

LEARNING DICTIONARY FOR EFFICIENT SIGNAL COMPRESSION

Afshin Abdi, Ali Payani, Faramarz Fekri

Georgia Institute of Technology

ABSTRACT

We consider the problem of learning dictionaries for data compression. Different from ordinary learning methods, the objective is to design a dictionary such that the signal has a low entropy representation in the basis of the dictionary, rather than giving a sparse or low-energy representation. To achieve this goal, we need to consider the effect of quantization on the rate-distortion curve as well as an estimation of the distributions of the coefficients. Based on this probability estimation, the coefficients are computed, quantized and then entropy-coded. As such, we have developed algorithms for different classes of dictionaries; orthonormal, union of orthonormals and general dictionaries with unit-norm atoms, to iteratively learn the dictionary and the distribution models of the coefficients. A mixture of Gaussians is adopted to estimate the probability and is updated using the expectation maximization algorithm together with the dictionary learning. Simulation results on the real seismic data show the effectiveness of the proposed algorithm compared to ordinary dictionary learning methods.

Index Terms— Dictionary Learning, Signal Compression, Low-Entropy Signal Representation

1. INTRODUCTION

Nowadays, modern sensing applications generate excessive amount of data which require a large bandwidth for transmission and a large repository to store. For example, a typical seismic survey may generate tens of terabytes of data. Therefore, compression at the sensor before transmission is emerging as a critical part in many data acquisition applications, especially in wireless systems.

Numerous algorithms have been developed for efficient signal compression. A large class of such algorithms relies on finding an appropriate transform to represent data more efficiently, in a way that discarding the 'least important' information in the transform domain has an insignificant effect on the quality of the data. Examples of such transforms include discrete cosine transform (DCT) and wavelets which have been successfully used for image compression in JPEG and JPEG2000 standards. In these methods, the dictionary is fixed and not adapted to a specific desired class of signals. In recent years, many algorithms have been proposed for designing signal-dependent (overcomplete) dictionaries, especially for the sparse representation. These methods have been shown

to produce great results in applications such as denoising [1,2], face detection [3], image classification and restoration [4,5] and recently for image and video compression [6–8].

In this paper, the goal is to design dictionaries that are adapted to the available data and result in good compression performance, especially for the class of seismic signals. Our approach is different in the sense that the main objective is designing a dictionary and devising an algorithm to find a representation of the signal in the dictionary domain such that the representation requires less bits for compression. This is a departure from conventional dictionary learning methods, whose objective is to design dictionaries for sparse signal representation instead of directly optimizing for compression.

In the next section, we formulate the problem more rigorously and state the assumptions that help in designing a dictionary-based compression algorithm. In section III, the algorithms for different classes of dictionaries are presented and finally, through simulations, we demonstrate the effectiveness of the proposed approach.

1.1. Notations

We denote vectors by bold-faced lower-case letters and matrices by bold-faced upper-case letters. For a vector \mathbf{a} , its i -th entry is denoted by a_i and for a matrix \mathbf{A} , $\mathbf{a}_{:,i}$ or \mathbf{a}_i represents the i -th column, $\mathbf{a}_{i,:}$ the i -th row of the matrix and $a_{i,j}$ the element at row i and column j . $\|\mathbf{A}\|_F$ is the Frobenius norm of matrix \mathbf{A} and $\|\mathbf{a}\|_2$ is the ℓ_2 -norm of vector \mathbf{a} . \log_2 denotes the logarithm in base 2 and \ln is the natural logarithm. $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes a multivariate Gaussian random variable with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.

2. PROBLEM STATEMENT

In this paper, we investigate the problem of dictionary domain compression of data $\mathbf{X} = [\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_N]$. It involves designing an appropriate dictionary \mathbf{D} , from a set of desired dictionaries $\mathcal{D} \subset \mathbb{R}^{n \times m}$, and an algorithm to find discrete coefficients $\mathbf{w}^{(q)}$ from a (possibly unknown) set $\mathcal{W} \subset \mathbb{R}^m$, such that $\mathbf{x} \approx \mathbf{D}\mathbf{w}^{(q)}$ and the rate to transmit $\mathbf{w}^{(q)}$ is as low as possible. Minimizing the total rate, R , of transmitting $\mathbf{W} = [\mathbf{w}_1 \mathbf{w}_2, \dots, \mathbf{w}_N]$ such that the error remains less than a given threshold, ε , can be expressed as

$$\begin{aligned} & \min_{\mathbf{D} \in \mathcal{D}, \mathbf{W} \in \mathcal{W}} R(\mathbf{W}) \\ & \text{s.t. } \|\mathbf{X} - \mathbf{D}\mathbf{W}\|_F^2 \leq \varepsilon. \end{aligned} \quad (\text{P1})$$

Alternatively, for a suitable parameter λ , we desire to

$$\min_{\mathbf{D} \in \mathcal{D}, \mathbf{W} \in \mathcal{W}} \|\mathbf{X} - \mathbf{D}\mathbf{W}\|_F^2 + \lambda R(\mathbf{W}). \quad (\text{P2})$$

Note that we have ignored the cost of transmitting dictionary, \mathbf{D} , as it is constant and independent of the learning algorithm. In the limit, by increasing the number of data samples, the fixed cost of transmitting \mathbf{D} becomes negligible compared to the rate of all coefficients. For example, in our simulations, it contributes to less than 10% of the total bit-rate. Or, as is the case in the memory-assisted compression [9], the dictionary is learned at both the encoder and decoder, via training the dictionary over available data in the common memory from past transmissions.

Two possible approaches to optimize the rate are

1. Directly optimizing for discrete values of $\mathbf{w}^{(q)}$. This requires optimizing for a discrete (finite or infinite) set \mathcal{W} and a probability distribution $P(\cdot)$ on each element of \mathcal{W} . However, the increase in the number of parameters makes the optimization problem difficult to solve in most cases and may result in poor local minimum.
2. An alternative method would be to let the coefficients take any value, but use an (optimum) quantizer to discretize and then encode for compression.

In this paper, we focus on the second approach. The coefficients $\mathbf{W} \in \mathbb{R}^{m \times N}$ are found such that with an appropriate quantizer $Q(\cdot)$, $\mathbf{W}_q = Q(\mathbf{W})$, $\|\mathbf{X} - \mathbf{D}\mathbf{W}_q\|_F^2 \leq \varepsilon$ and the rate $R(\mathbf{W}_q)$ is minimized.

2.1. Quantizer Design

Although it is possible to use Lloyd-Max [10, 11] approach to design a quantizer with minimum distortion, for any number of quantization bins, it does not necessarily give the minimum rate for the required distortion. On the other hand, entropy-constrained algorithms to design the quantizer are more complex to optimize and in our experiments with real data on seismic signals and at our desired SNRs (20-40 dB), they did not perform much better. Therefore, we decided to use uniform quantizers to discretize the coefficients. In addition to its simplicity, uniform quantizer has the advantages of (i)

1. The scalar uniform quantizer is asymptotically efficient under weak assumptions on the density function, i.e., its output entropy is asymptotically smaller than the rate of any other quantizer for the given distortion criteria [12].
2. The error due to the quantization is uniformly bounded for all values of \mathbf{w} , i.e., if $\mathbf{w}^{(q)} = Q(\mathbf{w})$, then $\|\mathbf{x} - \mathbf{D}\mathbf{w}^{(q)}\|_2 \leq \|\mathbf{x} - \mathbf{D}\mathbf{w}\|_2 + c$ where c is a constant that only depends on the quantization step-size δ and dictionary \mathbf{D} , independent of \mathbf{w} and \mathbf{x} . This helps controlling the SNR of each individual block of data, \mathbf{x} .

3. In quantizing a continuous random variable z with small quantization step-size, δ , the rate is $R(z^{(q)}) \approx H_d(z) - \log_2 \delta$, where $H_d(\cdot)$ is the differential entropy. Therefore, optimizing for coefficients in $\mathbb{R}^{m \times N}$ and using a uniform quantizer adds only a constant to the cost function, independent of the final solution.

2.2. Modeling the Distribution of Coefficients

The next step is estimating the probability distribution of the coefficients to design a (near optimum) codebook and compute the rate. Although using non-parametric density estimation methods gives more flexibility, it requires more information to be transmitted to have the probability model at both encoder and decoder. Furthermore, optimizing for the distribution to reduce the rate is not an easy task and it has the potential of over-fitting the distribution. Hence, for the new data, the model may not be a good representative. Therefore, parametric density estimation is the method of choice for rate optimization.

In this paper, we used Gaussian Mixture Model (GMM) to estimate the distribution of coefficients, i.e., $p(\mathbf{w}) = \sum_{\mathbf{s}} \pi(\mathbf{s})p(\mathbf{w}|\mathbf{s})$ where $\pi(\mathbf{s})$ is the weight of sources indexed by \mathbf{s} and $p(\cdot|\mathbf{s}) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{s}}, \boldsymbol{\Sigma}_{\mathbf{s}})$ is a multivariate Gaussian distribution of the same dimension as \mathbf{w} . We choose GMM mainly because the proposed algorithm will have a closed form solution and becomes less complex. Furthermore, for large enough number of sources in the mixture, GMM results in the same distribution for the *quantized* coefficients, $\mathbf{w}^{(q)}$.

To further reduce the complexity of the compression algorithm, we encode each coefficient separately. This implies that the underlying distributions for w_i 's, used in the compression, be independent, i.e., $p(\mathbf{w}) = \prod_i p_i(w_i)$. This is equivalent to assuming diagonal covariance matrices in the mixture model.

2.3. Outline of the Dictionary Learning Algorithm

As simultaneously optimizing for all parameters is not an easy task, we divide the optimization problem (P1) or (P2) into three steps such that at each step, only one set of parameters are updated, while the remainings are kept fixed.

As the coefficients are encoded separately, $R(\mathbf{W}) = \sum_{i=1}^N R(\mathbf{w}_i) = -\sum_{i=1}^N \log_2 p(\mathbf{w}_i)$. The rate can be upper bounded as $-\log p(\mathbf{w}_i) \leq -\log \pi(\hat{\mathbf{s}}_i) - \log p(\mathbf{w}_i|\hat{\mathbf{s}}_i)$, where $\hat{\mathbf{s}}_i$ is the maximum a posteriori (MAP) estimation of source index, \mathbf{s} , from the coefficients \mathbf{w}_i . This can be interpreted as the encoder at the transmitter estimates the source that have generated the coefficients and use the specific codebook designed for that source to compress the coefficients. This helps to further simplify the optimization problem. Thus, the general steps of iterative dictionary learning for compression can be summarized as following;

1. For the current coefficients, update dictionary \mathbf{D} to minimize the error

$$\min_{\mathbf{D} \in \mathcal{D}} \|\mathbf{X} - \mathbf{D}\mathbf{W}\|_F^2. \quad (1)$$

- Fixing the distribution model $\Theta = \{(\pi(\mathbf{s}), \boldsymbol{\mu}_s, \boldsymbol{\Sigma}_s)\}$, and dictionary \mathbf{D} , the sub-problem of finding the best compressible coefficient for each data vector \mathbf{x} is

$$\min_{\mathbf{s}, \mathbf{w}} \|\mathbf{x} - \mathbf{D}\mathbf{w}\|_2^2 - \lambda (\log \pi(\mathbf{s}) + \log p(\mathbf{w}|\mathbf{s})), \quad (2)$$

or

$$\begin{aligned} \min_{\mathbf{s}} \max_{\mathbf{w}} \log \pi(\mathbf{s}) + \log p(\mathbf{w}|\mathbf{s}) \\ \text{s.t. } \|\mathbf{x} - \mathbf{D}\mathbf{w}\|_2^2 \leq \varepsilon. \end{aligned} \quad (3)$$

- Update the distribution to fit the coefficients and further reduce the bit rate.
- Iterate the above steps until convergence.

In the next section, we consider different classes of dictionaries and address each of the above subproblems.

3. LEARNING DICTIONARY FOR COMPRESSION

In this section, we consider three different classes of dictionaries: orthonormal, union of orthonormal dictionaries, and general dictionaries with unit-norm atoms.

3.1. Updating Dictionary

As mentioned earlier, at each iteration of the algorithm, the dictionary is updated to minimize the error $\|\mathbf{X} - \mathbf{D}\mathbf{W}\|_F^2$ for a fixed \mathbf{W} . Here, we briefly review the adopted methods to update the dictionary in each desired class of dictionaries.

3.1.1. Orthonormal dictionary

It is well known that the orthonormal dictionary \mathbf{D} that minimizes $\|\mathbf{X} - \mathbf{D}\mathbf{W}\|_F^2$ for fixed \mathbf{X} and \mathbf{W} is given by the singular value decomposition:

Lemma 1 (see e.g. [13]). *Let USV^T be the svd of $\mathbf{X}\mathbf{W}^T$. The orthonormal dictionary that minimizes $\|\mathbf{X} - \mathbf{D}\mathbf{W}\|_F^2$ is given by $\mathbf{D} = UV^T$.*

3.1.2. Union of Orthonormal Dictionaries

Let $\mathbf{D} = [\mathbf{D}_1 \mathbf{D}_2 \dots \mathbf{D}_L]$ be the concatenation of L orthonormal dictionaries \mathbf{D}_l . By decomposing the coefficients \mathbf{W} accordingly, the error can be rewritten as $\mathbf{X} - \mathbf{D}\mathbf{W} = \mathbf{X} - \sum_{l=1}^L \mathbf{D}_l \mathbf{W}_l$. In this case, block coordinate relaxation methods [13, 14] have been used successfully to update the sub-dictionaries, one at a time. Let $\mathbf{E}_l = \mathbf{X} - \sum_{k \neq l} \mathbf{D}_k \mathbf{W}_k$. Fixing all sub-dictionaries except \mathbf{D}_l , Lemma 1 can be used to find the optimum dictionary that minimizes the error, i.e., $\mathbf{D}_l = \arg\min_{\mathbf{D}} \|\mathbf{E}_l - \mathbf{D}\mathbf{W}_l\|_F^2$.

3.1.3. Dictionaries with Unit-Norm Atoms

Assume that $\mathcal{D} = \{\mathbf{D} \in \mathbb{R}^{n \times m} : \|\mathbf{d}_i\|_2 = 1, i = 1, 2, \dots, m\}$. Unlike the orthonormal case, there is no closed-form solution to update \mathbf{D} . However, different algorithms such as MOD and projection [15], K-SVD [16] and Lagrange dual [17] have been proposed to find the dictionary that minimizes the error.

As the complexity of the above methods is relatively high, we propose using the following simple step-by-step algorithm to update the dictionary atoms iteratively. Note that $\mathbf{D}\mathbf{W} = \sum_j \mathbf{d}_j \mathbf{w}_{j,:}$. Therefore, fixing all atoms of the dictionary except \mathbf{d}_i , the error will be $\mathbf{E}_i = \mathbf{X} - \sum_{j \neq i} \mathbf{d}_j \mathbf{w}_{j,:}$ and the optimum choice for \mathbf{d}_i for the given coefficients is

$$\begin{aligned} \mathbf{d}_i^* = \arg\min_{\mathbf{d}} \|\mathbf{E}_i - \mathbf{d}\mathbf{w}_{i,:}\|_F^2 \\ \text{s.t. } \|\mathbf{d}\|_2 = 1, \end{aligned}$$

whose solution is $\mathbf{d}_i^* = \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$, where $\mathbf{v} = \mathbf{E}_i \mathbf{w}_{i,:}^T$.

Note that in the rate-constrained representation of the signal, (P1) or (P2), since the coefficients no longer minimize only the distortion term $\|\mathbf{X} - \mathbf{D}\mathbf{W}\|_F^2$, they will not be orthogonal to the error, i.e. $\mathbf{E}_i \mathbf{w}_{i,:}^T \neq \mathbf{0}$ and \mathbf{d}_i^* is usually well-defined.

3.2. Computing Coefficients

Recall that for a given source index s in the mixture model, the coefficients are found using (see (2))

$$\begin{aligned} \mathbf{w}^*(s) = \arg\min_{\mathbf{w}} \|\mathbf{x} - \mathbf{D}\mathbf{w}\|_2^2 - \lambda (\log \pi(\mathbf{s}) + \log p(\mathbf{w}|\mathbf{s})) \\ = \boldsymbol{\mu}_s + (\mathbf{D}^T \mathbf{D} + \frac{\lambda}{2} \boldsymbol{\Sigma}_s^{-1})^{-1} \mathbf{D}^T (\mathbf{x} - \mathbf{D}\boldsymbol{\mu}_s), \end{aligned} \quad (4)$$

where $\boldsymbol{\mu}_s$ and $\boldsymbol{\Sigma}_s$ are the mean and covariance matrix of the multivariate Gaussian source determined by index s .

The resulting minimum total cost is obtained as

$$\begin{aligned} \mathcal{J}(s, \mathbf{w}^*(s)) = \lambda \left[(\mathbf{x} - \mathbf{D}\boldsymbol{\mu}_s)^T (\lambda \mathbf{I} + 2\mathbf{D}\boldsymbol{\Sigma}_s \mathbf{D}^T)^{-1} \right. \\ \left. (\mathbf{x} - \mathbf{D}\boldsymbol{\mu}_s) + 0.5 \ln \det(\boldsymbol{\Sigma}_s) - \ln \pi(s) \right]. \end{aligned} \quad (5)$$

For a given dictionary \mathbf{D} , finding the best sources s^* for (P2), which minimizes (5), requires enumerating all possible choices of s and comparing the values of the total cost function. Alternatively, for (P1), it requires an additional optimization step to find λ such that $\|\mathbf{x} - \mathbf{D}\mathbf{w}^*(s)\|_2^2 \leq \varepsilon$ prior to comparing the rates for all possible choices of s .

However, since the covariance matrix is assumed to be diagonal¹, for the orthonormal dictionary, the source index and optimum value of the i -th coefficient are obtained by

$$s_i^* = \arg\min_s \frac{(\mu_{i,s} - y_i)^2}{\lambda + 2\sigma_{i,s}^2} + \ln \frac{\sigma_{i,s}}{\pi_i(s)}, \quad (6a)$$

$$w_i^* = w_i(s_i^*) = \frac{\lambda \mu_{i,s_i^*} + 2\sigma_{i,s_i^*}^2 y_i}{\lambda + 2\sigma_{i,s_i^*}^2}, \quad (6b)$$

where $\mathbf{y} = \mathbf{D}^T \mathbf{x}$, and $\mu_{i,s}$ and $\sigma_{i,s}$ are the mean and standard deviation of the s -th source in the GMM of the i -th coefficient.

For union of L orthonormal dictionaries, we can use block coordinate relaxation method to compute the coefficients efficiently. Let $\mathbf{D} = [\mathbf{D}_1 \mathbf{D}_2 \dots \mathbf{D}_L]$, and \mathbf{w}_l be the vector

¹Since each coefficient is compressed separately, they should have independent probability models.

of coefficients corresponding to D_l . Starting from an initial guess for w (using matching pursuit or pseudo-inverse of D), at each iteration, we fix all coefficients except w_l and compute $y_l = D_l^T \left(x - \sum_{k \neq l} D_k w_k \right)$. Then, s_l and w_l are updated by minimizing $\|y - D_l w_l\|_2^2 + \lambda R(w_l)$, given by (6).

3.3. Updating the Model

To update the models, the Expectation Maximization (EM) algorithm is exploited [18, 19]. At each iteration of the dictionary learning algorithm, few iterations of the EM algorithm is applied on each row of the coefficients to update the model.

4. SIMULATION RESULTS

We have evaluated the performance of our proposed method on real seismic traces available from [20] and [21], (referred to as DB1 and DB2 here). The data from each sensor was treated as a one-dimensional signal and was segmented into blocks of length n to create data vectors, x_i s. In the simulations $n = 16$ and $n = 32$ were used. To model the distributions of the coefficients, different number of sources, K , in the Gaussian mixture model were tested and based on the results, we found out that there is no significant improvement in the performance by increase K beyond 5. So, here, only the results for $K = 5$ mixtures for each coefficient are presented.

The rate-distortion (R-D) curves of the proposed method is compared against standard DCT (for orthonormal dictionaries) and [DCT I] (for over-complete dictionaries) and standard dictionary learning algorithms for sparse representation (LASSO-based optimization with different weights). To obtain the R-D curves, we ran all algorithms for different values of the parameters and used appropriate uniform quantizer to get the best compression rate for the SNR from 20 to 40 dB.

Figures 1 and 2 show the results for orthonormal dictionaries of size 32×32 on DB1 and DB2, respectively. As DB1 has more measurement noise, the performance gap (in Fig. 1) is not as significant as those of DB2 (in Fig. 2). Furthermore, since the existing sparse dictionary learning method is not specifically designed with the compression gain as the objective, it has a poor performance and sometimes, in our simulations we observed that it performs even worse than standard DCT-based methods. Similar results were obtained for dictionary size 16×16 which is omitted here.

In Fig. 3, the performance of the proposed algorithm for learning union of two orthonormal dictionaries ($n = 16$ and $m = 32$) is compared versus the [DCT I] dictionary and sparse dictionary learning method, where the coefficients of [DCT I] dictionary is computed via orthogonal matching pursuit [22].

5. CONCLUSION

In this paper, we have developed a dictionary learning algorithm such that the desired class of signals has a low-entropy representation in the basis of the dictionary and hence, is more suitable for compression. To do so, we first studied the effect

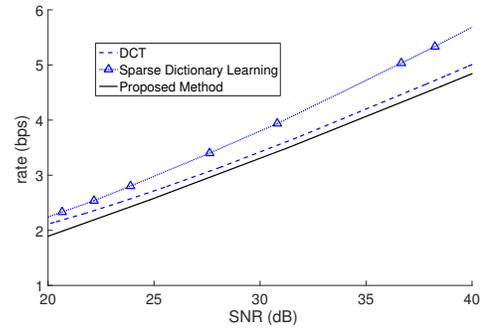


Fig. 1. R-D of dictionaries for database DB1, $n = 32$, $k = 5$

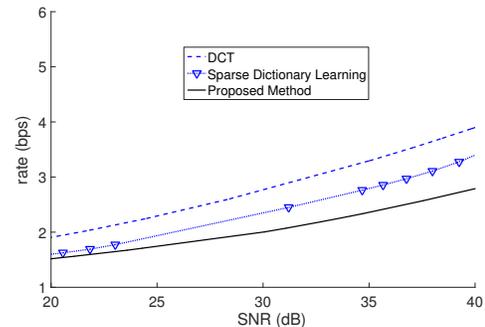


Fig. 2. R-D of dictionaries for database DB2, $n = 32$, $k = 5$

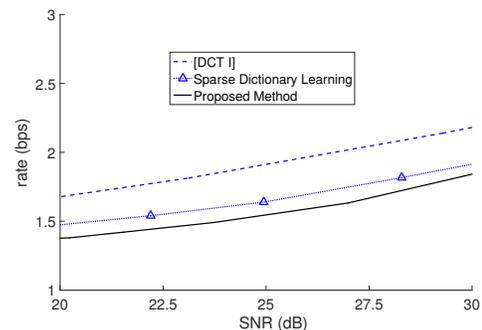


Fig. 3. R-D (per coefficient) for different dictionaries of size 16×32 and $K = 5$, applied on DB2

of quantizer on the rate-distortion and concluded that for the real signals (seismic traces) and the desired range of SNR (20dB and above), the uniform quantizer performs close to optimum. To further simplify the compression algorithm, the encoding of each coefficient is done separately and in a memoryless manner. We used different Gaussian mixture models to approximate the distributions of the coefficients. Based on these assumptions, the dictionary learning algorithm consists of iteratively updating the probability models, computing compressible coefficients and updating the dictionary to reduce the distortion and rate. Simulation results showed that the proposed algorithm performs well for the compression of real seismic traces, outperforming ordinary learning methods for the sparsity and standard DCT transform.

6. REFERENCES

- [1] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Transactions on Image Processing*, vol. 15, no. 12, pp. 3736–3745, Dec 2006.
- [2] Y. M. Huang, L. Moisan, M. K. Ng, and T. Zeng, "Multiplicative noise removal via a learned dictionary," *IEEE Transactions on Image Processing*, vol. 21, no. 11, pp. 4534–4543, Nov 2012.
- [3] Q. Zhang and B. Li, "Discriminative k-svd for dictionary learning in face recognition," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, June 2010, pp. 2691–2698.
- [4] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman, "Non-local sparse models for image restoration," in *2009 IEEE 12th International Conference on Computer Vision*, Sept 2009, pp. 2272–2279.
- [5] I. Ramirez, P. Sprechmann, and G. Sapiro, "Classification and clustering via dictionary learning with structured incoherence and shared features," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, June 2010, pp. 3501–3508.
- [6] K. Skretting and K. Engan, "Image compression using learned dictionaries by RLS-DLA and compared with K-SVD," in *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2011, pp. 1517–1520.
- [7] K. Kreutz-Delgado, J. Murray, B. Rao, K. Engan, T. Lee, and T. Sejnowski, "Dictionary learning algorithms for sparse representation," *Neural Computation*, vol. 15, no. 2, pp. 349–396, Feb 2003.
- [8] J. F. Murray and K. Kreutz-Delgado, "Sparse image coding using learned overcomplete dictionaries," in *Machine Learning for Signal Processing, 2004. Proceedings of the 2004 14th IEEE Signal Processing Society Workshop*, Sept 2004, pp. 579–588.
- [9] A. Beirami and F. Fekri, "Memory-assisted universal source coding," in *Data Compression Conference (DCC), 2012*, April 2012, pp. 392–392.
- [10] J. Max, "Quantizing for minimum distortion," *IRE Transactions on Information Theory*, vol. 6, no. 1, pp. 7–12, March 1960.
- [11] S. Lloyd, "Least squares quantization in PCM," *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 129–137, Mar 1982.
- [12] H. Gish and J. Pierce, "Asymptotically efficient quantizing," *IEEE Transactions on Information Theory*, vol. 14, no. 5, pp. 676–683, Sep 1968.
- [13] S. Lesage, R. Gribonval, F. Bimbot, and L. Benaroya, "Learning unions of orthonormal bases with thresholded singular value decomposition," in *Proceedings. (ICASSP '05). IEEE International Conference on Acoustics, Speech, and Signal Processing, 2005.*, vol. 5, March 2005, pp. v/293–v/296 Vol. 5.
- [14] S. Sardy, A. G. Bruce, and P. Tseng, "Block coordinate relaxation methods for nonparametric wavelet denoising," *Journal of computational and graphical statistics*, vol. 9, no. 2, pp. 361–379, 2000.
- [15] K. Engan, S. O. Aase, and J. H. Husoy, "Method of optimal directions for frame design," in *Acoustics, Speech, and Signal Processing, 1999. Proceedings., 1999 IEEE International Conference on*, vol. 5, 1999, pp. 2443–2446 vol.5.
- [16] M. Aharon, M. Elad, and A. Bruckstein, "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Transactions on Signal Processing*, vol. 54, no. 11, pp. 4311–4322, Nov 2006.
- [17] H. Lee, A. Battle, R. Raina, and A. Y. Ng, "Efficient sparse coding algorithms," in *Advances in Neural Information Processing Systems 19*, B. Schölkopf, J. C. Platt, and T. Hoffman, Eds. MIT Press, 2007, pp. 801–808.
- [18] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the royal statistical society. Series B (methodological)*, pp. 1–38, 1977.
- [19] J. Bilmes, "A gentle tutorial of the EM algorithm and its application to parameter estimation for gaussian mixture and hidden markov models," International Computer Science Institute, Tech. Rep. TR-97-021, 1998.
- [20] "UTAM seismic data library," <http://utam.gg.utah.edu/SeismicData/SeismicData.html>, [Online].
- [21] "USGS seismic data library," <http://energy.usgs.gov/GeochemistryGeophysics/>, [Online].
- [22] S. G. Mallat and Z. Zhang, "Matching pursuits with time-frequency dictionaries," *IEEE Transactions on Signal Processing*, vol. 41, no. 12, pp. 3397–3415, Dec 1993.