eq. (2)] and LDs (learned dictionaries) were generated with the K-SVD algorithm (Table II), using OMP for the sparse coding stage. The number of non-zero coefficients solved with OMP for each dictionary was held fixed at exactly  $T$  non-zero coefficients. The initial dictionary  $\mathbf{Q}^0$  was populated using randomly selected examples from the training sets  $\mathbf{Y}$ .

#### A. Learning SSP dictionaries from data

Here, LDs and EOFs were generated using the full SSP data from HF-97 ( $M = 1000$ ) and SCS ( $M = 755$ ). The EOFs and LDs from HF-97 are shown in Figs. 4 and 5 and from the SCS in Fig. 6. The HF-97 LD, with  $N = K$  and  $T = 1$ , is compared to the EOFs ( $K = 30$ ) in Fig. 4. Only the leading order EOFs [Fig. 4(a)] are informative of ocean SSP variability whereas all shape functions in the LD [Fig. 4(b)] are informative [Figs. 4(c)–4(d)]. This behavior is also evident

TABLE II. The K-SVD Algorithm (Ref. 12).

Given: $\mathbf{Y} \in \mathbb{R}^{K \times M}$ , $\mathbf{Q}^0 \in \mathbb{R}^{K \times N}$ , $T \in \mathbb{N}$ , and $i = 0$	
Repeat until convergence:	
1.	Sparse coding
	for $m = 1: M$
	solve Eq. (6) using any sparse solver
a:	$\hat{\mathbf{x}}_m = \arg \min_{\mathbf{x}_m \in \mathbb{R}^N} \ \mathbf{y}_m - \mathbf{Q}^i \mathbf{x}_m\ _2$ subject to $\ \mathbf{x}_m\ _0 \leq T$
	end
b:	$\mathbf{X} = [\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_M]$
2.	Dictionary update
	for $j = 1: N$
a:	compute reconstruction error $\mathbf{E}_j$ as
	$\mathbf{E}_j = \mathbf{Y} - \sum_{n \neq j} \mathbf{q}_n^i \mathbf{x}_T^n$
b:	obtain $\mathbf{E}_j^R$ , $\mathbf{x}_R^j$ corresponding to nonzero $\mathbf{x}_T^j$
c:	apply SVD to $\mathbf{E}_j^R$
	$\mathbf{E}_j^R = \mathbf{USV}^T$
d:	update $\mathbf{q}_j^i$ : $\mathbf{q}_j^i = \mathbf{U}(:, 1)$
e:	update $\mathbf{x}_R^j$ : $\mathbf{x}_R^j = \mathbf{V}(:, 1)\mathbf{S}(1, 1)$
	end
f:	$\mathbf{Q}^{i+1} = \mathbf{Q}^i$
	$i = i + 1$

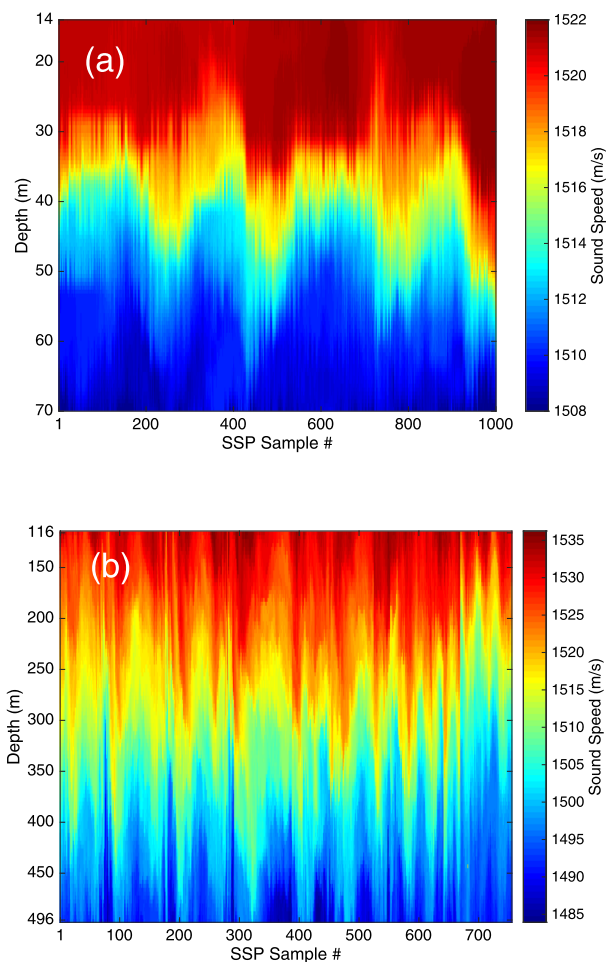


FIG. 3. (Color online) Sound speed profile (SSP) data from (a) HF-97 and (b) SCS experiments.

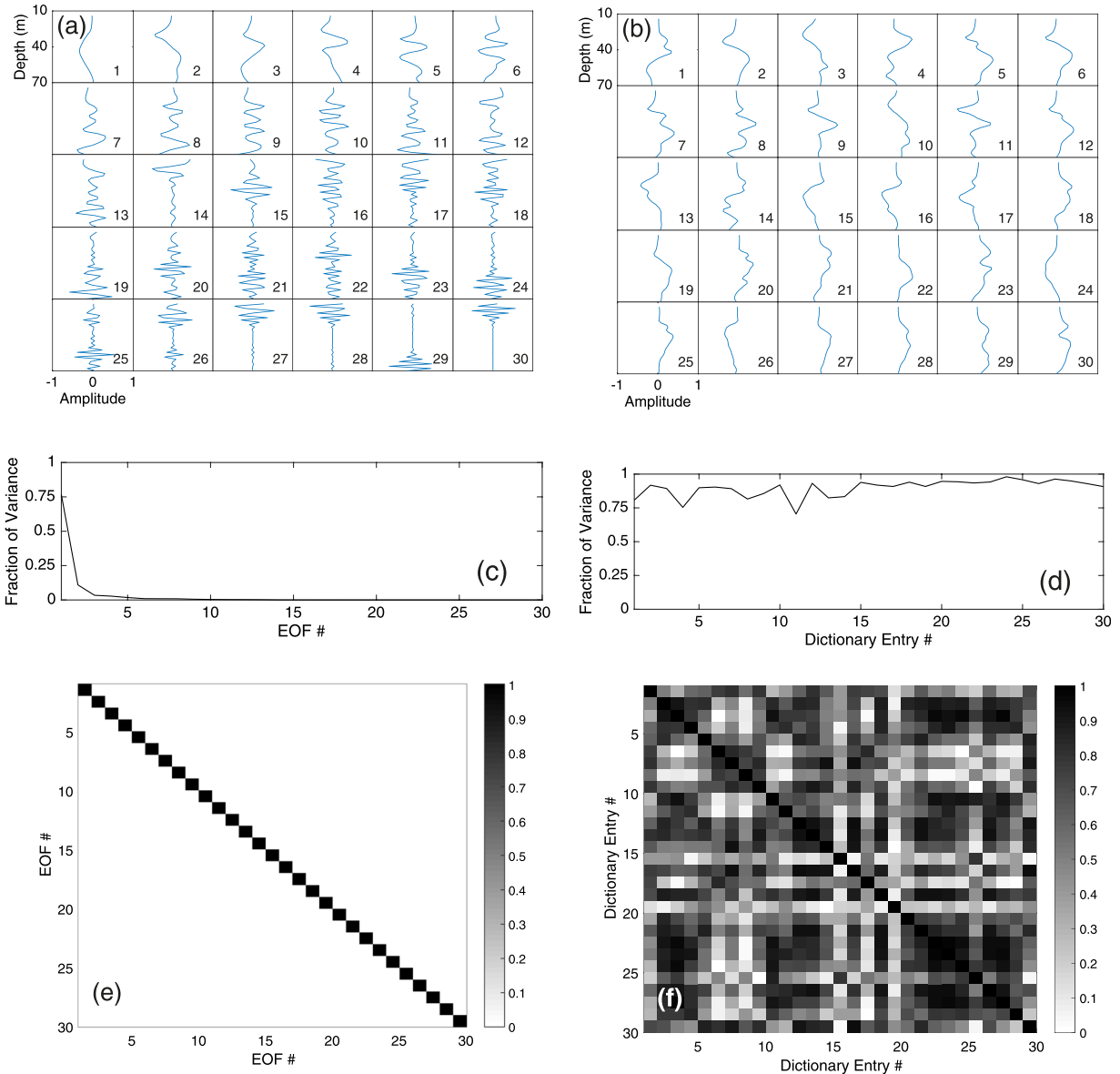


FIG. 4. (Color online) HF-97: (a) EOFs and (b) LD entries ( $N=K$  and  $T=1$ , sorted by variance  $\sigma_{q_n}^2$ ). Fraction of (c) total SSP variance explained by EOFs and (d) SSP variance explained for examples using LD entries. Coherence of (e) EOFs and (f) LD entries.

for the SCS data set (Fig. 6). The EOFs ( $K=50$ ) calculated from the full training set are shown in Fig. 6(a), and the LD entries for  $N=50$  and  $T=1$  sparse coefficient are shown in Fig. 6(b). The overcomplete LDs for the HF-97 data shown in Fig. 5 and for the SCS data in Fig. 6(c).

As illustrated in Fig. 1, by relaxing the requirement of orthogonality for the shape functions, the shape functions can better fit the data and thereby achieve greater compression. The Gram matrix  $\mathbf{G}$ , which gives the coherence of matrix columns, is defined for a matrix  $\mathbf{A}$  with unit  $\ell_2$ -norm columns as  $\mathbf{G} = |\mathbf{A}^T \mathbf{A}|$ . The Gram matrix for the EOFs [Fig. 4(e)] shows the shapes in the EOF dictionary are orthogonal ( $\mathbf{G} = \mathbf{I}$ , by definition), whereas those of the LD [Fig. 4(f)] are not.

## B. Reconstruction of SSP training data

In this section, EOFs and LDs are trained on the full SSP data sets  $\mathbf{Y} = [y_1, \dots, y_M]$ . Reconstruction performance

of the EOF and LDs are then evaluated on SSPs within the training set, using a mean error metric.

The coefficients for the learned  $\mathbf{Q}$  and initial  $\mathbf{Q}^0$  dictionaries  $\hat{\mathbf{x}}_m$  are solved from the sparse objective [Eq. (6)] using OMP. The least squares (LS) solution for the  $T$  leading-order coefficients  $\mathbf{x}_L \in \mathbb{R}^T$  from the EOFs  $\mathbf{P}$  were solved by Eq. (5). The best combination of  $T$  EOF coefficients was solved from the sparse objective [Eq. (6)] using OMP. Given the coefficients  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]$  describing examples  $\mathbf{Y} = [y_1, \dots, y_m]$ , the reconstructed examples  $\hat{\mathbf{Y}} = [\hat{y}_1, \dots, \hat{y}_m]$  are given by  $\hat{\mathbf{Y}} = \mathbf{Q}\hat{\mathbf{X}}$ . The mean reconstruction error (ME) for the training set is then

$$\text{ME} = \frac{1}{KM} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_1. \quad (19)$$

We here use the  $\ell_1$ -norm to stress the robustness of the LD reconstruction.

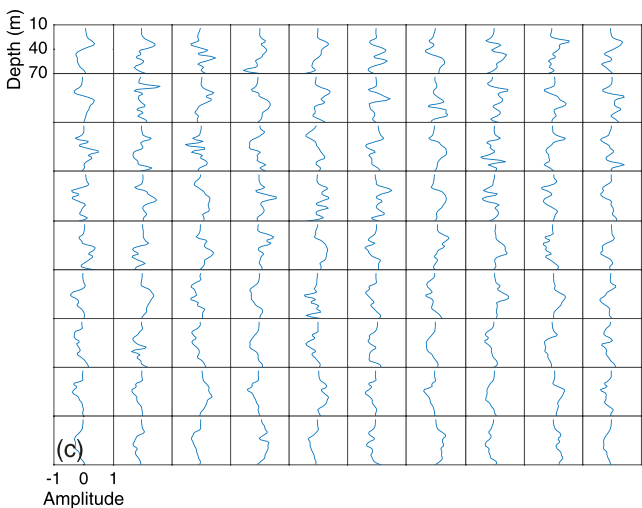
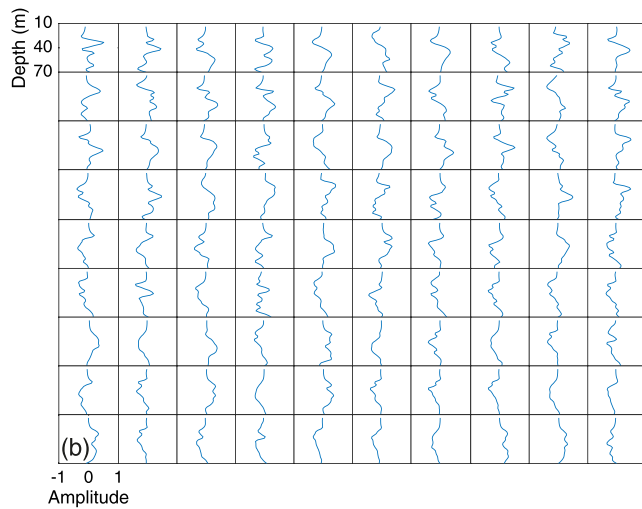
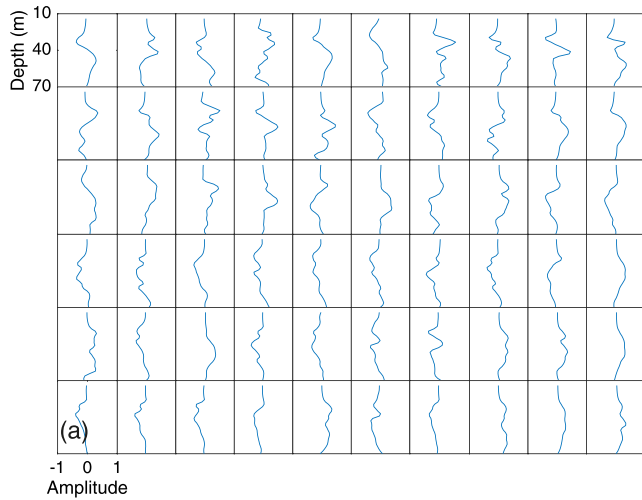


FIG. 5. (Color online) HF-97: LD entries (a)  $N=60$  and  $T=1$ , (b)  $N=90$  and  $T=1$ , and (c)  $N=90$  and  $T=5$ . Dictionary entries are sorted in descending variance  $\sigma_{q_n}^2$ .

To illustrate the optimality of LDs for SSP compression, the K-SVD algorithm was run using EOFs as the initial dictionary  $\mathbf{Q}^0$  for  $T=1$  non-zero coefficient. The convergence of ME for the K-SVD iterations is shown in Fig. 7(a). After

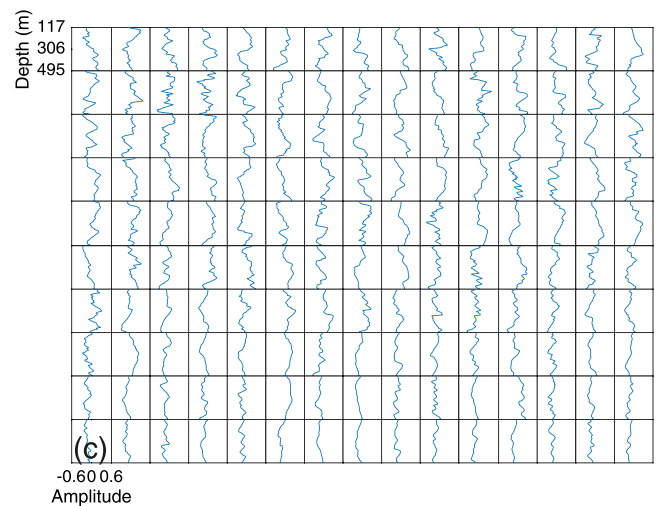
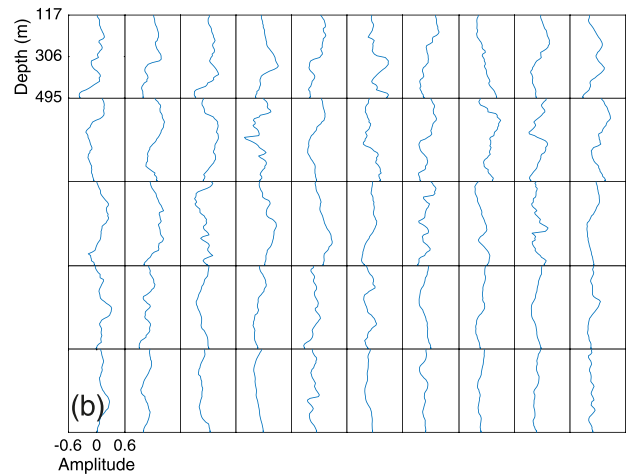
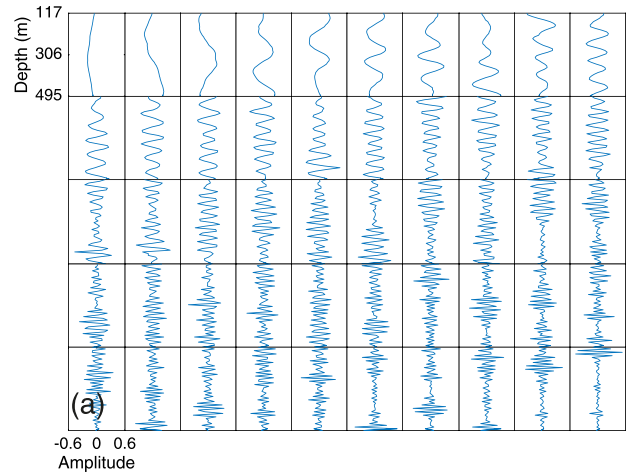


FIG. 6. (Color online) SCS: EOFs (a) and LD entries, (b)  $N=K=50$  and  $T=1$ , and (c)  $N=150$  and  $T=1$ . Dictionary entries are sorted in descending variance  $\sigma_{q_n}^2$ .

30 K-SVD iterations, the mean error of the  $M=1000$  profile training set is decreased by nearly half. The convergence is much faster for  $\mathbf{Q}^0$  consisting of randomly selected examples from  $\mathbf{Y}$ .

For LDs, increasing the number of entries  $N$  or increasing the number of sparse coefficients  $T$  will always reduce the reconstruction error ( $N$  and  $T$  are decided with computational

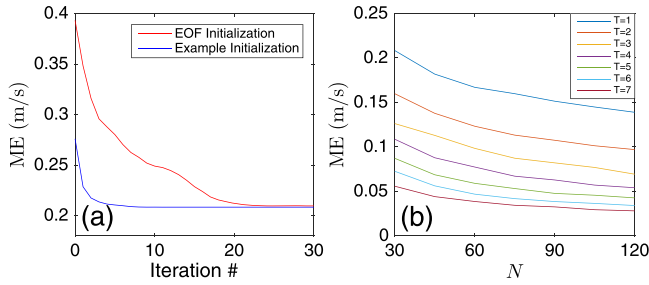


FIG. 7. (Color online) HF-97: (a) Convergence of LD ( $N = 30$ ,  $T = 1$ ) mean reconstruction error (ME), initialized using EOFs or  $N$  randomly selected examples from  $\mathbf{Y}$ . (b) ME versus non-zero coefficients  $T$  and number of dictionary entries  $N$ .

considerations). The effect of  $N$  and  $T$  on the mean reconstruction error for the HF-97 data is shown in Fig. 7(b). The errors are calculated for the range  $N = K$  to  $N = 4K$  and the dictionaries were optimized to use a fixed number non-zero coefficients ( $T$ ).

The reconstruction error using the EOF dictionary is compared to results from LDs  $\mathbf{Q}$  with  $N = 3K$ , using  $T$  non-zero coefficients. In Figs. 8(a) and 8(c) results are shown for the HF-97 ( $N = 90$ ) and SCS ( $N = 150$ ) data, respectively. Coefficients describing each example  $\mathbf{y}_m$ , were solved (1) from the LD  $\mathbf{Q}$ , (2) from  $\mathbf{Q}^0$ , the dictionary consisting of  $N$  randomly chosen examples from the training set (to illustrate improvements in reconstruction error made in the K-SVD iterations), (3) the leading order EOFs, and (4) the best combination of EOFs. The mean SSP reconstruction error using the LDs trained for each sparsity  $T$  is less than EOF reconstruction, for either leading order coefficients or best coefficient combination, for all values of  $T$  shown. The best combination of EOF coefficients, chosen approximately using OMP, achieves less error than the LS solution to the leading order EOFs, with added cost of search.

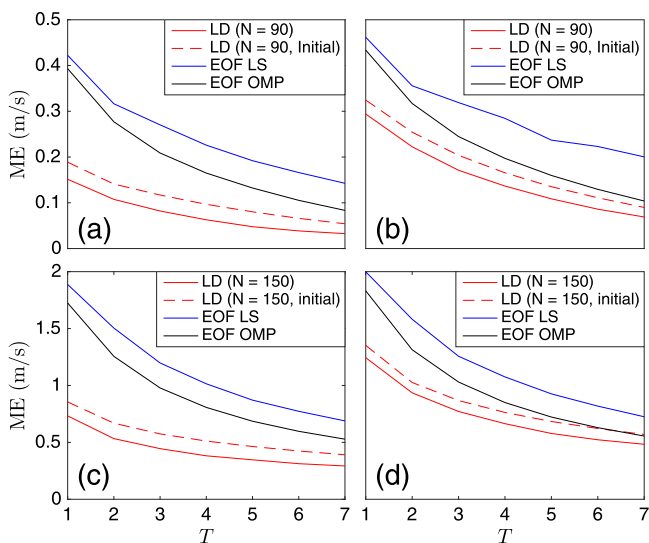


FIG. 8. (Color online) Mean reconstruction error (ME) versus  $T$  using EOFs (solved using LS and OMP) and LDs ( $N = 90$  for HF-97 and  $N = 150$  for SCS) for (a) HF-97 and (c) SCS. Mean reconstruction error  $ME_{CV}$  for out-of-sample data calculated with K-fold cross validation for  $J = 10$  folds, (b) HF-97 and (d) SCS.

Just one LD entry achieves the same ME as more than 6 leading order EOF coefficients, or greater than 4 EOF coefficients chosen by search [Figs. 8(a) and 8(c)]. To illustrate the representational power of the LD entries, both true and reconstructed SSPs are shown in Fig. 9(a) for the HF-97 data and in Fig. 9(b) for the SCS data. Nine true SSP examples from each training set, for HF-97 (SCS) taken at 100 (80) point intervals from  $m = 100 - 900$  (80 - 720), are reconstructed using one LD coefficient. It is shown for each case, that nearly all of the SSP variability is captured using a single LD coefficient.

### C. Cross-validation of SSP reconstruction

The out of sample SSP reconstruction performance of LDs and EOFs is tested using K-fold cross-validation.<sup>34</sup> The entire SSP data set  $\mathbf{Y}$  of  $M$  profiles, for each experiment, is divided into  $J$  subsets with equal numbers of profiles  $\mathbf{Y} = [\mathbf{Y}_1, \dots, \mathbf{Y}_J]$ , where the fold  $\mathbf{Y}_j \in \mathbb{R}^{K \times (M/J)}$ . For each of the  $J$  folds: (1)  $\mathbf{Y}_j$  is the set of out of sample test cases, and the training set  $\mathbf{Y}_{tr}$  is

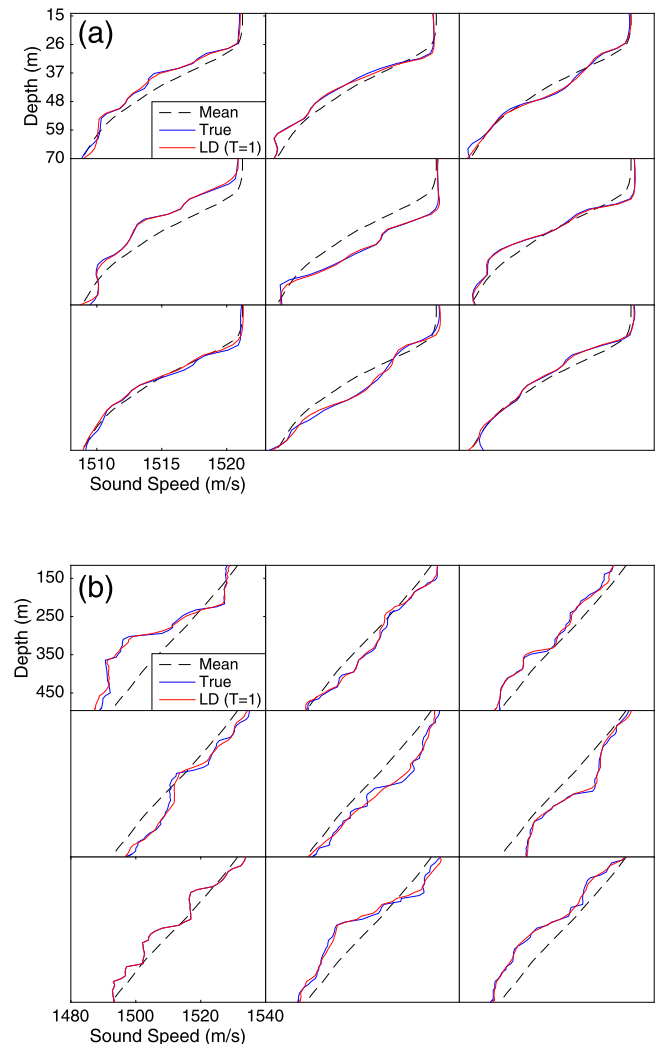


FIG. 9. (Color online) True SSP reconstruction of nine example profiles using one coefficient ( $T = 1$ ) from LD for (a) HF-97 ( $N = 90$ ) and (b) SCS ( $N = 150$ ).



$$\mathbf{Y}_{tr} = \{\mathbf{Y}_i | \forall i \neq j\}; \quad (20)$$

(2) the LD  $\mathbf{Q}_j$  and EOFs are derived using  $\mathbf{Y}_{tr}$ ; and (3) coefficients estimating test samples  $\mathbf{Y}_j$  are solved for  $\mathbf{Q}_j$  with sparse processor Eq. (6), and for EOFs by solving for leading order terms and by solving with sparse processor. The out of sample error from cross validation  $\text{ME}_{\text{CV}}$  for each method is then

$$\text{ME}_{\text{CV}} = \frac{1}{KM} \sum_{j=1}^J \|\mathbf{Y}_j - \hat{\mathbf{Y}}_j\|_1. \quad (21)$$

The out of sample reconstruction error  $\text{ME}_{\text{CV}}$  increases over the within-training-set estimates for both the learned and EOF dictionaries, as shown in Figs. 8(b) and 8(d) for  $J=10$  folds. The mean reconstruction error using the LDs, as in the within-training-set estimates, is less than the EOF dictionaries. For both the HF-97 (SCS) data, more than two (2) EOF coefficients, choosing best combination by search, or more than three (equal to 3) leading-order EOF coefficients solved with LS, are required to achieve the same out of sample performance as one LD entry.

#### D. Solution space for SSP inversion

Acoustic inversion for ocean SSP is a non-linear problem. One approach is coefficient search using genetic algorithms.<sup>1</sup> Discretizing each coefficient into  $H$  values, the number of candidate solutions for  $T$  fixed coefficients indices is

$$S_{\text{fixed}} = H^T. \quad (22)$$

If the coefficient indices for the solution can vary, as per dictionary learning with LD  $\mathbf{Q} \in \mathbb{R}^{K \times N}$ , the number of candidate solutions  $S_{\text{comb}}$  is

$$S_{\text{comb}} = H^T \frac{N!}{T!(N-T)!}. \quad (23)$$

Using a typical  $H=100$  point discretization of the coefficients, the number of possible solutions for fixed and combinatorial dictionary indices are plotted in Fig. 10. Assuming an unknown SSP similar to the training set, the SSP may be constructed up to acceptable resolution using one coefficient from the LD ( $10^4$  possible solutions, see Fig. 10). To achieve the similar ME, seven EOFs coefficients are required ( $10^{14}$  possible solutions, Fig. 10) using fixed indices and the best EOF combination requires five EOFs ( $10^{17}$  possible solutions, Fig. 10).

#### V. CONCLUSION

Given sufficient training data, dictionary learning generates optimal dictionaries for sparse reconstruction of a given signal class. Since these LDs are not constrained to be orthogonal, the entries fit the distribution of the data such that signal example is approximated using few LD entries. Relative to EOFs, each LD entry is informative to the signal variability.

The K-SVD dictionary learning algorithm is applied to ocean SSP data from the HF-97 and SCS experiments. It is

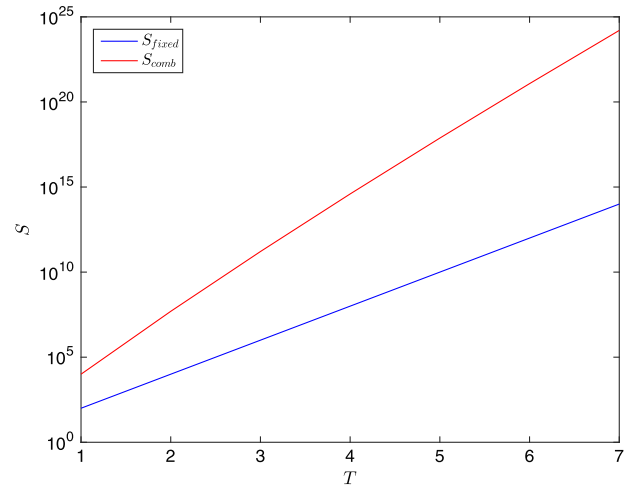


FIG. 10. (Color online) Number of candidate solutions  $S$  for SSP inversion versus  $T$ ,  $S_{\text{fixed}}$  using fixed indices and  $S_{\text{comb}}$  best combination of coefficients. Each coefficient is discretized with  $H=100$  for dictionary  $\mathbf{Q} \in \mathbb{R}^{K \times N}$  with  $N=100$ .

shown that the LDs generated describe ocean SSP variability with high resolution using fewer coefficients than EOFs. As few as one coefficient from a LD describes nearly all the variability in each of the observed ocean SSPs. This performance gain is achieved by the larger number of informative elements in the LDs over EOF dictionaries. Provided sufficient SSP training data are available, LDs can improve SSP inversion resolution with negligible computational expense. This could provide improvements to geoacoustic inversion,<sup>1</sup> matched field processing,<sup>36,37</sup> and underwater communications.<sup>31</sup>

#### ACKNOWLEDGMENTS

The authors would like to thank Dr. Robert Pintel for the use of the South China Sea CTD data. This work is supported by the Office of Naval Research, Grant No. N00014-11-1-0439.

<sup>1</sup>P. Gerstoft, "Inversion of seismoacoustic data using genetic algorithms and *a posteriori* probability distributions," *J. Acoust. Soc. Am.* **95**(2), 770–782 (1994).

<sup>2</sup>L. R. LeBlanc and F. H. Middleton, "An underwater acoustic sound velocity data model," *J. Acoust. Soc. Am.* **67**(6), 2055–2062 (1980).

<sup>3</sup>M. I. Taroudakis and J. S. Papadakis, "A modal inversion scheme for ocean acoustic tomography," *J. Comp. Acoust.* **1**(4), 395–421 (1993).

<sup>4</sup>P. Gerstoft and D. F. Gingras, "Parameter estimation using multifrequency range-dependent acoustic data in shallow water," *J. Acoust. Soc. Am.* **99**(5), 2839–2850 (1996).

<sup>5</sup>C. Park, W. Seong, P. Gerstoft, and W. S. Hodgkiss, "Geoacoustic inversion using backpropagation," *IEEE J. Ocean. Eng.* **35**(4), 722–731 (2010).

<sup>6</sup>B. A. Tan, P. Gerstoft, C. Yardim, and W. S. Hodgkiss, "Broadband synthetic aperture geoacoustic inversion," *J. Acoust. Soc. Am.* **134**(1), 312–322 (2013).

<sup>7</sup>C. F. Huang, P. Gerstoft, and W. S. Hodgkiss, "Effect of ocean sound speed uncertainty on matched-field geoacoustic inversion," *J. Acoust. Soc. Am.* **123**(6), EL162–EL168 (2008).

<sup>8</sup>R. Rubinstein, A. M. Bruckstein, and M. Elad, "Dictionaries for sparse representation modeling," *Proc. IEEE* **98**(6), 1045–1057 (2010).

<sup>9</sup>M. Elad, *Sparse and Redundant Representations* (Springer, New York, 2010).

<sup>10</sup>I. Tosic and P. Frossard, "Dictionary learning," *IEEE Sig. Proc. Mag.* **28**(2), 27–38 (2011).

- <sup>11</sup>K. Schnass, "On the identifiability of overcomplete dictionaries via the minimisation principle underlying K-SVD," *Appl. Comput. Harmonic Anal.* **37**(3), 464–491 (2014).
- <sup>12</sup>M. Aharon, M. Elad, and A. Bruckstein "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Trans. Signal Process.* **54**(11), 4311–4322 (2006).
- <sup>13</sup>K. Engan, S. O. Aase, and J. H. Husøy, "Multi-frame compression: Theory and design," *Signal Process.* **80**(10), 2121–2140 (2000).
- <sup>14</sup>A. Hyvärinen, J. Hurri, and P. O. Hoyer, *Natural Image Statistics: A Probabilistic Approach to Early Computational Vision* (Springer Science and Business Media, London, 2009).
- <sup>15</sup>C. Christopoulos, A. Skodras, and T. Ebrahimi, "The JPEG2000 still image coding system: An overview," *IEEE Trans. Cons. Elec.* **46**(4), 1103–1127 (2000).
- <sup>16</sup>A. Gersho and R. M. Gray, *Vector Quantization and Signal Compression* (Kluwer Academic, Norwell, MA, 1991).
- <sup>17</sup>H. L. Taylor, S. C. Banks, and J. F. McCoy, "Deconvolution with the  $\ell_1$ -norm," *Geophysics* **44**(1), 39–52 (1979).
- <sup>18</sup>E. Candès, "Compressive sampling," in *Proceedings of the International Congress of Mathematicians* (2006), Vol. 3, pp. 1433–1452.
- <sup>19</sup>G. Edelmann and C. Gaumont, "Beamforming using compressive sensing," *J. Acoust. Soc. Am.* **130**(4), EL232–EL237 (2011).
- <sup>20</sup>A. Xenaki, P. Gerstoft, and K. Mosegaard, "Compressive beamforming," *J. Acoust. Soc. Am.* **136**(1), 260–271 (2014).
- <sup>21</sup>P. Gerstoft, A. Xenaki, and C. F. Mecklenbräuker, "Multiple and single snapshot compressive beamforming," *J. Acoust. Soc. Am.* **138**(4), 2003–2014 (2015).
- <sup>22</sup>Y. Choo and W. Song, "Compressive spherical beamforming for localization of incipient tip vortex cavitation," *J. Acoust. Soc. Am.* **140**(6), 4085–4090 (2016).
- <sup>23</sup>C. Yardim, P. Gerstoft, W. S. Hodgkiss, and J. Traer, "Compressive geoaoustic inversion using ambient noise," *J. Acoust. Soc. Am.* **135**(3), 1245–1255 (2014).
- <sup>24</sup>M. Bianco and P. Gerstoft, "Compressive acoustic sound speed profile estimation," *J. Acoust. Soc. Am.* **139**(3), EL90–EL94 (2016).
- <sup>25</sup>S. Beckouche and J. Ma, "Simultaneous dictionary learning and denoising for seismic data," *Geophysics* **79**(3), A27–A31 (2014).
- <sup>26</sup>M. Taroudakis and C. Smaragdakis, "De-noising procedures for inverting underwater acoustic signals in applications of acoustical oceanography," in *Euronoise 2015 Maastricht* (2015), pp. 1393–1398.
- <sup>27</sup>T. Wang and W. Xu, "Sparsity-based approach for ocean acoustic tomography using learned dictionaries," in *OCEANS 2016 Shanghai IEEE*, pp. 1–6 (2016).
- <sup>28</sup>K. S. Alguri and J. B. Harley, "Consolidating guided wave simulations and experimental data: A dictionary leaning approach," *Proc. SPIE* **9805**, 98050Y (2016).
- <sup>29</sup>A. Hannachi, I. T. Jolliffe, and D. B. Stephenson, "Empirical orthogonal functions and related techniques in atmospheric science: A review," *Int. J. Climatol.* **27**(9), 1119–1152 (2007).
- <sup>30</sup>A. H. Monahan, J. C. Fyfe, M. H. Ambaum, D. B. Stephenson, and G. R. North, "Empirical orthogonal functions: The medium is the message," *J. Clim.* **22**(24), 6501–6514 (2009).
- <sup>31</sup>N. Carbone and W. S. Hodgkiss, "Effects of tidally driven temperature fluctuations on shallow-water acoustic communications at 18kHz," *IEEE J. Ocean. Eng.* **25**(1), 84–94 (2000).
- <sup>32</sup>W. S. Hodgkiss, W. A. Kuperman, and D. E. Ensberg, "Channel impulse response fluctuations at 6 kHz in shallow water," in *Impact of Littoral Environmental Variability of Acoustic Predictions and Sonar Performance* (Springer, Netherlands, 2002), pp. 295–302.
- <sup>33</sup>C. T. Liu, R. Pinkel, M. K. Hsu, J. M. Klymak, H. W. Chen, and C. Villanoy, "Nonlinear internal waves from the Luzon Strait," *Eos Trans. AGU* **87**(42), 449–451 (2006).
- <sup>34</sup>T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference and Prediction*, 2nd ed. (Springer, New York, 2009).
- <sup>35</sup>Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad, "Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition," in *IEEE Proc. 27th Annu. Asilomar Conf. Signals, Systems and Computers* (1993), pp. 40–44.
- <sup>36</sup>A. B. Baggeroer, W. A. Kuperman, and P. N. Mikhalevsky, "An overview of matched field methods in ocean acoustics," *IEEE J. Ocean. Eng.* **18**(4), 401–424 (1993).
- <sup>37</sup>C. M. Verlinden, J. Sarkar, W. S. Hodgkiss, W. A. Kuperman, and K. G. Sabra, "Passive acoustic source localization using sources of opportunity," *J. Acoust. Soc. Am.* **138**(1), EL54–EL59 (2015).