

Segmentation of Multivariate Mixed Data via Lossy Data Coding and Compression

Yi Ma, Harm Derksen, Wei Hong, and John Wright

*Coordinated Science Laboratory
1308 West Main Street, Urbana, IL 61801
University of Illinois at Urbana-Champaign*

REPORT DOCUMENTATION PAGE			Form Approved OMB NO. 0704-0188	
Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comment regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services.				
1. AGENCY USE ONLY (Leave blank)	2. REPORT DATE		3. REPORT TYPE AND DATES COVERED	
4. TITLE AND SUBTITLE Segmentation of Multivariate Mixed Data via Lossy Data Coding and Compression			5. FUNDING NUMBERS NSF IIS 03-47456 Career	
6. AUTHOR(S) Yi Ma, Harm Derksen, Wei Hong, and John Wright				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Coordinated Science Laboratory University of Illinois at Urbana-Champaign 1308 West Main Street Urbana, Illinois 61801-2307			8. PERFORMING ORGANIZATION REPORT NUMBER UILU-ENG-06-2216 DC-224	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) NSF Career IIs-0347456 NSF CRS-EHS-0509151 NSF CCF-TF-0514944 ONR YIP N00014-05-1-0644			10. SPONSORING/MONITORING AGENCY REPORT NUMBER	
11. SUPPLEMENTARY NOTES The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official position, policy, or decision, unless so designated by other documentation				
12a. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited.			12b. DISTRIBUTION CODE	
13. ABSTRACT (Maximum 200 words) In this paper, based on ideas from lossy data coding and compression, we present a simple but effective technique for segmenting multivariate mixed data that are drawn from a mixture of Gaussian distributions, which are allowed to be almost degenerate. The goal is to find the optimal segmentation that minimizes the overall coding length of the segmented data, subject to a given distortion. By analyzing the coding length/rate of mixed data, we formally establish some strong connections of data segmentation to many fundamental concepts in lossy data compression, rate distortion theory, and multiple-channel communications. We show that a deterministic segmentation is the (asymptotically) optimal solution for compressing mixed data. We propose a very simple and effective algorithm to find the optimal segmentation, which does not require any prior knowledge of the number or dimension of the groups, nor does it involve any parameter estimation. Simulation results reveal intriguing phase-transition behaviors of the number of segments when changing the level of distortion or the amount of outliers. Finally, we demonstrate how this technique can be readily applied to segment real imagery and bioinformatic data.				
14. SUBJECT TERMS multivariate mixed data, data segmentation, rate distortion, lossy data coding, data compression, image segmentation, microarray data clustering			15. NUMBER OF PAGES 35	
			16. PRICE CODE	
17. SECURITY CLASSIFICATION OF REPORT	18. SECURITY CLASSIFICATION OF THIS PAGE	19. SECURITY CLASSIFICATION OF ABSTRACT		20. LIMITATION OF ABSTRACT

Segmentation of Multivariate Mixed Data via Lossy Data Coding and Compression

Yi Ma, *Member, IEEE*, Harm Derksen, Wei Hong, John Wright

Abstract

In this paper, based on ideas from lossy data coding and compression, we present a simple but effective technique for segmenting multivariate mixed data that are drawn from a mixture of Gaussian distributions, which are allowed to be almost degenerate. The goal is to find the optimal segmentation that minimizes the overall coding length of the segmented data, subject to a given distortion. By analyzing the coding length/rate of mixed data, we formally establish some strong connections of data segmentation to many fundamental concepts in lossy data compression, rate distortion theory, and multiple-channel communications. We show that a deterministic segmentation is the (asymptotically) optimal solution for compressing mixed data. We propose a very simple and effective algorithm to find the optimal segmentation, which does not require any prior knowledge of the number or dimension of the groups, nor does it involve any parameter estimation. Simulation results reveal intriguing phase-transition behaviors of the number of segments when changing the level of distortion or the amount of outliers. Finally, we demonstrate how this technique can be readily applied to segment real imagery and bioinformatic data.

Index Terms

multivariate mixed data, data segmentation, rate distortion, lossy data coding, data compression, image segmentation, microarray data clustering.

Submitted to IEEE Transactions on Pattern Analysis and Machine Intelligence, April 2006

Yi Ma, Wei Hong, and John Wright are with the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, 1308 W Main St., Urbana, Illinois, 61801 {yima, weihong, jnwright}@uiuc.edu. This work was partially supported by NSF CAREER IIS-0347456, NSF CRS-EHS-0509151, NSF CCF-TF-0514955, ONR YIP N00014-05-1-0633.

Harm Derksen, Department of Mathematics, University of Michigan, 530 Church St., Ann Arbor, Michigan 48109, hderksen@umich.edu. This work was partially supported by NSF CAREER DMS-0349019.

I. INTRODUCTION

Data that arise from practical problems in such diverse fields as image/signal processing, pattern recognition, computer vision, and bioinformatics, are often characterized by complicated multi-modal, multivariate distributions. Segmentation (or clustering) is widely recognized as an important step in representing, analyzing, interpreting or compressing such mixed data.

Now the intriguing questions are: What does “segmentation” really mean and how to define it mathematically? What should the proper criterion for segmentation be and what do the segmentation results depend on? How should we measure the “gain” or “loss” of the segmentation? Last but not the least, why is segmentation the right thing to do? Answers to these questions to some extent have been complicated by the many approaches and solutions proposed in the literature for segmenting or modeling various types of mixed data (see [1], [2] and references therein for a review).

A somewhat traditional way of defining segmentation is to first choose a simple class of models which each subset is supposed to fit. Some of the popular models are either probabilistic distributions (e.g., Gaussian distributions) or geometric/algebraic sets (e.g., linear subspaces). Then the whole mixed data are assumed to be samples drawn from a mixture of such probabilistic distributions [3], [4] or geometric/algebraic sets [5]. The typical approach to segmenting the data then entails estimating the mixture of all the models and simultaneously or subsequently decomposing them into individual ones. In this way, data segmentation is essentially identified with a (mixture) model estimation problem. Segmenting the data and estimating the model are therefore strongly coupled together. Various approaches have been proposed to resolve the coupling in the literature:

- Iterate between the data segmentation and model estimation. Representative methods include the K-means algorithm [6]–[9] (or its variants [10]–[12]) and the Expectation Maximization (EM) algorithm [13], [14], which is essentially a greedy descent algorithm to find the maximum-likelihood (ML) estimate of a mixture of probabilistic (Gaussian) distributions [3], [4], [15].
- Resolve the coupling between data segmentation and model estimation by first estimating a mixture model that does not depend on the segmentation of the data and then decompose the mixture into individual components. Representative methods include Generalized Principal

Component Analysis (GPCA), in which the mixture model is assumed to be an arrangement of subspaces [5].

A common assumption behind all these approaches is that a good estimate of the underlying mixture model(s) is necessary for the segmentation of the data. In a sense, the goodness of the segmentation relies on how good the estimate is. For instance, the given data $W = (w_1, w_2, \dots, w_m)$ are commonly assumed to be drawn from a mixture of distributions: $p(x|\theta, \pi) \doteq \sum_{j=1}^k \pi_j p_j(x|\theta_j)$. When trying to obtain the optimal estimate of the mixture model, one usually chooses any of the model estimation criteria, e.g., the maximum likelihood (ML) estimate:

$$(\hat{\theta}, \hat{\pi})_{ML} = \arg \max_{\theta, \pi} \sum_{i=1}^m \log p(w_i|\theta, \pi), \quad (1)$$

where θ is the parameter of certain class of (mixture) distributions of interest. The EM algorithm [13] (or its variants [16]) is often used to optimize the likelihood function of such a mixture model. The ML criterion is equivalent to minimizing the negated log-likelihood: $\sum_i -\log p(w_i|\theta, \pi)$, which is approximately the expected coding length, $\text{Length}(W|\hat{\theta}, \hat{\pi})$, required to store the data using the optimal coding scheme for the distribution $p(x|\hat{\theta}, \hat{\pi})$ [17].

When the number of component models, k , is not given *a priori*, we must estimate it from the data, a difficult task that is further complicated when the data are corrupted by a significant amount of outliers. To some extent, almost all model selection criteria used to determine the number of component models are equivalent to minimizing the coding length needed to describe both the data and the model, i.e., the minimum description length (MDL) criterion [4], [18]–[20]:

$$(\hat{\theta}, \hat{\pi})_{MDL} = \arg \min_{\theta, \pi} L(W, \theta, \pi) = L(W|\theta, \pi) + L(\theta, \pi), \quad (2)$$

where the parameters θ, π are assumed to have certain distribution $p(\theta, \pi)$. In general, the length function $L(\cdot)$ is chosen according to the optimal Shannon coding [17]: $-\log p(W|\theta, \pi)$ for W and $-\log p(\theta, \pi)$ for θ, π . Incidentally, this objective function coincides with the maximum-likelihood estimate and hence the EM algorithm again becomes the method of choice [4].

However, ML and MDL only truly correspond to minimum coding lengths when the random variables to be encoded are discrete¹. For (multivariate) real-valued data, a finite coding length can only be obtained if we encode the data and model parameters *subject to a certain*

¹or for continuous random variables, in the limit as the quantization error goes to zero.

distortion $\varepsilon > 0$. To this end, [21] has studied the properties of *lossy maximum likelihood* (LML) and *lossy minimum description length* (LMDL) criteria. There, it is shown that (to first order, asymptotically) minimizing the coding rate of the data subject to the distortion ε :

$$(\hat{\theta}, \hat{\pi})_{LML} = \arg \min_{\theta, \pi} R(\hat{p}(W), \theta, \pi, \varepsilon), \quad (3)$$

$$(\hat{\theta}, \hat{\pi})_{LMDL} = \arg \min_{\theta, \pi} R(\hat{p}(W), \theta, \pi, \varepsilon) + L(\theta, \pi), \quad (4)$$

where $\hat{p}(W)$ is the empirical estimate of the probabilistic distribution from the data W , is equivalent to computing the LML or LMDL estimate, with desirable properties such as (strong) consistency as an estimator. In our context, the coding rate (subject to a distortion) provides a natural measure of the goodness of segmentation for real-valued mixed data. In fact, the goal of modeling and segmentation of mixed data should indeed be consistent with that of data coding/compression: If the data can be fit with better models after segmentation, the data should be represented or encoded more efficiently with respect to such models.

A. Contributions of This Paper

In this paper, we do not consider modeling and segmenting data that have arbitrary mixture distributions. We are only interested in data that consist of multiple Gaussian-like groups, which may have significantly different and anisotropic covariances. The covariances of the groups may be even nearly degenerate, in which case we essentially want to fit the data with *multiple subspaces, possibly of different dimensions*. In this context, vector quantization (VQ) can be viewed as the special case of fitting the data with zero-dimensional (affine) subspaces [10].

Our approach to segmenting such mixed data follows the spirit of (lossy) ML and MDL. Our goal is to *find the optimal segmentation of the mixed data which results in the shortest coding length subject to a given distortion of the data*. Our method however offers the following improvements over existing methods:

- 1) All of the estimates discussed above (ML,MDL,LML,LMDL) are optimal only in an *asymptotical* sense, i.e., for an infinite sequence of i.i.d. samples from the class of distributions of interest. In practice, however, we are often dealing with a finite (and often small) set of samples. Thus, we introduce a measure of coding length for each group that not only closely approximates the optimal rate-distortion function for a Gaussian source [17] but also gives a tight upper bound for any finite number of samples.

- 2) We will prove that with this choice of coding length/rate, the (asymptotically) optimal segmentation is *deterministic* – no probabilistic (or fuzzy) segmentation can further reduce the overall coding length. This provides a theoretical justification that (deterministic) segmentation is not only useful for pragmatic purposes, but also the optimal solution for compressing data that are mixture of Gaussians or subspaces.
- 3) An explicit formula for the coding length/rate function² allows one to directly evaluate goodness of the segmentation. The tightness of the formula for small data sets leads to an efficient³ “bottom-up” algorithm that minimizes the overall coding length by repeatedly merging small subsets, starting from individual data points. As we will show with extensive simulations and experiments, this approach resolves the difficult model selection issue [4] in an effective way, especially when *the number of groups is unknown or there is a significant amount of outliers*.
- 4) When the level of distortion (or the density of outliers) varies continuously, the number of groups typically exhibits a *phase transition* behavior similar to that in statistical physics, with the “correct” segmentation corresponding to one of the stable phases. Our simulations show that the number of segments need not be a monotonic function of the distortion.⁴

B. Organization of This Paper

We provide a summary of the basic ideas and the resulting algorithm of our approach in Section II. In Section III, based on ideas from the rate distortion theory in information theory, we introduce a formula for the coding rate/length needed to encode a set of vectors subject to a given distortion. An alternative verification of the formula is given in Appendix I, and Appendix II shows how the formula should be modified when the data is nonzero mean. In Section IV, we study properties of the overall coding rate/length of mixed data after being segmented into multiple groups. Extensive simulation and experimental results of the proposed algorithm on synthetic and real data are given in Section V.

²This is the case for Gaussian sources. In general computing the rate-distortion function for an arbitrary distribution is a difficult problem although many numerical methods exist in the literature (see [22] and references therein).

³The complexity of the proposed algorithm is polynomial in both the size and dimension of the data.

⁴A different phase transition has been noticed in vector quantization using deterministic annealing, where the number of clusters increases monotonically when the annealing temperature decreases [10].

II. BASIC IDEAS AND ALGORITHM

In this section, we give a self-contained summary of the main ideas and algorithm of this paper and leave more detailed mathematical analysis and justification to Section III and IV. Readers who are interested only in the algorithm and experiments may bypass the next two sections and skip to Section V without any loss of continuity.

A. Lossy Coding of Multivariate Data

A lossy coding scheme maps a set of vectors $V = (v_1, v_2, \dots, v_m) \in \mathbb{R}^{n \times m}$ to a sequence of binary bits, such that the original vectors can be recovered up to an allowable distortion $\mathbb{E}[\|v_i - \hat{v}_i\|^2] \leq \varepsilon^2$. The length of the encoded sequence is denoted as the function $L(V) : \mathbb{R}^{n \times m} \rightarrow \mathbb{Z}_+$.

In general, the coding scheme and the associated $L(\cdot)$ function can be chosen to be optimal for any family of distributions of interest. In the case where the data are i.i.d. samples from a zero-mean⁵ multivariate Gaussian distribution $\mathcal{N}(0, \Sigma)$, the function $R = \frac{1}{2} \log_2 \det(I + \frac{n}{\varepsilon^2} \Sigma)$ provides a good approximation to the optimal rate-distortion function [17].⁶ As $\hat{\Sigma} = \frac{1}{m} VV^T$ is an estimate of the covariance Σ , the average number of bits needed per vector is:

$$R(V) \doteq \frac{1}{2} \log_2 \det \left(I + \frac{n}{\varepsilon^2 m} VV^T \right). \quad (5)$$

For readers who are less familiar with the rate-distortion theory, we will give an intuitive explanation of this formula in Section III.

Representing the m vectors of V therefore requires $mR(V)$ bits. Since the optimal codebook is adaptive to the data V , we must also represent it with an additional $nR(V)$ bits⁷, yielding an

⁵For simplicity, in the main text, we will derive and present our main results with the zero-mean assumption. However, all the formulae, results, and algorithms can be readily extended to the nonzero mean case, as shown in Appendix II.

⁶Strictly speaking, the rate-distortion function for the Gaussian source $\mathcal{N}(0, \Sigma)$ is $R = \frac{1}{2} \log_2 \det \left(\frac{n}{\varepsilon^2} \Sigma \right)$ when $\frac{\varepsilon^2}{n}$ is smaller than the smallest eigenvalue of Σ . Thus the approximation is good only when the distortion ε is relatively small. However, when $\frac{\varepsilon^2}{n}$ is larger than some eigenvalues of Σ , the rate distortion function becomes more complicated [17]. Nevertheless, the approximate formula $R = \frac{1}{2} \log_2 \det \left(I + \frac{n}{\varepsilon^2} \Sigma \right)$ can be viewed as the rate-distortion of the ‘‘regularized’’ source that works for all range of ε . Furthermore, as we will show in Appendix I, the same formula gives a tight upper bound of the coding rate for any finite number of samples.

⁷This can be viewed as the cost of coding the n principal axes of the data covariance $\frac{1}{m} VV^T$. A more detailed explanation of $L(V)$ is given in Section III.

overall coding length of

$$L(V) \doteq (m+n)R(V) = \frac{m+n}{2} \log_2 \det \left(I + \frac{n}{\varepsilon^2 m} VV^T \right). \quad (6)$$

We will study the properties of this function in Section III. For purposes of segmentation, it suffices to note that in addition to being (approximately) asymptotically optimal for Gaussian data, $L(V)$ also provides a tight bound on the number of bits needed to code a finite number of vectors (that span a subspace), regardless of the underlying probability distribution (see Appendix I for a proof).

B. Segmentation via Data Compression

Given a set of samples, $W = (w_1, w_2, \dots, w_m) \in \mathbb{R}^{n \times m}$, one can always view them as drawn from a single Gaussian source and code W subject to distortion ε^2 using $L(W)$ bits. However, if the samples are drawn from a mixture of Gaussian distributions or subspaces, it may be more efficient to code W as the union of multiple (disjoint) groups: $W = W_1 \cup W_2 \cup \dots \cup W_k$. If each group is coded separately, the total number of bits needed is

$$L^s(W_1, W_2, \dots, W_k) \doteq \sum_{i=1}^k L(W_i) + |W_i| \left(-\log_2(|W_i|/m) \right), \quad (7)$$

where $|W_i|$ indicates the cardinality (i.e. number of vectors) of the group W_i . In the above expression, the term $\sum_{i=1}^k |W_i| \left(-\log_2(|W_i|/m) \right)$ is the number of bits needed to code (losslessly) the membership of the m samples in the k groups (e.g. using the Huffman coding [17]).⁸

Then, given a fixed coding scheme with its associated coding length function $L(\cdot)$, an optimal segmentation is one which minimizes the segmented coding length, $L^s(\cdot)$, over all possible partitions of W . Moreover, we will see that due to the properties of the rate-distortion function (5) for Gaussian data, softening the objective function (7) by allowing probabilistic (or fuzzy) segmentation does *not* further reduce the (expected) overall coding length (see Theorem 3 of Section IV).

Notice that the above objective (7) is a function of the distortion ε . In principle, one may add a ‘‘penalty’’ term,⁹ so as to determine the optimal distortion ε^* . The resulting objective will then

⁸Here we assume that the ordering of the samples is random and entropy coding is the best we can do to code the membership. However, if the samples are ordered such that nearby samples more likely belong to the same group (e.g., in segmenting pixels of an image), the second term can and should be replaced by a tighter estimate.

⁹For instance, we may add the term $mn \log \varepsilon$ to the overall coding length L^s .

corresponds to an optimal coding length that only depends on the data. However, we here leave ε as a free parameter to be set by the user. In practice, this allows the user to obtain potentially hierarchical segmentation of the data at different scales of quantization. We will thoroughly examine how the value of ε affects the final segmentation through experiments in Section V.

C. Minimizing the Coding Length

Finding the global minimum of the overall coding length L^s over all partitions of the dataset is a daunting combinatorial optimization problem, intractable for large data sets. Nevertheless, the coding length can be effectively minimized in a steepest descent fashion, as outlined in Algorithm 1. The minimization proceeds in a “bottom-up” fashion: initially, every sample is treated as its own group. At each iteration, two groups S_1 and S_2 are chosen so that merging them results in the greatest decrease in the coding length. The algorithm terminates when the coding length cannot be further reduced by merging any pair of groups.¹⁰ A simple implementation which maintains a table containing $L^s(S_i \cup S_j)$ for all i, j requires $O(m^3 + m^2n^3)$ time, where m is the number of samples and n the dimension of the space.

Algorithm 1 (Pairwise Steepest Descent of Coding Length).

- 1: **input:** the data $W = (w_1, w_2, \dots, w_m) \in \mathbb{R}^{n \times m}$ and a distortion $\varepsilon^2 > 0$.
 - 2: initialize $\mathcal{S} := \{S = \{w\} \mid w \in W\}$.
 - 3: **while** $|\mathcal{S}| > 1$ **do**
 - 4: choose distinct sets $S_1, S_2 \in \mathcal{S}$ such that $L^s(S_1 \cup S_2) - L^s(S_1, S_2)$ is minimal.
 - 5: **if** $L^s(S_1 \cup S_2) - L^s(S_1, S_2) \geq 0$ **then break;**
 - 6: **else** $\mathcal{S} := (\mathcal{S} \setminus \{S_1, S_2\}) \cup \{S_1 \cup S_2\}$.
 - 7: **end**
 - 8: **output:** \mathcal{S}
-

Extensive simulations and experiments demonstrate that this algorithm is consistently and remarkably effective in segmenting data that are a mixture of Gaussians or subspaces (see Section V). It tolerates significant amounts of outliers, and requires no prior knowledge of the number

¹⁰In the supplementary material, we have included a video showing the convergence of this algorithm on data drawn from mixtures of subspaces in \mathbb{R}^3 .

of groups. As a greedy descent scheme, the algorithm does not guarantee to always find the globally optimal segmentation for any given (W, ε) .¹¹ From our experience, we found that the main factor affecting the global convergence of the algorithm seems to be the density of the samples relative to the distortion ε^2 . In Section V we will give strong empirical evidences for the convergence of the algorithm over a wide range of ε .

III. LOSSY CODING OF MULTIVARIATE DATA

In this section, we give a more detailed justification and analysis of the coding rate/length functions introduced in the previous section. In the next section, we provide a more thorough justification of the compression-based approach to data segmentation. Readers who are less concerned with such technical details may skip these two sections at first read, without much loss of continuity.

If the given data $w_i \in \mathbb{R}^n$ are i.i.d. samples of a random vector w with the probabilistic distribution $p(w)$, the optimal coding scheme and the optimal coding rate of such a random vector w have been well studied in *information theory* (see [17] and references therein). However, here we are dealing with a finite set of vectors $W = (w_1, w_2, \dots, w_m)$. Such a data set can be viewed as a non-parametric distribution itself – each vector w_i in W occurs with an equal probability $1/m$. The optimal coding scheme for the distribution $p(w)$ is no longer optimal for W and the formula for the coding length no longer accurate. Nevertheless, some of the basic ideas of deriving the optimal coding rate can still be extended to the non-parametric setting. In this section, borrowing ideas from information theory, we derive a tight bound of the coding length or rate for the given data W . In Appendix I, we give an alternative derivation of the bound. Although both approaches essentially arrive at the same estimate, they together reveal that the derived coding length/rate function holds under different conditions:

- 1) The derivation in this section shows that for small ε the formula for $R(W)$ gives a good approximation to the (asymptotically) optimal rate-distortion function of a Gaussian source.
- 2) The derivation in Appendix I shows that the same coding length/rate formula works for any finite set of vectors W that span a subspace.

¹¹However, it may be possible to improve the convergence by using more complicated split-and-merge strategies [16]. In addition, due to Theorem 1 of Section IV, the globally (asymptotically) optimal segmentation can also be computed via concave optimization [23], at the cost of potentially exponential computation time.

A. The Rate Distortion Function

For simplicity, we here assume that the given data are zero mean, i.e. $\mu \doteq \frac{1}{m} \sum_i w_i = 0$. The reader may refer to Appendix II for the case in which the mean is not zero. Let ε^2 be the squared error allowable for encoding every vector w_i . That is, if \hat{w}_i is an approximation of w_i , we allow $\mathbb{E}[\|w_i - \hat{w}_i\|^2] \leq \varepsilon^2$. In other words, on average, the allowable squared error for each entry of w_i is ε^2/n .

The solution to coding the vectors in W , subject to the mean squared error ε^2 , can be explained by *sphere packing*, which is normally adopted in information theory [17]. Here we are allowed to perturb each vector $w_i \in W$ within a sphere of radius ε in \mathbb{R}^n . In other words, we are allowed to distort each entry of w_i with an (independent) random variable of variance ε^2/n . Without loss of generality, we may model the error as an independent additive Gaussian noise:

$$\hat{w}_i = w_i + z_i, \quad \text{with } z_i \sim \mathcal{N}\left(0, \frac{\varepsilon^2}{n} I\right). \quad (8)$$

Then the covariance matrix of the vectors \hat{w}_i is:

$$\hat{\Sigma} \doteq \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m \hat{w}_i \hat{w}_i^T\right] = \frac{\varepsilon^2}{n} I + \frac{1}{m} W W^T \in \mathbb{R}^{n \times n}. \quad (9)$$

The volume of the region spanned by these vectors is proportional to (the square root of the determinant of the covariance matrix): $\text{vol}(\hat{W}) \propto \sqrt{\det\left(\frac{\varepsilon^2}{n} I + \frac{1}{m} W W^T\right)}$. Similarly, the volume spanned by each random vector z_i is proportional to $\text{vol}(z) \propto \sqrt{\det\left(\frac{\varepsilon^2}{n} I\right)}$.

In order to encode each vector, we can partition the region spanned by all the vectors into non-overlapping spheres of radius ε . When the volume of the region $\text{vol}(\hat{W})$ is significantly larger than the volume of the sphere, the total number of spheres that we can pack into the region is approximately equal to

$$\#\text{of spheres} = \text{vol}(\hat{W})/\text{vol}(z). \quad (10)$$

Thus, to know each vector w_i with an accuracy up to ε^2 , we only need to specify which sphere w_i is in (see Figure 1). If we use binary numbers to label all the spheres in the region of interest, the number of bits needed is

$$R(W) \doteq \log_2(\#\text{of spheres}) = \log_2\left(\text{vol}(\hat{W})/\text{vol}(z)\right) = \frac{1}{2} \log_2 \det\left(I + \frac{n}{m\varepsilon^2} W W^T\right), \quad (11)$$

where the last equality uses the fact $\det(A)/\det(B) = \det(B^{-1}A)$.

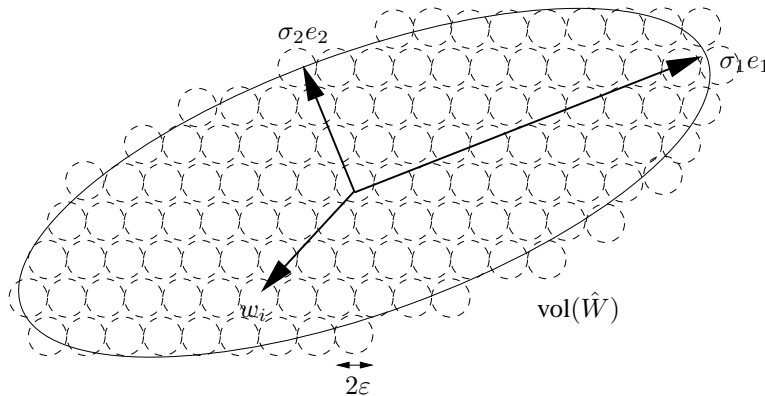


Fig. 1. Coding of a set of vectors in a region in \mathbb{R}^n with an accuracy up to ε^2 . To know the vector w_i , we only need to know the label of the corresponding sphere. e_1, e_2 represent the singular vectors of the matrix \hat{W} and σ_1, σ_2 the singular values.

If the samples w_i are drawn from a Gaussian source $\mathcal{N}(0, \Sigma)$, then $\frac{1}{m}WW^T$ converges to the covariance Σ of the Gaussian source. Thus, we have $R(W) \rightarrow \frac{1}{2} \log_2 \det \left(I + \frac{n}{\varepsilon^2} \Sigma \right)$ as $m \rightarrow \infty$. When $\frac{\varepsilon^2}{n} \leq \lambda_{\min}(\Sigma)$, the optimal rate-distortion for a parallel i.i.d. $\mathcal{N}(0, \Sigma)$ source is $\frac{1}{2} \log_2 \det \left(\frac{n}{\varepsilon^2} \Sigma \right)$, to which (11) provides a good approximation. In general, the optimal rate-distortion is a complicated formula given by reverse-waterfilling on the eigenvalues of Σ (see Theorem 13.3.3 of [17]). The approximation (11) provides an upper bound which holds for all ε , and is tight when ε is small relative to the eigenvalues of the covariance.

The formula for $R(W)$ can also be viewed as the rate-distortion of the source W regularized by a noise of variance $\frac{\varepsilon^2}{n}$ as in equation (8). The covariance $\hat{\Sigma}$ of the perturbed vectors \hat{w}_i always satisfies $\frac{\varepsilon^2}{n} \leq \lambda_{\min}(\hat{\Sigma})$, allowing for a simple, analytic expression for the rate distortion for all range of ε . This regularized rate-distortion has the further advantage of agreeing with the bound for the coding length of finitely many vectors that span a subspace, derived in Appendix I. In addition, this formula resembles the channel capacity of an MIMO Gaussian channel. The interested reader may refer to Appendix III.

Notice that the formula for $R(W)$ is accurate only in the asymptotical sense, i.e., when we are dealing with a large number of samples and the error ε is small (relative to the magnitude of the data W). We want to emphasize that the above derivation of the coding rate does not give an actual coding scheme. The construction of efficient coding schemes which achieve the optimal rate-distortion bound is itself a difficult problem (see, for example, [24] and references

therein). However, for the purpose of measuring the quality of segmentation and compression, all that matters is that *in principle* a scheme attaining the optimal rate $R(W)$ exists.

B. The Coding Length Function

Given the coding rate $R(W)$, the total number of bits needed to encode the m vectors in W is

$$mR(W) = \frac{m}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} WW^T \right). \quad (12)$$

From the communication point of view, $mR(W)$ bits are already sufficient as both the transmitter and the receiver share the same code book – that is they both know the region spanned by W in \mathbb{R}^n . However, from the data representation or compression point of view, we need more bits to represent the code book itself. This is equivalent to specifying all the principal axes of the region spanned by the data, i.e. the singular values/vectors of W , see Figure 1. As the number of principal axes is n , we need $nR(W)$ additional bits to encode them. Therefore, the total number of bits needed to encode the m vectors in $W \subset \mathbb{R}^n$ subject to the squared error ε^2 is¹²

$$L(W) \doteq (m+n)R(W) = \frac{m+n}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} WW^T \right). \quad (13)$$

Appendix I provides an alternative derivation of the same coding length function $L(W)$, as an upper bound for a finite number of samples. If the data W have a non-zero mean, we need more bits to encode the mean too. See in Appendix II how the coding length function should be properly modified in that case.

C. Properties of the Coding Length Function

1) *Commutative Property:* Since $WW^T \in \mathbb{R}^{n \times n}$ and $W^TW \in \mathbb{R}^{m \times m}$ have the same non-zero eigenvalues, the coding length function can also be expressed as:

$$L(W) = \frac{m+n}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} WW^T \right) = \frac{m+n}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} W^TW \right).$$

Thus, if $n \ll m$, the second expression will be less costly for computing the coding length. The matrix W^TW , which depends only on the inner products between pairs of data vectors, is

¹²Compared to the MDL criterion (2), if the term $mR(W)$ corresponds to the coding length for the data, the term $nR(W)$ then corresponds to the coding length for the model parameter θ .

known in the statistical learning literature as the *kernel matrix*. This property suggests that the ideas and the algorithm presented in Section II can be readily extended to segment data sets that have *nonlinear* structures, by choosing a proper kernel function.

2) *Invariant Property*: Notice that in the zero-mean case, the coding length function $L(W)$ is invariant under an orthogonal transformation of the data W . That is, for any orthogonal matrix $U \in O(n)$ or $V \in O(m)$, we have

$$L(UW) = L(W) = L(WV). \quad (14)$$

In other words, the length function depends only on the singular values of W (or eigenvalues of WW^T). This equality suggests that one may choose any orthonormal basis (e.g., Fourier, wavelets) to represent and encode the data and the number of bits needed should always be the same. This agrees with the fact that the chosen coding length (or rate) is optimal for a Gaussian source. However, if the data are non-Gaussian or nonlinear, a proper transformation can still be useful for compressing the data.¹³ In this paper we are essentially seeking a partition, rather than a transformation, of the non-Gaussian (or nonlinear) data set, such that each subset is sufficiently Gaussian (or subspace-like) and hence cannot be compressed any further, either by (orthogonal) transformation or segmentation.

IV. CODING LENGTH OF SEGMENTED DATA

Now suppose we have partitioned the set of m vectors $W = (w_1, w_2, \dots, w_m)$ into k non-overlapping groups $W = W_1 \cup W_2 \cup \dots \cup W_k$. Then the total number of bits needed to encode the segmented data is $L^s(W_1, W_2, \dots, W_k) = \sum_{i=1}^k L(W_i) + |W_i|(-\log_2(|W_i|/m))$. Here the superscript “s” is used to indicate the coding length after segmentation.

A. Segmentation and Compression

To better understand under what conditions a set of data should or should not be segmented so that the overall coding length/rate becomes smaller, we here provide two representative examples. In the examples, we want to study whether a data set should be partitioned into two subsets of

¹³For a more thorough discussion on why some transformations (such as wavelets) are useful for data compression, the reader may refer to [25].

an equal number of vectors: $W_1, W_2 \in \mathbb{R}^{n \times m}$. To simplify the analysis, we assume $m \gg n$ so that we can ignore the asymptotically insignificant terms in the coding length/rate function.

Example 1 (Uncorrelated Subsets): Notice that in general, we have

$$\begin{aligned} L(W_1) + L(W_2) &= \frac{m}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} W_1 W_1^T \right) + \frac{m}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} W_2 W_2^T \right) \\ &\leq \frac{2m}{2} \log_2 \det \left(I + \frac{n}{2m\varepsilon^2} (W_1 W_1^T + W_2 W_2^T) \right) = L(W_1 \cup W_2), \end{aligned}$$

where the inequality is from the concavity of the function $\log_2 \det(\cdot)$ (see Theorem 7.6.7 of [26]). Thus, if the difference $L(W_1 \cup W_2) - (L(W_1) + L(W_2))$ is large, the overhead needed to encode the membership of the segmented data (here one bit per vector) becomes insignificant. If we further assume that W_2 is a rotated version of W_1 , i.e. $W_2 = UW_1$ for some $U \in O(n)$, one can show that the difference $L(W_1 \cup W_2) - (L(W_1) + L(W_2))$ is (approximately) maximized when W_2 becomes orthogonal to W_1 . We call two groups W_1, W_2 *uncorrelated* if $W_1^T W_2 = 0$. Thus, segmenting the data into uncorrelated groups typically reduces the overall coding length. From the viewpoint of sphere packing, Figure 2 explains the reason.

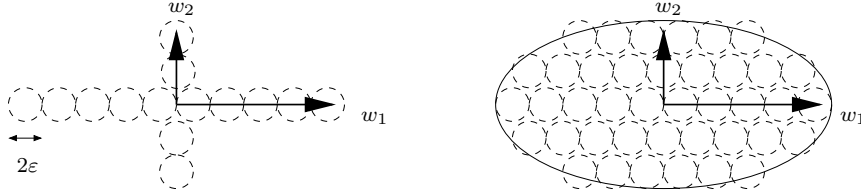


Fig. 2. The number of spheres (code words) of two different schemes for coding two orthogonal vectors. Left: encoding the two vectors separately; Right: encoding the two vectors together.

Example 2 (Strongly Correlated Subsets): We say two groups W_1, W_2 are strongly correlated if they span the same subspace in \mathbb{R}^n . Or somewhat equivalently, we may assume that W_1 and W_2 have approximately the same covariance $W_2 W_2^T \approx W_1 W_1^T$. Thus we have

$$\begin{aligned} L(W_1) + L(W_2) &= \frac{m}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} W_1 W_1^T \right) + \frac{m}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} W_2 W_2^T \right) \\ &\approx \frac{2m}{2} \log_2 \det \left(I + \frac{n}{2m\varepsilon^2} (W_1 W_1^T + W_2 W_2^T) \right) = L(W_1 \cup W_2). \end{aligned}$$

Since $L^s(W_1, W_2) = L(W_1) + L(W_2) + H(|W_1|, |W_2|)$, the overhead needed to encode the membership becomes significant and the segmented data require more bits than the unsegmented.

B. Optimality of Deterministic Segmentation

So far, we have considered only partitioning the data W into k non-overlapping groups. That is, each vector is assigned to a group of probability either 0 or 1. We call such a segmentation “deterministic.” In this section, we examine an important question: *Is there a probabilistic segmentation of the data that can achieve an even lower coding rate?* That is, we consider a more general class of segmentations in which we assign each vector w_i to the group j according to a probability $\pi_{ij} \in [0, 1]$, with $\sum_{j=1}^k \pi_{ij} = 1$ for all $i = 1, 2, \dots, m$.

To facilitate counting the expected coding length of such (probabilistically) segmented data, we introduce a matrix Π_j that collects the membership of the m vectors in group j :

$$\Pi_j \doteq \begin{pmatrix} \pi_{1j} & 0 & \cdots & 0 \\ 0 & \pi_{2j} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \pi_{mj} \end{pmatrix} \in \mathbb{R}^{m \times m}. \quad (15)$$

These matrices satisfy the constraint: $\sum_{j=1}^k \Pi_j = I_{m \times m}$, $\Pi_j \succeq 0$.

Obviously, the j th group has an expected number of $\mathbf{tr}(\Pi_j)$ vectors and the expected covariance is $\frac{1}{\mathbf{tr}(\Pi_j)} W \Pi_j W^T$. If viewed as a Gaussian source, the coding rate of the j th group is bounded by: $R(W_j) \doteq \frac{1}{2} \log_2 \det \left(I + \frac{n}{\mathbf{tr}(\Pi_j) \varepsilon^2} W \Pi_j W^T \right)$. If for each vector w_i , we code it using the coding scheme for the j th group with probability π_{ij} , then the expected total number of bits required to encode the data W according to the segmentation $\Pi = \{\Pi_j\}$ is bounded by¹⁴

$$L^s(W, \Pi) \doteq \sum_{j=1}^k \frac{\mathbf{tr}(\Pi_j) + n}{2} \log_2 \det \left(I + \frac{n}{\mathbf{tr}(\Pi_j) \varepsilon^2} W \Pi_j W^T \right) + \mathbf{tr}(\Pi_j) \left(-\log_2 \frac{\mathbf{tr}(\Pi_j)}{m} \right). \quad (16)$$

Similarly, the expected number of bits needed to encode each vector is bounded by

$$R^s(W, \Pi) \doteq \frac{1}{m} L^s(W, \Pi) = \sum_{j=1}^k \frac{\mathbf{tr}(\Pi_j)}{m} \left(R(W_j) - \log_2 \frac{\mathbf{tr}(\Pi_j)}{m} \right) + \frac{n}{m} R(W_j). \quad (17)$$

Thus, one may consider that the optimal segmentation Π^* is the global minimum of the expected overall coding length $L^s(W, \Pi)$, or equivalently the average coding rate $R^s(W, \Pi)$. To

¹⁴Strictly speaking, the formula is an upper bound for the expected coding length because $L^s(W, \Pi)$ is essentially a concave function of the group assignment Π (see the proof of Theorem 3). Hence, $L^s(W, \mathbb{E}[\Pi]) \geq \mathbb{E}[L^s(W, \Pi)]$ (using that $f(\mathbb{E}[x]) \geq \mathbb{E}[f(x)]$ for concave functions).

some extent, one can view the minimum value of $R^s(W, \Pi)$ as a good approximation to the actual entropy of the given data set W .¹⁵

Notice that the second term in the expression of $R^s(W, \Pi)$, $\frac{n}{m}R(W_j)$, is insignificant when the number of samples is large $m \gg n$. Nevertheless, this term, as well as the term that encodes the membership of the vectors, gives a *tight* bound on the coding length even for small sets of samples. This essentially allows us to find the optimal segmentation in a bottom-up manner by merging small subsets of samples, which is effectively harnessed by the greedy algorithm introduced in Section II. That said, for the rest of this section, we examine more carefully the asymptotic properties of the coding length/rate function.

The first term in the expression of $R^s(W, \Pi)$ is the only part that matters asymptotically (i.e. when the number of vectors in each group goes to infinity) and we denote it as:

$$R^{s,\infty}(W, \Pi) \doteq \sum_{j=1}^k \frac{\mathbf{tr}(\Pi_j)}{2m} \log_2 \det \left(I + \frac{n}{\varepsilon^2 \mathbf{tr}(\Pi_j)} W \Pi_j W^T \right) - \frac{\mathbf{tr}(\Pi_j)}{m} \log_2 \left(\frac{\mathbf{tr}(\Pi_j)}{m} \right).$$

Thus, the global minimum of $R^{s,\infty}(W, \Pi)$ determines the optimal segmentation when the sample size is large.

Theorem 3: The asymptotic part $R^{s,\infty}(W, \Pi)$ of the rate distortion function $R^s(W, \Pi)$ is a concave function of Π in the convex domain $\Omega \doteq \{\Pi : \sum_{j=1}^k \Pi_j = I, \Pi_j \succeq 0\}$.

Proof: Let \mathcal{S} be the set of all $m \times m$ non-negative definite symmetric matrices. We will show that $R^{s,\infty}(W, \Pi)$ is concave as a function from $\mathcal{S}^k \rightarrow \mathbb{R}$, and so is it when restricted to the domain of interest, $\Omega \subset \mathcal{S}^k$.

First consider the second term of $R^{s,\infty}(W, \Pi)$. Notice that $\sum_{j=1}^k \mathbf{tr}(\Pi_j) = m$ is a constant. So we only need to show the concavity of the function $g(P) \doteq -\mathbf{tr}(P) \log_2 \mathbf{tr}(P)$ for $P \in \mathcal{S}$. The function, $f(x) = -x \log_2 x$ is concave, and $g(P) = f(\mathbf{tr}(P))$. So for $\lambda \in [0, 1]$,

$$\begin{aligned} g(\lambda P_1 + (1 - \lambda)P_2) &= f(\lambda \mathbf{tr}(P_1) + (1 - \lambda) \mathbf{tr}(P_2)) \\ &\geq \lambda f(\mathbf{tr}(P_1)) + (1 - \lambda) f(\mathbf{tr}(P_2)) = \lambda g(P_1) + (1 - \lambda) g(P_2). \end{aligned}$$

Thus, $g(P)$ is concave in P .

¹⁵Especially when the data W indeed consist of a mixture of subsets and each group is a typical set of samples from a (almost degenerate) Gaussian distribution.

Now consider the first term of $R^{s,\infty}(W, \Pi)$. Let

$$h(\Pi_j) \doteq \mathbf{tr}(\Pi_j) \log_2 \det \left(I + \frac{n}{\varepsilon^2 \mathbf{tr}(\Pi_j)} W \Pi_j W^T \right).$$

It is well-known in information theory that the function $q(P) \doteq \log_2 \det(P)$ is concave for $P \in \mathcal{S}$ and $P \succ 0$ (see Theorem 7.6.7 of [26]). Now define $r : \mathcal{S} \rightarrow \mathbb{R}$ to be

$$r(\Pi_j) \doteq \log_2 \det(I + \alpha W \Pi_j W^T) = q(I + \alpha W \Pi_j W^T).$$

Since r is just the concave function q composed with an affine transformation $\Pi_j \mapsto I + \alpha W \Pi_j W^T$, r is concave (see Section 3.2.3 of [27]). Let $\psi : \mathcal{S} \times \mathbb{R}_+ \rightarrow \mathbb{R}$ as

$$\psi(\Pi_j, t) \doteq t \cdot \log_2 \det \left(I + \frac{n}{\varepsilon^2 t} W \Pi_j W^T \right) = t \cdot r\left(\frac{1}{t} \Pi_j\right).$$

According to Theorem 3.2.6 of [27], ψ is concave. Notice that $H \doteq \{(\Pi_j, t) : t = \mathbf{tr}(\Pi_j)\}$ is a linear subspace in the product space of \mathbb{R} and the space of all symmetric matrices. So, $H \cap (\mathcal{S} \times \mathbb{R}_+)$ is a convex set, and the desired function, $h(\Pi_j) = \psi(\Pi_j, \mathbf{tr}(\Pi_j))$, is just the restriction of ψ to this convex set. Thus, h is concave.

Since $R^{s,\infty}(W, \Pi)$ is a sum of concave functions in Π_j , it is concave as a function from \mathcal{S}^k to \mathbb{R} , and so is its restriction to the convex set Ω in \mathcal{S}^k . ■

Since $R^{s,\infty}(W, \Pi)$ is concave, its global minimum Π^* is always reached at the boundary, or more precisely, at a vertex of the convex domain Ω , as shown in Figure 3. At the vertex of Ω , the entries π_{ij} of Π^* are either 0s or 1s. It means that even if we allow soft assignment of each point to the k groups according to any probabilistic distribution, the optimal solution with the minimal coding length can always be achieved by assigning each point to one of the groups with probability one! This is the reason why Algorithm 1 does not consider any probabilistic segmentation and is still able to produce (approximately) optimal segmentation.

Another implication of the above theorem is that the problem of minimizing the coding length is essentially a concave optimization problem. Many effective concave optimization algorithms can be adopted to find the globally optimal segmentation, such as the simplex algorithm [23]. However, such generic concave optimization algorithms typically have high (potentially exponential) complexity. In the next section, we will show with extensive simulations and experiments that the greedy algorithm proposed in Section II is already effective in minimizing the coding length.

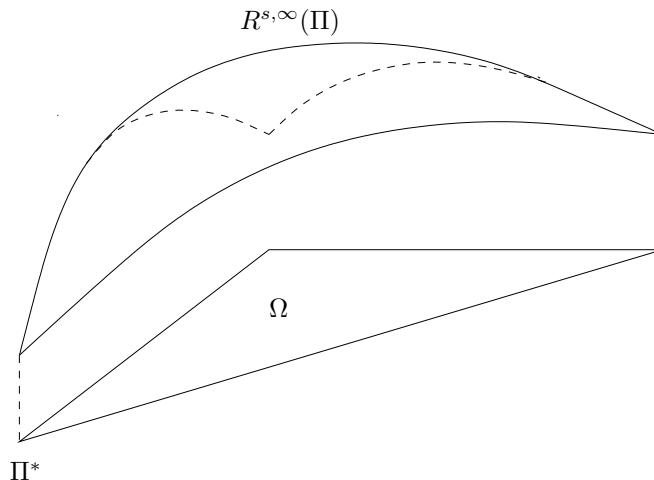


Fig. 3. The function $R^{s,\infty}(W, \Pi)$ is a concave function of Π over a convex domain Ω , which is in fact a polytope in the space \mathbb{R}^{mk} . The minimal coding length is achieved at a vertex Π^* of the polytope.

Interestingly, in multiple-channel communications, the goal is instead to *maximize* the channel capacity, which has very much the same formula as the coding rate function. See Appendix III for more detail. The above theorem suggests that a higher channel capacity may be achieved inside the convex domain Ω , i.e. by probabilistically assigning the transmitters into certain number of groups. As the coding rate function is concave, the maximal channel capacity can be very easily computed via convex optimization [27].

V. SIMULATION AND EXPERIMENTAL RESULTS

In this section, we conduct simulations on a variety of challenging data sets to examine the effectiveness of the proposed coding length function as well the performance of the steepest descent algorithm. In the end, we will also demonstrate some experimental results of applying the algorithm to segment imagery and bioinformatic data.

A. Simulations

1) *Segmentation of Linear Subspaces of Different Dimensions:* We first demonstrate the ability of the algorithm to segment noisy samples drawn from a mixture of linear subspaces of different dimensions. Figure 4 summarizes the configurations tested. For every d -dimensional subspace, $d \times 100$ samples are drawn uniformly from a ball of diameter 1 lying on the subspace. Each

sample is corrupted with independent Gaussian noise of standard deviation $\varepsilon_0 = 0.04$. The segmentation is computed using Algorithm 1, with $\varepsilon = \varepsilon_0$.

Subspace dimensions	Identified dimensions	Classification (%) (Algorithm 1)	Classification (%) (E-M)
$(2, 1, 1)$ in \mathbb{R}^3	2, 1, 1	96.62	39.33
$(2, 2, 1)$ in \mathbb{R}^3	2, 2, 1	90.00	68.98
$(4, 2, 2, 1)$ in \mathbb{R}^5	4, 2, 2, 1	98.53	43.36
$(6, 3, 1)$ in \mathbb{R}^7	6, 3, 1	99.77	66.16
$(7, 5, 2, 1, 1)$ in \mathbb{R}^8	7, 5, 2, 1, 1	98.04	42.29

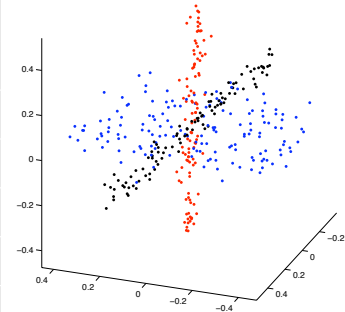


Fig. 4. Left: Simulation results for data drawn from mixtures of noisy linear subspaces. Classification percentages are averaged over 25 trials. Our algorithm correctly identifies the number and dimension of the subspaces in all 25 trials, for all configurations. Far right column: results using Expectation-Maximization with random initialization. Right: the computed segmentation for $(2, 1, 1)$ in \mathbb{R}^3 is displayed.

In each case, the algorithm stops at the correct number of groups, and the dimensions of the segments W_i match those of the generating subspaces.¹⁶ The correctness of the segmentation is further corroborated by the high percentage of points correctly classified (by comparing the segments with the *a priori* groups). For all five configurations, the average percentage of samples assigned to the correct group was at least 90.0%. The main cause of classification error is points which lie near the intersection of multiple subspaces. Due to noise, it may actually be more efficient to code such points according to the optimal coding scheme for one of the other subspaces. We compare our method to Expectation-Maximization, seeded with a random initialization. In all cases, Algorithm 1 dramatically outperforms EM in terms of classification error, despite requiring less prior knowledge.

Since in practice, ε_0 is not known, it is important to investigate the sensitivity of the results to the choice of ε . For each of the examples in Figure 4, Figure 5 gives the range of ε for which Algorithm 1 converges to the a-priori number and dimension of subspaces. Notice that for each of the configurations considered, there exists a significant range of ε for which the greedy algorithm converges.

¹⁶The dimension of each segment W_i is identified using principal component analysis (PCA) by thresholding the singular values of W_i with respect to ε .

Subspace dimensions	(2, 1, 1) in \mathbb{R}^3	(2, 2, 1) in \mathbb{R}^3	(4, 2, 2, 1) in \mathbb{R}^5	(6, 3, 1) in \mathbb{R}^7	(7, 5, 2, 1, 1) in \mathbb{R}^8
$\log_{10} \varepsilon_{max} - \log_{10} \varepsilon_{min}$	2.5	1.75	2.0	2.0	.75

Fig. 5. The size of the range of $\log \varepsilon$ for which the greedy algorithm converges to the correct number and dimension of groups, for each of the arrangements considered in Figure 4.

2) *Global Convergence*: Empirically, we find that Algorithm 1 does not suffer many of the difficulties with local minima that plague other clustering algorithms such as EM. The convergence appears to depend mostly on the density of the samples relative to the distortion ε . For example, if the number of samples is fixed at $m = 1200$, and the data are drawn from three $\lceil \frac{n}{2} \rceil$ -dimensional subspaces in \mathbb{R}^n , the algorithm converges to the correct solution for $n = 2$ upto $n = 56$. Here, we choose $\varepsilon = \varepsilon_0 = 0.008$. Beyond $n = 56$, the algorithm fails to converge to the three *a priori* subspaces as the samples have become too sparse. For $n > 56$, the computed segmentation gives a higher coding length than the *a priori* segmentation.

The same observation occurs for subspaces with different dimensions. For example, we randomly draw 800 noisy ($\varepsilon_0 = 0.14$) samples from four subspaces of dimension 20, 15, 15, 10 in \mathbb{R}^{40} . The results of the greedy algorithm at different distortion ε are shown in Figure 6. As we see from the results, when the distortion ε is very small, the greedy algorithm does not necessarily converge to the optimal coding length. Nevertheless, the number of groups, 4, is still identified correctly by the algorithm when ε becomes relatively large.

3) *Robustness to Outliers*: We test the robustness of Algorithm 1 to outliers on the easily visualized example of two lines and a plane in \mathbb{R}^3 . 158 samples are drawn uniformly from a 2-D disc of diameter 1. 100 samples are drawn uniformly from each of the two line segments of length 1. The additive noise level is $\varepsilon_0 = 0.03$. The data set is contaminated with m_o outliers, whose three coordinates are uniformly distributed on $[-0.5, 0.5]$.

As the number of outliers increases, the segmentation exhibits several distinct phases. For $m_o \leq 300$ (45.6% outliers), the algorithm always finds the correct segmentation. The outliers are merged into a single (three-dimensional) group. From $m_o = 400$ (52.8% outliers) upto $m_o = 1100$ (75.4% outliers), the two lines are correctly identified, but samples on the plane are merged with the outliers. For $m_o = 1200$ (77.4% outliers) and higher, all of the data samples

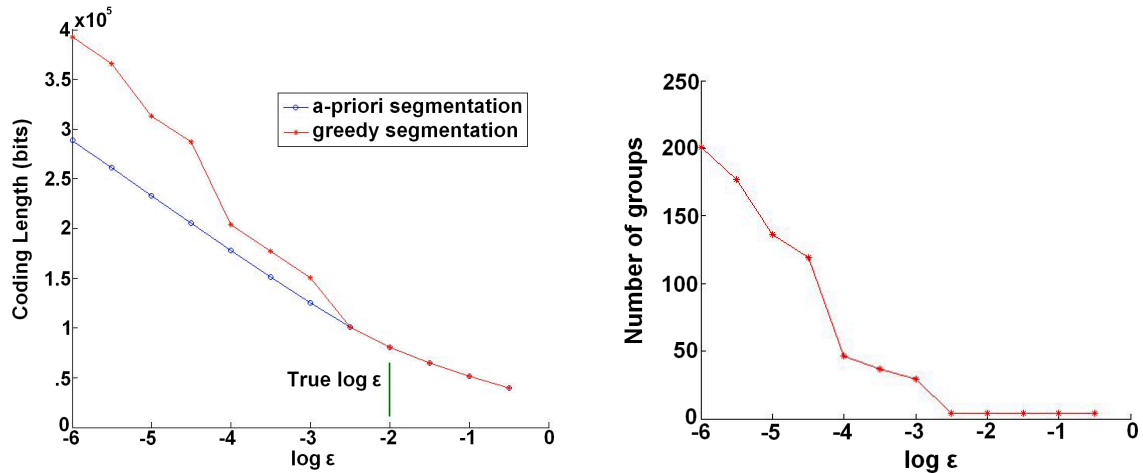


Fig. 6. Left: the coding length found by the greedy algorithm (the red curve) compared to the ground truth (the blue curve) for data drawn from four linear subspaces of dimension 20, 15, 15, 10 in \mathbb{R}^{40} . Right: the number of groups found by the greedy algorithm – it converges to the correct number 4 when the distortion is relatively large.

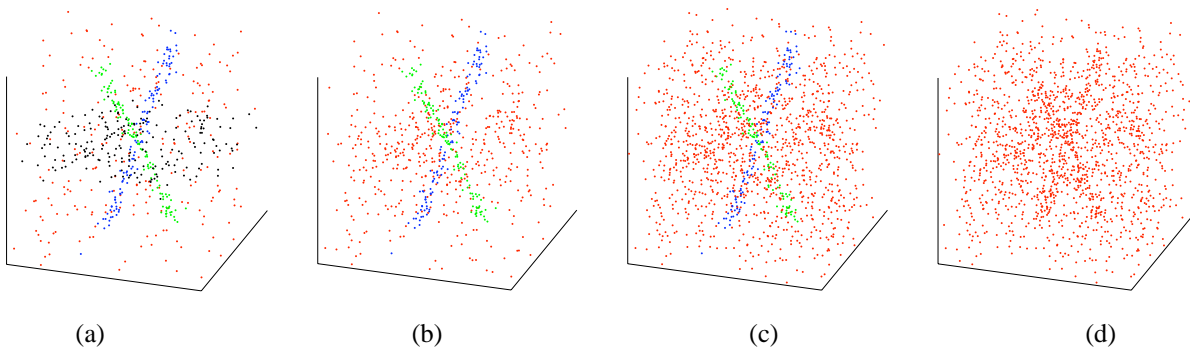


Fig. 7. Segmentation results for data drawn from three linear subspaces, corrupted with various numbers of outliers, m_o . (a) $m_o = 300$ (45.6% outliers). (b) $m_o = 400$ (52.8% outliers). (c) $m_o = 1100$ (75.4% outliers). (d) $m_o = 1200$ (77.0% outliers).

are merged into one group, as the distribution of data has become essentially random in the ambient space. Figure 7 shows the results for $m_o = 300, 400, 1100, 1200$. Notice that the effect of adding the outliers resembles the effect of ice (the lines and the plane) being melted away by warm water. This suggests a similarity between the artificial process of data clustering and the physical process of phase transition.

4) *Number of Segments versus Distortion Level:* Figure 8 shows how the number of segments changes as ϵ varies. $m = 358$ points are drawn from two lines and a plane, as in the previous

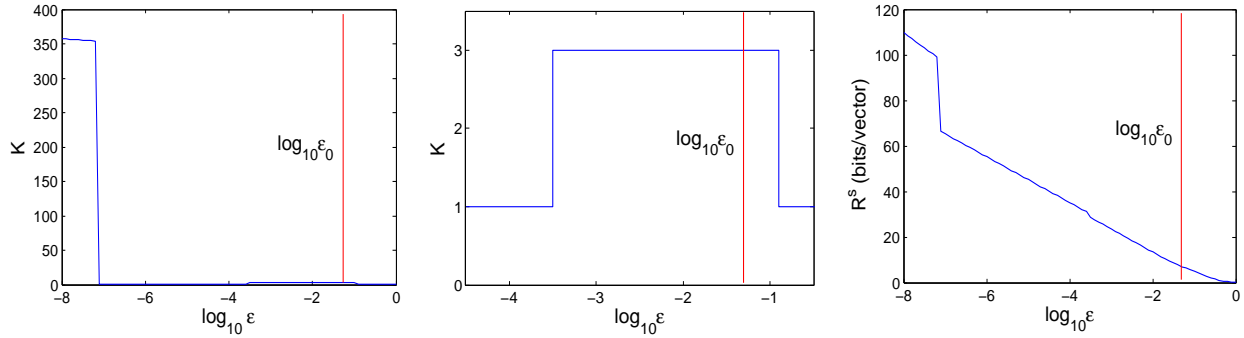


Fig. 8. The effect of varying ε , with $\varepsilon_0 = 0.05$. Left: number of groups, k , versus $\log(\varepsilon)$. Center: detail of k versus $\log(\varepsilon)$ around $\log(\varepsilon_0)$. Right: the coding rate (bits per vector) versus $\log(\varepsilon)$.

experiment, and then perturbed with noise of standard deviation $\varepsilon_0 = 0.05$. Notice that the number of groups experiences distinct phases, with abrupt transitions around several critical values of ε . For sufficiently small ε , each data point is grouped by itself. However, as ε increases, the cost of coding the group membership begins to dominate, and all the points are grouped together in a single three-dimensional subspace (the ambient space). Around the true noise level, ε_0 , there is another stable phase, corresponding to the three *a priori* subspaces. Finally, as ε becomes large, the number of segments reverts to 1, as it becomes most efficient to represent the points using a single zero-dimensional subspace (the origin).

This behavior contrasts with the phase transition discussed in [10]. There, the number of segments increases monotonically throughout the simulated annealing process. Because our formulation allows the dimension of the segments to vary, the number of segments does not decrease monotonically with ε . Notice, however, that the phase corresponding to the “correct” (*a priori*) segmentation is stable over several orders of magnitude of the parameter ε . This is important since in practice the true noise level ε_0 is usually unknown.

Another interesting thing to notice is that the coding rate $R^s(W)$ in many regions is mostly a linear function of $-\log_{10} \varepsilon$:

$$R^s(W) \approx -\beta \log_{10} \varepsilon + \alpha, \quad (18)$$

for some constants $\alpha, \beta > 0$, which is a typical characteristic of the rate-distortion function of Gaussians.

For this data set, the algorithm typically takes 11 seconds to run in Matlab on a 1.6GHz PC.

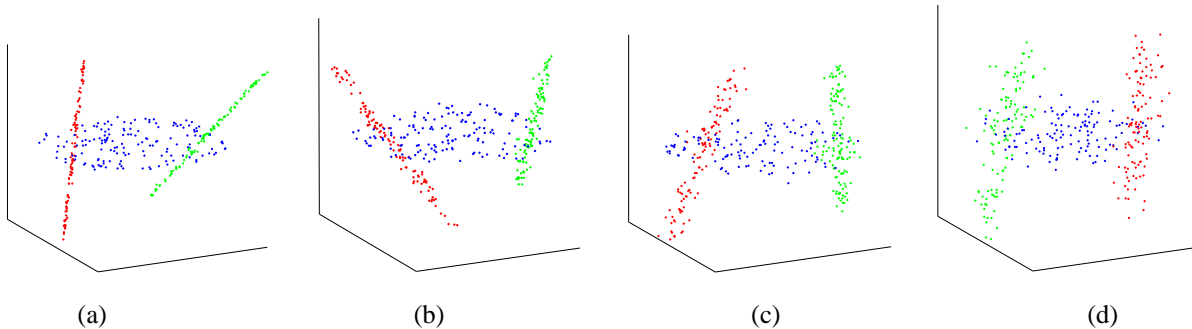


Fig. 9. The segmentation results for data drawn from 3 affine subspaces at different noise level ε_0 . The ε in the algorithm is chosen to be $\varepsilon = \varepsilon_0$. (a) $\varepsilon_0 = 0.01$, (b) $\varepsilon_0 = 0.03$, (c) $\varepsilon_0 = 0.05$, (d) $\varepsilon_0 = 0.08$.

5) *Segmentation of Affine Subspaces:* Appendix II shows how the coding length function should be properly modified in the case when the data are not zero-mean. Here, we show how the modified algorithm works for affine subspaces. A number of samples are drawn from three linear subspaces in \mathbb{R}^3 and their centers are translated to $[2.1, 2.2, 2]^T$, $[2.4, 1.9, 2.1]^T$, $[1.9, 2.5, 1.9]^T$, respectively.

Figure 9 shows the segmentation results at different noise level, with the distortion level chosen as $\varepsilon = \varepsilon_0$. For $10^{-7} < \varepsilon < 0.1$, the algorithm always identifies the correct number of subspaces with $\varepsilon = \varepsilon_0$. When $\varepsilon \leq 10^{-7}$, the density of the samples within the subspace becomes more important than the distortion orthogonal to the subspace, and the algorithm no longer converges with $\varepsilon = \varepsilon_0$. However, for such small distortion, there always exists a large stable phase (with respect to changing ε) corresponding to the correct number of subspaces, $k = 3$. When $\varepsilon_0 > 0.1$, the algorithm starts to fail and merge the data samples into one or two groups.

We now fix the Gaussian noise at $\varepsilon_0 = 0.02$, and add m_o outliers whose three coordinates are uniformly distributed in the range of $[1.5, 2.5]$, which is the same as the range of the inliers. When the number of outliers is $\leq m_o = 200$ (35.8% outliers), the algorithm finds the correct segmentation, and all the outlying samples are segmented into one group. From $m_o = 300$ (45.6% outliers) to $m_o = 700$ (66.2% outliers), the algorithm still identifies the two lines and one plane. However, the outliers above and below the plane are clustered into two separate groups. For more than $m_o = 800$ (69.1% outliers), the algorithm identifies the two lines, but samples from the plane are merged with the outliers into one group. Figure 10 shows the segmentation results for $m_o = 200, 300, 700, 800$, respectively.

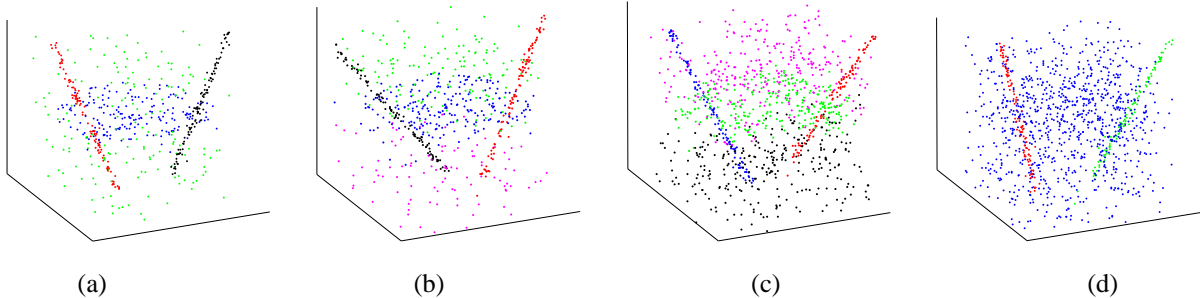


Fig. 10. The segmentation results for data drawn from 3 affine subspaces with different number of outliers m_o . The ε in the algorithm is $\varepsilon = \varepsilon_0 = 0.02$. (a) $m_o = 200$ (35.8% outliers), (b) $m_o = 300$ (45.6% outliers), (c) $m_o = 700$ (66.2% outliers), (d) $m_o = 800$ (69.1% outliers).

In this example, the algorithm takes 18 seconds to 9 minutes to run the simulation for different m_o 's in Matlab on a 1.6GHz PC.

6) *Model Selection for Affine Subspaces and Nonzero-Mean Gaussians:* We compare the performance of Algorithm 1 to that of [4] on mixed data drawn from affine subspaces and nonzero mean Gaussians. We test the algorithms' performance over multiple trials for three different types of data distribution. The first is three affine subspaces: two lines and one plane, with noise standard deviation $\varepsilon_0 = 0.01$ and no outliers. Samples are drawn as in the previous examples. The means of the three groups are fixed (as in the previous examples), but the orientations of the two lines are chosen randomly. The second distribution tested is three affine subspaces: two planes and one line, with 158 points drawn from each plane and 100 from the line, again with $\varepsilon_0 = 0.01$. The orientations of one plane and of the line are chosen randomly. The final distribution tested is a mixture of $K = 3$ full-rank Gaussians in \mathbb{R}^2 , with means $[2, 0]$, $[0, 0]$, $[0, 2]$ and covariance $\text{diag}(2, 0.2)$ (this is Figure 3 of [4]). 900 points are sampled (with uniform probability) from the three Gaussians.

For the two subspace examples, we run Algorithm 1 with $\varepsilon = \varepsilon_0 = 0.01$. For the third example, we set $\varepsilon = 0.2$. We repeat each trial 50 times. Figure 11 shows a histogram of the number of groups arrived at by the two algorithms. For both algorithms, all of the segmentations with $K = 3$ are essentially correct (classification error $< 4\%$). However, for degenerate, or subspace-like data (Figure 11(a) and Figure 11(b)), Algorithm 1 is much more likely to converge to the a-priori group number. For full-rank Gaussians (Figure 11(c)), Algorithm 1 performs quite well

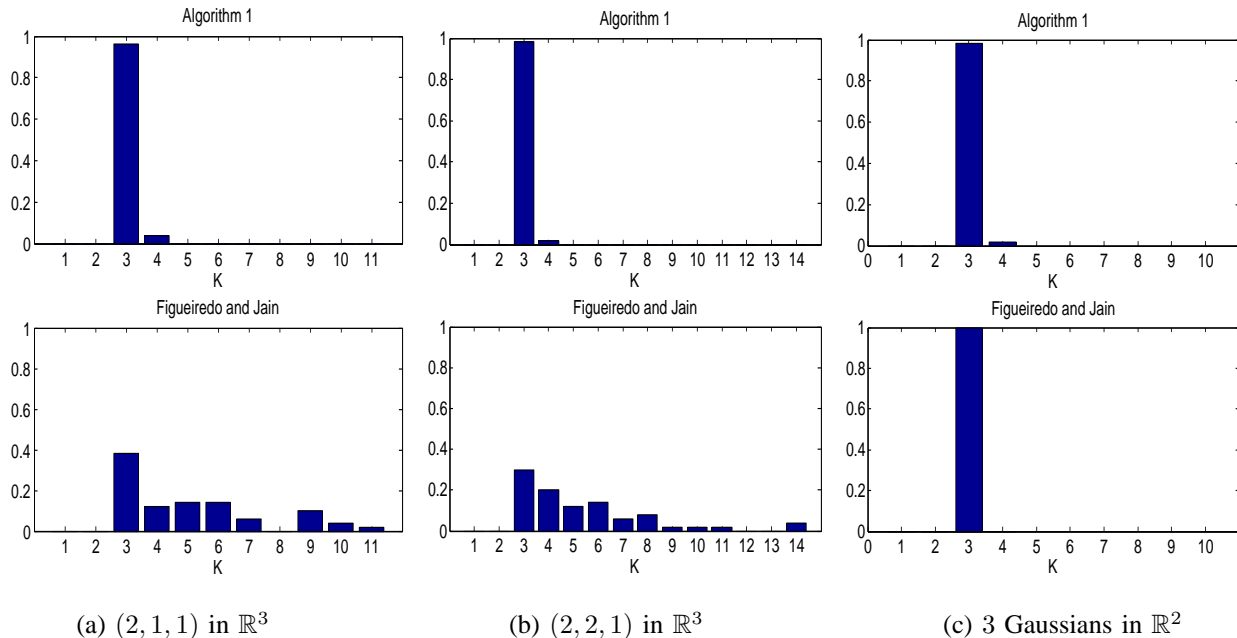


Fig. 11. Frequency of occurrence for various K in 50 trials. Top row: Algorithm 1. Bottom row: [4]. The left and center columns show results for randomly generated arrangements of affine subspaces. The right column shows results for datasets generated from three full-rank Gaussians, as in [4]. For all cases, the correct number of groups is $K = 3$.

(converging to the true group number in 49 out of 50 runs), but is slightly outperformed by [4], which finds the correct segmentation in all 50 trials. The single failure of Algorithm 1 occurs because the greedy descent converges to a local minimum of the coding length, rather than the global minimum.

B. Experiments on Real Data

In this section, we test the proposed segmentation method and algorithm on real imagery and bioinformatic data. Our goal here is to demonstrate that our method is capable of finding visually appealing structures in real data. However, we emphasize that it does not provide a complete solution to either of these practical problems. Such a solution usually entails a significant amount of domain-specific knowledge and engineering. Nevertheless, from these preliminary results with images and microarray data, we believe that the method presented in this paper provides a generic solution for segmenting mixed data that is simple and effective enough to be easily customized

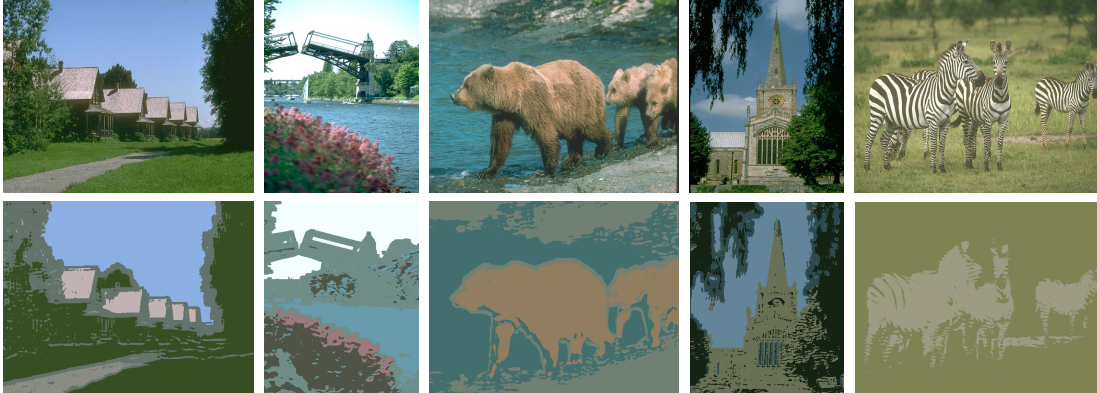


Fig. 12. Image segmentation (via the L formula for nonzero-mean Gaussian data) with $\varepsilon = 1$. Top row: original images. Bottom row: computed segmentations. Each segment is painted with its mean color.

for a broad range of practical problems.¹⁷

1) *Image Segmentation*: Figure 12 shows the segmentation of several images from the Berkeley image segmentation database via Algorithm 1 (using $L(\cdot)$ for nonzero-mean Gaussian data in Appendix II). The size of all the images is 480×320 pixels. We select an 8×8 window around each pixel to use as a feature vector¹⁸ for segmentation. A random subset of 1,000 vectors are selected. PCA is applied to these vectors, and they are projected onto their first 8 principal components. Subsampling and projection are necessary here due to the sheer volume of data. For a 480×320 color image, we are dealing with 153,600 vectors in an $8 \times 8 \times 3 = 192$ dimensional space, beyond the computational power and memory of a personal computer. The subsampled and projected vectors are clustered using Algorithm 1 with $\varepsilon = 1$. The remaining vectors are then grouped to the nearest segment. Figure 12 displays the results, without any further pre- or post-processing.

The segmentation can be further improved by first breaking the image into many small, homogeneous regions via a superpixel step. We compute the superpixel oversegmentation using the publically available code of [30]. We use its grouping to initialize the steepest descent procedure. To each pixel, we associate a 8×8 Gaussian-weighted window as a feature vector. Spatially adjacent groups are then repeatedly merged so as to achieve the

¹⁷At the time this paper is being prepared, we have also tested our algorithm on other mixed data such as speech and handwritten digits. The results are equally encouraging.

¹⁸Raw pixel values provide a simple and intuitive feature for testing our approach on real data. More visually appealing segmentations might be obtained with more sophisticated features (e.g. filterbanks [28], [29]). We leave this to future work.

greatest decrease in the coding length at each step. Figure 13 shows some representative results from the Berkeley segmentation database. The results for the entire database are available online at <http://perception.csl.uiuc.edu/~jnwright/hidden2/berkeleyResult.html>. A quantitative comparison of the performance of our method and several popular image segmentation algorithms will appear in future work.



Fig. 13. Segmentation results for greedily merging adjacent segments to decrease the coding length. Here, the merging process is initialized via a superpixel over-segmentation. $\varepsilon = 0.02$ for all images.

2) *Clustering of Microarray Data*: Figure 14 shows the result of applying Algorithm 1 to gene expression data. The dataset¹⁹ consists of 13,872 vectors in \mathbb{R}^{19} , each of which describes the expression level of a single gene at different time points during an experiment on anthrax sporulation. A random subset of 600 vectors is visualized in figure 14(a). Here, rows correspond to genes and columns to time points. We cluster these vectors without any preprocessing, using Algorithm 1 with $\varepsilon = 1$. The algorithm finds three distinct clusters, which are displayed in figure 14(b) by reordering the rows.

¹⁹GDS930, available at <http://www.ncbi.nlm.nih.gov/projects/geo>.

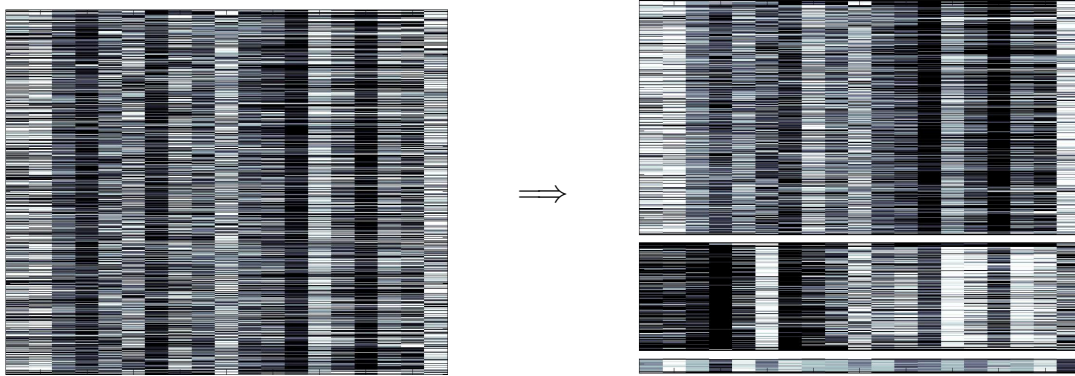


Fig. 14. Segmentation of microarray data. Left: raw data. Each row represents the expression level of a single gene. Right: Three distinct clusters are found, visualized by reordering the rows.

Figure 15 shows clustering results on two additional gene expression datasets²⁰. The first consists of 8,448 vectors in \mathbb{R}^5 , describing the expression levels of yeast genes at 5 different time points during a heat shock experiment. Figure 15(a) shows expression levels for a randomly selected subset of 1,200 genes. We cluster these vectors using Algorithm 1, with $\varepsilon = 0.1$. Our algorithm discovers a number of visually coherent clusters, shown in Figure 15(b). The second dataset consists of 45,101 vectors in \mathbb{R}^{10} , each of which corresponds to the expression level of a single gene under varying experimental conditions (this experiment investigated Down Syndrome-related leukemias). We run Algorithm 1 with $\varepsilon = 1$ on a subset of 800 of these vectors (shown in Figure 15). Three large, distinct clusters emerge, visualized in Figure 15(d) by reordering the rows of the data.

VI. CONCLUSIONS AND DISCUSSIONS

In this paper, we have proposed a new approach to segment multivariate mixed data from a lossy data coding/compression viewpoint. Unlike most conventional model-based top-down approaches to segmenting the data, our work leads to a data-driven bottom-up approach to obtain the optimal segmentation. In addition, this new approach allows us to examine explicitly the effect of a varying distortion on the segmentation result. From our experience, we find the lossy data compression based approach and the proposed greedy algorithm have the following attractive features:

²⁰GDS34 (left) and GDS1316 (right), also available at <http://www.ncbi.nlm.nih.gov/projects/geo>.

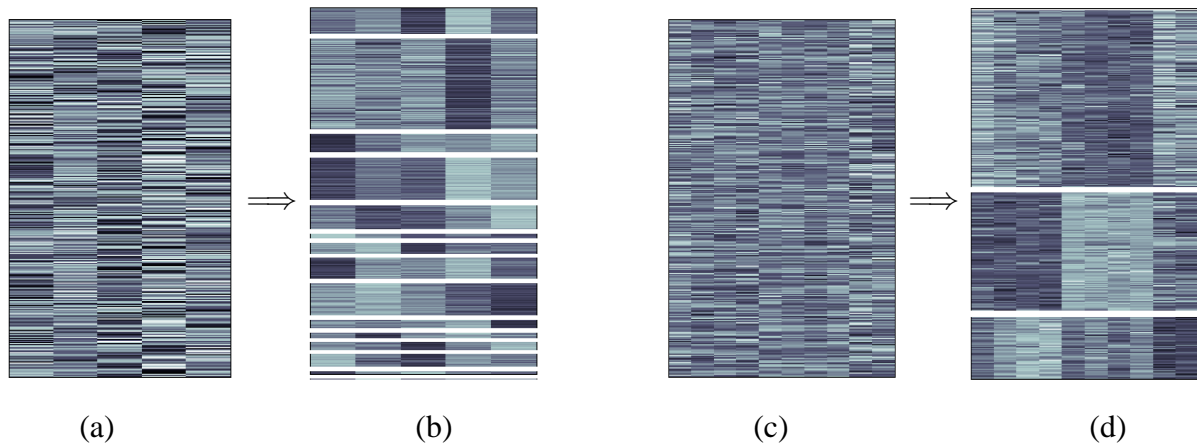


Fig. 15. Results on two microarray datasets. (a) raw yeast data. (b) segmentation, visualized by reordering rows. The greedy algorithm discovers a number of distinct clusters of varying size. (c) raw leukemia data. (d) segmentation. Three clusters are found.

- 1) The minimum coding length objective and the proposed greedy algorithm together deal with difficult issues such as outliers and model (number or dimension) selection. The segmentation result is very stable with respect to the distortion and the noise, and is very robust with respect to outliers.
- 2) The gain or loss of segmentation/merging is measured by a *physically meaningful* quantity – binary bits, and the (simulation) results even resemble the physical phenomenon of phase transition.
- 3) The greedy algorithm harnesses the tightness of the proposed coding length function for small sets of samples and takes a *bottom-up* approach that starts from merging the data one at a time. Thus it needs no initialization and the optimal segmentation is obtained without knowing anything about the (underlying) subspace(s) or Gaussian model(s).
- 4) The greedy algorithm is *scalable* – its complexity is polynomial in both the number of samples and the dimension of the data. The algorithm usually does converge to the optimal solution as long as the distortion is not too small relative to the density of the samples in each subspace.

Our analysis has shown connections of data segmentation with many fundamental concepts and results in information theory. The simulations and experiments have suggested potential connections with phase transition in statistical physics. From a theoretical standpoint, it would

be highly desirable to obtain analytical conditions on the critical values of the distortion and the sampling density of outliers that can explain and predict the phase transition of the number of segments.

Moreover, we do not see any technical difficulty in extending this approach to supervised-learning for purposes such as detection, classification, and recognition. It may also be extended to segment other types of structures, such as non-Gaussian probabilistic distributions and nonlinear manifolds. As we have mentioned earlier in the paper, there are many possible ways to improve the efficiency or convergence of the greedy algorithm or even develop (or employ) new optimization algorithms to minimize the coding length function. We will investigate such possibilities in the future.

VII. ACKNOWLEDGMENT

We would like to thank Professor Bin Yu of the Statistics Department at the University of California at Berkeley for pointing out relevant references about minimum description length (MDL) and mixture distributions.

We would also like to thank Dr. Allen Yang for his contribution to the application to image segmentation and the new segmentation results shown in Figure 13.

REFERENCES

- [1] A. Jain, M. Murty, and P. Flynn, "Data clustering: a review," *ACM Computing Surveys*, vol. 31, no. 3, pp. 264–323, 1999.
- [2] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*. Springer, 2001.
- [3] M. Tipping and C. Bishop, "Mixtures of probabilistic principal component analyzers," *Neural Computation*, vol. 11, no. 2, pp. 443–482, 1999.
- [4] M. A. T. Figueiredo and A. K. Jain, "Unsupervised learning of finite mixture models," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 3, pp. 381–396, 2002.
- [5] R. Vidal, Y. Ma, and S. Sastry, "Generalized principal component analysis (GPCA)," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 12, pp. 1–15, 2005.
- [6] S. Lloyd, "Least squares quantization in PCM," *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 129–137, March 1982.
- [7] E. Forgy, "Cluster analysis of multivariate data: efficiency vs. interpretability of classifications (abstract)," *Biometrics*, vol. 21, pp. 768–769, 1965.
- [8] R. Jancey, "Multidimensional group analysis," *Austral. J. Botany*, vol. 14, pp. 127–130, 1966.
- [9] J. MacQueen, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, 1967, pp. 281–297.

- [10] K. Rose, “Deterministic annealing for clustering, compression, classification, regression, and related optimization problems,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2210–2239, 1998.
- [11] E. P. Xing, A. Y. Ng, M. I. Jordan, and S. Russell, “Distance metric learning, with application to clustering with side information,” in *Proceedings of NIPS*, 2002.
- [12] J. Ho, M. Yang, J. Lim, K. Lee, and D. Kriegman, “Clustering appearances of objects under varying illumination conditions,” in *Proceedings of International Conference on Computer Vision and Pattern Recognition*, 2003.
- [13] A. Dempster, N. Laird, and D. Rubin, “Maximum likelihood from incomplete data via the EM algorithm,” *Journal of the Royal Statistical Society*, vol. 39(B), pp. 1–38, 1977.
- [14] G. McLachlan and T. Krishnan, *The EM algorithm and extensions*. John Wiley & Sons, 1997.
- [15] Z. Ghahramani and G. E. Hinton, “The EM algorithm for mixtures of factor analyzers,” *Technical Report CRG-TR-96-1, Department of Computer Science, University of Toronto*, 1996.
- [16] N. Ueda, R. Nakan, and Z. Ghahramani, “SMEM algorithm for mixture models,” *Neural Computation*, vol. 12, pp. 2109–2128, 2000.
- [17] T. Cover and J. Thomas, *Elements of Information Theory*. Wiley Series in Telecommunications, 1991.
- [18] J. Rissanen, “Modeling by shortest data description,” *Automatica*, vol. 14, pp. 465–471, 1978.
- [19] A. Barron, J. Rissanen, and B. Yu, “The minimum description length principle in coding and modeling,” *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2743–2760, 1998.
- [20] M. Hansen and B. Yu, “Model selection and the principle of minimum description length,” *Journal of American Statistical Association*, vol. 96, pp. 746–774, 2001.
- [21] M. Madiman, M. Harrison, and I. Kontoyiannis, “Minimum description length vs. maximum likelihood in lossy data compression,” in *Proceedings of the 2004 IEEE International Symposium on Information Theory*, 2004.
- [22] K. Rose, “A mapping approach to rate-distortion computation and analysis,” *IEEE Transactions on Information Theory*, vol. 40, no. 6, pp. 1939–1952, 1994.
- [23] H. Benson, “Concave minimization: Theory, applications and algorithms,” in R. Horst and P.M. Pardalos, eds., *Handbook of Global Optimization*, 1994.
- [24] J. Hamkins and K. Zeger, “Gaussian source coding with spherical codes,” *IEEE Transactions on Information Theory*, vol. 48, no. 11, pp. 2980–2989, 2002.
- [25] D. Donoho, M. Vetterli, R. DeVore, and I. Daubechies, “Data compression and harmonic analysis,” *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2435–2476, 1998.
- [26] R. A. Horn and C. R. Johnson, *Matrix Analysis*. Cambridge University Press, 1985.
- [27] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [28] J. Malik, S. Belongie, T. Leung, and J. Shi, “Contour and texture analysis for image segmentation,” vol. 43, no. 1, pp. 7–27, 2001.
- [29] S. Zhu, C. Guo, Y. Wu, and Y. Wang, “What are textons,” in *Proceedings of European Conference on Computer Vision*, 2002, pp. 793–807.
- [30] G. Mori, “Guiding model search using segmentation,” in *Proceedings of IEEE International Conference on Computer Vision*, vol. 2, 2005, pp. 1417–1423.
- [31] D. Tse and P. Viswanath, *Fundamentals of Wireless Communications*. Cambridge University Press, 2005.

APPENDIX I

LOSSY CODING OF SUBSPACE-LIKE DATA

In Section III, we have shown that in principle, one can construct a coding scheme for a given set of data $W \in \mathbb{R}^{n \times m}$ such that the average number of bits needed to encode each vector is bounded by

$$R(W) = \frac{1}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} WW^T \right), \quad (19)$$

if W is drawn from a multivariate Gaussian distribution of covariance $\Sigma = \frac{1}{m} WW^T$. However, we do not know in the non-parametric setting (i.e. with finite number of samples), whether the above coding length is still of any good. In this appendix, we provide a constructive proof that $L(W) = (m+n)R(W)$ indeed gives a tight upper bound for the number of bits needed to encode W . One interesting feature of the construction is that the coding scheme apparently relies on coding the subspace spanned by the vectors (i.e., the singular vectors) and the coordinates of the vectors with respect to the subspace. Thus geometrically, minimizing the coding length (via segmentation) is essentially to reduce the “dimension” (of each subset) of the data.

Consider the singular value decomposition (SVD) of the data matrix $W = U\Sigma V^T$. Let $B = (b_{ij}) = \Sigma V^T$. The column vectors of $U = (u_{ij})$ form a basis for the subspace spanned by vectors in W , and the column vectors of B are the coordinates of the vectors with respect to this basis.

For coding purpose, we store the approximated matrices $U + \delta U$ and $B + \delta B$. The matrix W can be recovered as

$$W + \delta W \doteq (U + \delta U)(B + \delta B) = UB + \delta UB + U\delta B + \delta U\delta B. \quad (20)$$

Then $\delta W \approx \delta UB + U\delta B$ as entries of $\delta U\delta B$ are negligible when ε is small (relative to the data W). The squared error introduced to the entries of W are

$$\sum_{i,j} \delta w_{ij}^2 = \mathbf{tr}(\delta W \delta W^T) \approx \mathbf{tr}(U\delta B \delta B^T U^T + \delta U B B^T \delta U^T + \delta U B \delta B^T U^T + U \delta B B^T \delta U^T).$$

We may further assume that the coding errors δU and δB are zero-mean independent random variables. Using the fact that $\mathbf{tr}(AB) = \mathbf{tr}(BA)$, the expected squared error becomes

$$\mathbb{E}(\mathbf{tr}(\delta W \delta W^T)) = \mathbb{E}(\mathbf{tr}(\delta B \delta B^T)) + \mathbb{E}(\mathbf{tr}(\Sigma^2 \delta U^T \delta U)).$$

Now, let us encode each entry b_{ij} with a precision $\varepsilon' = \frac{\varepsilon}{\sqrt{n}}$ and u_{ij} with a precision $\varepsilon'' = \frac{\varepsilon\sqrt{m}}{\sqrt{\lambda_j n}}$, where λ_j is the j th eigenvalue of WW^T .²¹ This is equivalent to assume that the error δb_{ij} is uniformly distributed in the interval $[-\frac{\varepsilon}{\sqrt{n}}, \frac{\varepsilon}{\sqrt{n}}]$ and δu_{ij} is uniformly distributed in the interval $[-\frac{\varepsilon\sqrt{m}}{\sqrt{\lambda_j n}}, \frac{\varepsilon\sqrt{m}}{\sqrt{\lambda_j n}}]$. Under such a coding precision, it is easy to verify that

$$\mathbb{E}(\mathbf{tr}(\delta W \delta W^T)) \leq \frac{2\varepsilon^2 m}{3} < \varepsilon^2 m. \quad (21)$$

Then the mean squared error per vector in W is

$$\frac{1}{m} \mathbb{E}(\mathbf{tr}(\delta W \delta W^T)) < \varepsilon^2. \quad (22)$$

The number of bits to store the coordinates b_{ij} with precision $\varepsilon' = \frac{\varepsilon}{\sqrt{n}}$ is

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^m \frac{1}{2} \log_2 \left(1 + \left(\frac{b_{ij}}{\varepsilon'} \right)^2 \right) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m \log_2 \left(1 + \frac{b_{ij}^2 n}{\varepsilon^2} \right) \\ & \leq \frac{m}{2} \sum_{i=1}^n \log_2 \left(1 + \frac{n \sum_{j=1}^m b_{ij}^2}{m \varepsilon^2} \right) = \frac{m}{2} \sum_{i=1}^n \log_2 \left(1 + \frac{n \lambda_i}{m \varepsilon^2} \right). \end{aligned}$$

In the above inequality, we have applied the following inequality:

$$\frac{\log(1 + a_1) + \log(1 + a_2) + \cdots + \log(1 + a_n)}{n} \leq \log \left(1 + \frac{a_1 + a_2 + \cdots + a_n}{n} \right) \quad (23)$$

for nonnegative real numbers $a_1, a_2, \dots, a_n \geq 0$.

Similarly, the number of bits to store the entries of the singular vectors u_{ij} with precision $\varepsilon'' = \frac{\varepsilon\sqrt{m}}{\sqrt{\lambda_i n}}$ is

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^n \frac{1}{2} \log_2 \left(1 + \left(\frac{u_{ij}}{\varepsilon''} \right)^2 \right) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \log_2 \left(1 + \frac{u_{ij}^2 n^2 \lambda_j}{m \varepsilon^2} \right) \\ & \leq \frac{n}{2} \sum_{j=1}^n \log_2 \left(1 + \frac{n^2 \lambda_j \sum_{i=1}^n u_{ij}^2}{m \varepsilon^2} \right) = \frac{n}{2} \sum_{j=1}^n \log_2 \left(1 + \frac{n \lambda_j}{m \varepsilon^2} \right). \end{aligned}$$

Thus, for U and B together, we need a total of

$$L(W) = \frac{m+n}{2} \sum_{i=1}^n \log_2 \left(1 + \frac{n \lambda_i}{m \varepsilon^2} \right) = \frac{m+n}{2} \log_2 \det \left(I + \frac{n}{m \varepsilon^2} WW^T \right). \quad (24)$$

We thus have proved the statement given in the beginning of this section: $L(W) = (m+n)R(W)$ gives a good upper bound on the number of bits needed to encode W .

²¹Notice that ε'_j normally does not increase with the number of vectors m , because λ_j increases proportionally to m .

APPENDIX II

NON-ZERO MEAN DISTRIBUTION

In the above analysis, we have assumed that the given vectors $W = (w_1, w_2, \dots, w_m)$ are zero-mean. In general, these vectors may have a non-zero mean. In other words, the points represented by these vectors may lie in an affine subspace, instead of a linear subspace.

In case W is not zero mean, let $\mu \doteq \frac{1}{m} \sum_{i=1}^m w_i \in \mathbb{R}^n$ and define the matrix

$$V \doteq \mu \cdot \mathbf{1}_{1 \times m} = (\mu, \mu, \dots, \mu) \in \mathbb{R}^{n \times m}. \quad (25)$$

Then $\bar{W} \doteq W - V$ is a matrix whose column vectors have zero mean. We may apply the same coding scheme in the previous section to \bar{W} .

Let $\bar{W} = U\Sigma V^T \doteq UB$ be the singular value decomposition of \bar{W} . Let $\delta U, \delta B, \delta \mu$ be the error in coding U, B, μ , respectively. Then the error induced on the matrix W is

$$\delta W = \delta \mu \cdot \mathbf{1}_{1 \times m} + U\delta B + \delta U B. \quad (26)$$

Assuming that $\delta U, \delta B, \delta \mu$ are zero-mean independent random variables, the expected total squared error is

$$\mathbb{E}(\mathbf{tr}(\delta W \delta W^T)) = m\mathbb{E}(\delta \mu^T \delta \mu) + \mathbb{E}(\mathbf{tr}(\delta B \delta B^T)) + \mathbb{E}(\mathbf{tr}(\Sigma \delta U^T \delta U)). \quad (27)$$

We encode entries of B and U with the same precision as before. We encode each entry μ_i of the mean vector μ with the precision $\varepsilon' = \frac{\varepsilon}{\sqrt{n}}$ and assume that the error $\delta \mu_i$ is a uniform distribution in the interval $[-\frac{\varepsilon}{\sqrt{n}}, \frac{\varepsilon}{\sqrt{n}}]$. Then we have $m\mathbb{E}(\delta \mu^T \delta \mu) = \frac{m\varepsilon^2}{3}$. Using equation (21) for the zero-mean case, the total squared error satisfies

$$\mathbb{E}(\mathbf{tr}(\delta W \delta W^T)) \leq \frac{m\varepsilon^2}{3} + \frac{2m\varepsilon^2}{3} = m\varepsilon^2. \quad (28)$$

Then the mean squared error per vector in W is still bounded by ε^2 :

$$\frac{1}{m}\mathbb{E}(\mathbf{tr}(\delta W \delta W^T)) \leq \varepsilon^2. \quad (29)$$

Now in addition to the $L(\bar{W})$ bits needed to encode U and B , the number of bits needed to encode the mean vector μ with precision $\varepsilon' = \frac{\varepsilon}{\sqrt{n}}$ is

$$\sum_{i=1}^n \frac{1}{2} \log_2 \left(1 + \left(\frac{\mu_i}{\varepsilon'} \right)^2 \right) = \frac{1}{2} \sum_{i=1}^n \log_2 \left(1 + \frac{n\mu_i^2}{\varepsilon^2} \right) \leq \frac{n}{2} \log_2 \left(1 + \frac{\mu^T \mu}{\varepsilon^2} \right), \quad (30)$$

where the last inequality is from the inequality (23).

Thus, the total number bits needed to store W is

$$L(W) = \frac{m+n}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} \bar{W} \bar{W}^T \right) + \frac{n}{2} \log_2 \left(1 + \frac{\mu^T \mu}{\varepsilon^2} \right). \quad (31)$$

Notice that if W is actually zero-mean, we have $\mu = 0$, $\bar{W} = W$, and the above expression for $L(W)$ is exactly the same as before.

APPENDIX III

RELATION TO MULTIPLE-CHANNEL CAPACITY

In wireless communication, the relationship between m transmitters and n receivers is often modeled as a fading multiple-input-multiple-output (MIMO) channel:

$$y = Wx + z, \quad (32)$$

where $y, z \in \mathbb{R}^n$ and $x \in \mathbb{R}^m$. z is a random vector that models the (additive) channel noise. It is often assumed that z has a Gaussian distribution $\mathcal{N}(0, \sigma^2 I)$. Then the model is known as the additive white Gaussian noise (AWGN) channel.

It is known in multiple-channel communications [31] that in the high signal-to-noise ratio (SNR) regime, the channel capacity is given by

$$C(W) \doteq \frac{1}{2} \log_2 \det \left(I + \frac{P}{m\sigma^2} W W^T \right), \quad (33)$$

where P is the total transmission power of the m transmitters [31]. The ratio P/σ^2 is the common SNR at each receiving antenna.

We could not help but notice a striking resemblance between the coding rate $R(W)$ in (11) and the wireless channel capacity $C(W)$ in (33). Notice that the noise variance σ^2 corresponds to the (entry-wise) mean squared error ε^2/n , and the power P is often assumed to be a constant and we may normalize it to be 1. Then the capacity becomes exactly the coding rate of W :

$$C(W) = R(W) = \frac{1}{2} \log_2 \det \left(I + \frac{n}{m\varepsilon^2} W W^T \right).$$

Thus, the concavity of the coding rate function $R^{s,\infty}(W, \Pi)$ (Theorem 3 in Section IV) suggests that an even higher channel capacity may be achieved by probabilistically assigning the transmitters into multiple groups. The capacity of such a probabilistic transmitting channel is a concave function in Π :

$$C(W, \Pi) \doteq \sum_{j=1}^k \frac{\mathbf{tr}(\Pi_j)}{2m} \log_2 \det \left(I + \frac{n}{\varepsilon^2 \mathbf{tr}(\Pi_j)} W \Pi_j W^T \right),$$

which has a unique maximum (for any given k).