

Feature Selection for Unsupervised Learning

Jennifer G. Dy

JDY@ECE.NEU.EDU

*Department of Electrical and Computer Engineering
Northeastern University
Boston, MA 02115, USA*

Carla E. Brodley

BRODLEY@ECN.PURDUE.EDU

*School of Electrical and Computer Engineering
Purdue University
West Lafayette, IN 47907, USA*

Editor: Stefan Wrobel

Abstract

In this paper, we identify two issues involved in developing an automated feature subset selection algorithm for unlabeled data: the need for finding the number of clusters in conjunction with feature selection, and the need for normalizing the bias of feature selection criteria with respect to dimension. We explore the feature selection problem and these issues through FSSEM (Feature Subset Selection using Expectation-Maximization (EM) clustering) and through two different performance criteria for evaluating candidate feature subsets: scatter separability and maximum likelihood. We present proofs on the dimensionality biases of these feature criteria, and present a cross-projection normalization scheme that can be applied to any criterion to ameliorate these biases. Our experiments show the need for feature selection, the need for addressing these two issues, and the effectiveness of our proposed solutions.

Keywords: clustering, feature selection, unsupervised learning, expectation-maximization

1. Introduction

In this paper, we explore the issues involved in developing automated feature subset selection algorithms for unsupervised learning. By unsupervised learning we mean unsupervised classification, or clustering. Cluster analysis is the process of finding “natural” groupings by grouping “similar” (based on some similarity measure) objects together.

For many learning domains, a human defines the features that are potentially useful. However, not all of these features may be relevant. In such a case, choosing a subset of the original features will often lead to better performance. Feature selection is popular in supervised learning (Fukunaga, 1990; Almuallim and Dietterich, 1991; Cardie, 1993; Kohavi and John, 1997). For supervised learning, feature selection algorithms maximize some function of predictive accuracy. Because we are given class labels, it is natural that we want to keep only the features that are related to or lead to these classes. But in unsupervised learning, we are not given class labels. Which features should we keep? Why not use all the information we have? The problem is that not all features are important. Some of the features may be redundant, some may be irrelevant, and some can even misguide clustering results. In addition, reducing the number of features increases comprehensibility and ameliorates the problem that some unsupervised learning algorithms break down with high dimensional data.

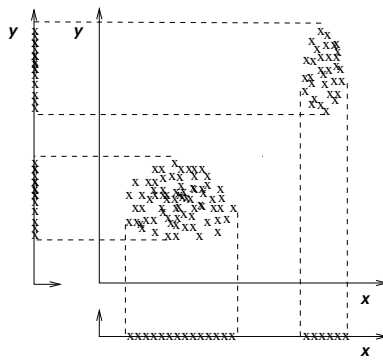


Figure 1: In this example, features x and y are redundant, because feature x provides the same information as feature y with regard to discriminating the two clusters.

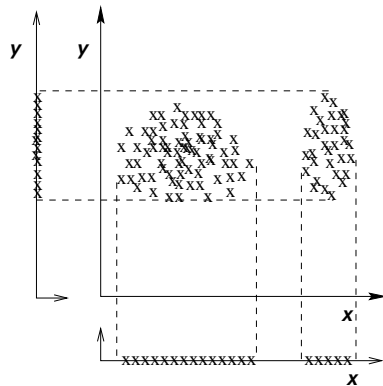


Figure 2: In this example, we consider feature y to be irrelevant, because if we omit x , we have only one cluster, which is uninteresting.

Figure 1 shows an example of feature redundancy for unsupervised learning. Note that the data can be grouped in the same way using only either feature x or feature y . Therefore, we consider features x and y to be redundant. Figure 2 shows an example of an irrelevant feature. Observe that feature y does not contribute to cluster discrimination. Used by itself, feature y leads to a single cluster structure which is uninteresting. Note that irrelevant features can misguide clustering results (especially when there are more irrelevant features than relevant ones). In addition, the situation in unsupervised learning can be more complex than what we depict in Figures 1 and 2. For example, in Figures 3a and b we show the clusters obtained using the feature subsets: $\{a, b\}$ and $\{c, d\}$ respectively. Different feature subsets lead to varying cluster structures. Which feature set should we pick?

Unsupervised learning is a difficult problem. It is more difficult when we have to simultaneously find the relevant features as well. A key element to the solution of any problem is to be able to precisely define the problem. In this paper, we define our task as:

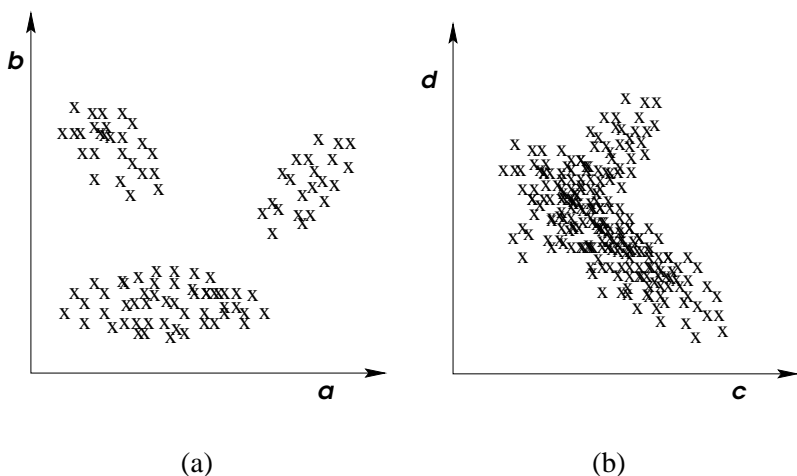


Figure 3: A more complex example. Figure a is the scatterplot of the data on features a and b . Figure b is the scatterplot of the data on features c and d .

The goal of feature selection for unsupervised learning is to find the smallest feature subset that best uncovers “interesting natural” groupings (clusters) from data according to the chosen criterion.

There may exist multiple redundant feature subset solutions. We are satisfied in finding any one of these solutions. Unlike supervised learning, which has class labels to guide the feature search, in unsupervised learning we need to define what “interesting” and “natural” mean. These are usually represented in the form of criterion functions. We present examples of different criteria in Section 2.3.

Since research in feature selection for unsupervised learning is relatively recent, we hope that this paper will serve as a guide to future researchers. With this aim, we

1. Explore the wrapper framework for unsupervised learning,
2. Identify the issues involved in developing a feature selection algorithm for unsupervised learning within this framework,
3. Suggest ways to tackle these issues,
4. Point out the lessons learned from this endeavor, and
5. Suggest avenues for future research.

The idea behind the wrapper approach is to cluster the data as best we can in each candidate feature subspace according to what “natural” means, and select the most “interesting” subspace with the minimum number of features. This framework is inspired by the supervised wrapper approach (Kohavi and John, 1997), but rather than wrap the search for the best feature subset around a supervised induction algorithm, we wrap the search around a clustering algorithm.

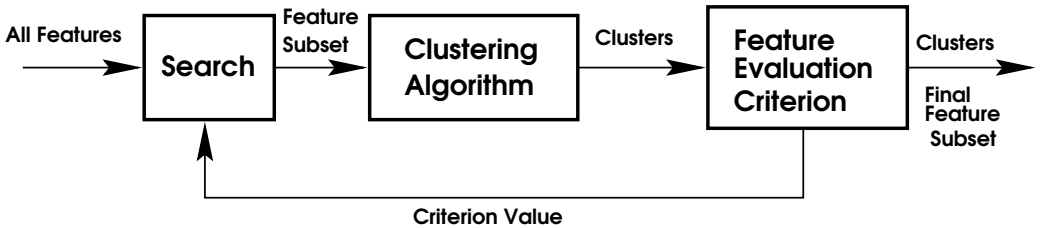


Figure 4: Wrapper approach for unsupervised learning.

In particular, this paper investigates the wrapper framework through FSSEM (feature subset selection using EM clustering) introduced in (Dy and Brodley, 2000a). Here, the term “EM clustering” refers to the expectation-maximization (EM) algorithm (Dempster et al., 1977; McLachlan and Krishnan, 1997; Moon, 1996; Wolfe, 1970; Wu, 1983) applied to estimating the maximum likelihood parameters of a finite Gaussian mixture. Although we apply the wrapper approach to EM clustering, the framework presented in this paper can be applied to any clustering method. FSSEM serves as an example. We present this paper such that applying a different clustering algorithm or feature selection criteria would only require replacing the corresponding clustering or feature criterion.

In Section 2, we describe FSSEM. In particular, we present the search method, the clustering method, and the two different criteria we selected to guide the feature subset search: scatter separability and maximum likelihood. By exploring the problem in the wrapper framework, we encounter and tackle two issues:

1. different feature subsets have different numbers of clusters, and
2. the feature selection criteria have biases with respect to feature subset dimensionality.

In Section 3, we discuss the complications that finding the number of clusters brings to the simultaneous feature selection/clustering problem and present one solution (FSSEM-k). Section 4 presents a theoretical explanation of why the feature selection criterion biases occur, and Section 5 provides a general normalization scheme which can ameliorate the biases of any feature criterion toward dimension.

Section 6 presents empirical results on both synthetic and real-world data sets designed to answer the following questions: (1) Is our feature selection for unsupervised learning algorithm better than clustering on all features? (2) Is using a fixed number of clusters, k , better than using a variable k in feature search? (3) Does our normalization scheme work? and (4) Which feature selection criterion is better? Section 7 provides a survey of existing feature selection algorithms. Section 8 provides a summary of the lessons learned from this endeavor. Finally, in Section 9, we suggest avenues for future research.

2. Feature Subset Selection and EM Clustering (FSSEM)

Feature selection algorithms can be categorized as either filter or wrapper (John et al., 1994) approaches. The filter approach basically pre-selects the features, and then applies the selected feature subset to the clustering algorithm. Whereas, the wrapper approach incorporates the clustering algorithm in the feature search and selection. We choose to explore the problem in the wrapper frame-

work because we are interested in understanding the interaction between the clustering algorithm and the feature subset search.

Figure 4 illustrates the wrapper approach. Our input is the set of all features. The output is the selected features and the clusters found in this feature subspace. The basic idea is to search through feature subset space, evaluating each candidate subset, F_t , by first clustering in space F_t using the clustering algorithm and then evaluating the resulting clusters and feature subset using our chosen feature selection criterion. We repeat this process until we find the best feature subset with its corresponding clusters based on our feature evaluation criterion. The wrapper approach divides the task into three components: (1) feature search, (2) clustering algorithm, and (3) feature subset evaluation.

2.1 Feature Search

An exhaustive search of the 2^d possible feature subsets (where d is the number of available features) for the subset that maximizes our selection criterion is computationally intractable. Therefore, a greedy search such as sequential forward or backward elimination (Fukunaga, 1990; Kohavi and John, 1997) is typically used. Sequential searches result in an $O(d^2)$ worst case search. In the experiments reported, we applied sequential forward search. Sequential forward search (SFS) starts with zero features and sequentially adds one feature at a time. The feature added is the one that provides the largest criterion value when used in combination with the features chosen. The search stops when adding more features does not improve our chosen feature criterion. SFS is not the best search method, nor does it guarantee an optimal solution. However, SFS is popular because it is simple, fast and provides a reasonable solution. For the purposes of our investigation in this paper, SFS would suffice. One may wish to explore other search methods for their wrapper approach. For example, Kim et al. (2002) applied evolutionary methods. Kittler (1978), and Russell and Norvig (1995) provide good overviews of different search strategies.

2.2 Clustering Algorithm

We choose EM clustering as our clustering algorithm, but other clustering methods can also be used in this framework. Recall that to cluster data, we need to make assumptions and define what “natural” grouping means. We apply the standard assumption that each of our “natural” groups is Gaussian. This assumption is not too limiting because we allow the number of clusters to adjust to our data, i.e., aside from finding the clusters we also find the number of “Gaussian” clusters. In Section 3, we discuss and present a solution to finding the number of clusters in conjunction with feature selection. We provide a brief description of EM clustering (the application of EM to approximate the maximum likelihood estimate of a finite mixture of multivariate Gaussians) in Appendix A. One can obtain a detailed description of EM clustering in (Fraley and Raftery, 2000; McLachlan and Krishnan, 1997). The Gaussian mixture assumption limits the data to continuous valued attributes. However, the wrapper framework can be extended to other mixture probability distributions (McLachlan and Basford, 1988; Titterton et al., 1985) and to other clustering methods, including graph theoretic approaches (Duda et al., 2001; Fukunaga, 1990; Jain and Dubes, 1988).

2.3 Feature Subset Selection Criteria

In this section, we investigate the feature subset evaluation criteria. Here, we define what “interestingness” means. There are two general views on this issue. One is that the criteria defining “interestingness” (feature subset selection criteria) should be the criteria used for clustering. The other is that the two criteria need not be the same. Using the same criteria for both clustering and feature selection provides a consistent theoretical optimization formulation. Using two different criteria, on the other hand, presents a natural way of combining two criteria for checks and balances. Proof on which view is better is outside the scope of this paper and is an interesting topic for future research. In this paper, we look at two feature selection criteria (one similar to our clustering criterion and the other with a different bias).

Recall that our goal is to find the feature subset that best discovers “interesting” groupings from data. To select an optimal feature subset, we need a measure to assess cluster quality. The choice of performance criterion is best made by considering the goals of the domain. In studies of performance criteria a common conclusion is: “Different classifications [clusterings] are right for different purposes, so we cannot say any one classification is best.” – Hartigan, 1985 .

In this paper, we do not attempt to determine the best criterion (one can refer to Milligan (1981) on comparative studies of different clustering criteria). We investigate two well-known measures: scatter separability and maximum likelihood. In this section, we describe each criterion, emphasizing the assumptions made by each.

Scatter Separability Criterion: A property typically desired among groupings is cluster separation. We investigate the scatter matrices and separability criteria used in discriminant analysis (Fukunaga, 1990) as our feature selection criterion. We choose to explore the scatter separability criterion, because it can be used with any clustering method.¹ The criteria used in discriminant analysis assume that the features we are interested in are features that can group the data into clusters that are unimodal and separable.

S_w is the within-class scatter matrix and S_b is the between class scatter matrix, and they are defined as follows:

$$S_w = \sum_{j=1}^k \pi_j E\{(X - \mu_j)(X - \mu_j)^T | \omega_j\} = \sum_{j=1}^k \pi_j \Sigma_j, \quad (1)$$

$$S_b = \sum_{j=1}^k \pi_j (\mu_j - M_o)(\mu_j - M_o)^T, \quad (2)$$

$$M_o = E\{X\} = \sum_{j=1}^k \pi_j \mu_j, \quad (3)$$

where π_j is the probability that an instance belongs to cluster ω_j , X is a d -dimensional random feature vector representing the data, k the number of clusters, μ_j is the sample mean vector of cluster ω_j , M_o is the total sample mean, Σ_j is the sample covariance matrix of cluster ω_j , and $E\{\cdot\}$ is the expected value operator.

1. One can choose to use the non-parametric version of this criterion measure (Fukunaga, 1990) for non-parametric clustering algorithms.

S_w measures how scattered the samples are from their cluster means. S_b measures how scattered the cluster means are from the total mean. We would like the distance between each pair of samples in a particular cluster to be as small as possible and the cluster means to be as far apart as possible with respect to the chosen similarity metric (Euclidean, in our case). Among the many possible separability criteria, we choose the $trace(S_w^{-1}S_b)$ criterion because it is invariant under any nonsingular linear transformation (Fukunaga, 1990). Transformation invariance means that once m features are chosen, any nonsingular linear transformation on these features does not change the criterion value. This implies that we can apply weights to our m features or apply any nonsingular linear transformation or projection to our features and still obtain the same criterion value. This makes the $trace(S_w^{-1}S_b)$ criterion more robust than other variants. $S_w^{-1}S_b$ is S_b normalized by the average cluster covariance. Hence, the larger the value of $trace(S_w^{-1}S_b)$ is, the larger the normalized distance between clusters is, which results in better cluster discrimination.

Maximum Likelihood (ML) Criterion: By choosing EM clustering, we assume that each grouping or cluster is Gaussian. We maximize the likelihood of our data given the parameters and our model. Thus, maximum likelihood (ML) tells us how well our model, here a Gaussian mixture, fits the data. Because our clustering criterion is ML, a natural criterion for feature selection is also ML. In this case, the “interesting” groupings are the “natural” groupings, i.e., groupings that are Gaussian.

3. The Need for Finding the Number of Clusters (FSSEM-k)

When we are searching for the best subset of features, we run into a new problem: that the number of clusters, k , depends on the feature subset. Figure 5 illustrates this point. In two dimensions (shown on the left) there are three clusters, whereas in one-dimension (shown on the right) there are only two clusters. Using a fixed number of clusters for all feature sets does not model the data in the respective subspace correctly.

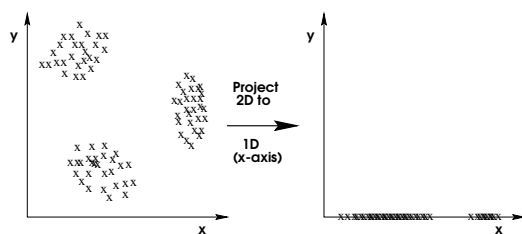


Figure 5: The number of cluster components varies with dimension.

Unsupervised clustering is made more difficult when we do not know the number of clusters, k . To search for k for a given feature subset, FSSEM-k currently applies Bouman et al.’s method (1998) for merging clusters and adds a Bayesian Information Criterion (BIC) (Schwarz, 1978) penalty term to the log-likelihood criterion. A penalty term is needed because the maximum likelihood estimate increases as more clusters are used. We do not want to end up with the trivial result wherein each data point is considered as an individual cluster. Our new objective function becomes: $F(k, \Phi) = \log(f(X|\Phi)) - \frac{1}{2}L\log(N)$ where N is the number of data points, L is the number of free

parameters in Φ , and $\log(f(X|\Phi))$ is the log-likelihood of our observed data X given the parameters Φ . Note that L and Φ vary with k .

Using Bouman et al.'s method (1998), we begin our search for k with a large number of clusters, K_{max} , and then sequentially decrement this number by one until only one cluster remains (a merge method). Other methods start from $k = 1$ and add more and more clusters as needed (split methods), or perform both split and merge operations (Ueda et al., 1999). To initialize the parameters of the $(k - 1)$ th model, two clusters from the k th model are merged. We choose the two clusters among all pairs of clusters in k , which when merged give the minimum difference between $F(k - 1, \Phi)$ and $F(k, \Phi)$. The parameter values that are not merged retain their value for initialization of the $(k - 1)$ th model. The parameters for the merged cluster (l and m) are initialized as follows:

$$\begin{aligned} \pi_j^{k-1,(0)} &= \pi_l + \pi_m; \\ \mu_j^{k-1,(0)} &= \frac{\pi_l \mu_l + \pi_m \mu_m}{\pi_l + \pi_m}; \\ \Sigma_j^{k-1,(0)} &= \frac{\pi_l (\Sigma_l + (\mu_l - \mu_j^{k-1,(0)}) (\mu_l - \mu_j^{k-1,(0)})^T) + \pi_m (\Sigma_m + (\mu_m - \mu_j^{k-1,(0)}) (\mu_m - \mu_j^{k-1,(0)})^T)}{\pi_l + \pi_m}; \end{aligned}$$

where the superscript $k - 1$ indicates the $k - 1$ cluster model and the superscript (0) indicates the first iteration in this reduced order model. For each candidate k , we iterate EM until the change in $F(k, \Phi)$ is less than ϵ (default 0.0001) or up to n (default 500) iterations. Our algorithm outputs the number of clusters k , the parameters, and the clustering assignments that maximize the $F(k, \Phi)$ criterion (our modified ML criterion).

There are myriad ways to find the “optimal” number of clusters k with EM clustering. These methods can be generally grouped into three categories: hypothesis testing methods (McLachlan and Basford, 1988), penalty methods like AIC (Akaike, 1974), BIC (Schwarz, 1978) and MDL (Rissanen, 1983), and Bayesian methods like AutoClass (Cheeseman and Stutz, 1996). Smyth (1996) introduced a new method called Monte Carlo cross-validation (MCCV). For each possible k value, the average cross-validated likelihood on M runs is computed. Then, the k value with the highest cross-validated likelihood is selected. In an experimental evaluation, Smyth showed that MCCV and AutoClass found k values that were closer to the number of classes than the k values found with BIC for their data sets. We chose Bouman et al.'s method with BIC, because MCCV is more computationally expensive. MCCV has complexity $O(MK_{max}^2 d^2 NE)$, where M is the number of cross-validation runs, K_{max} is the maximum number of clusters considered, d is the number of features, N is the number of samples and E is the average number of EM iterations. The complexity of Bouman et al.'s approach is $O(K_{max}^2 d^2 NE')$. Furthermore, for $k < K_{max}$, we do not need to re-initialize EM (because we merged two clusters from $k + 1$) resulting in $E' < E$. Note that in FSSEM, we run EM for each candidate feature subset. Thus, in feature selection, the total complexity is the complexity of each complete EM run times the feature search space. Recently, Figueiredo and Jain (2002) presented an efficient algorithm which integrates estimation and model selection for finding the number of clusters using minimum message length (a penalty method). It would be of interest for future work to examine these other ways for finding k coupled with feature selection.

4. Bias of Criterion Values to Dimension

Both feature subset selection criteria have biases with respect to dimension. We need to analyze these biases because in feature subset selection we compare the criterion values for subsets of dif-

ferent cardinality (corresponding to different dimensionality). In Section 5, we present a solution to this problem.

4.1 Bias of the Scatter Separability Criterion

The separability criterion prefers higher dimensionality; i.e., the criterion value monotonically increases as features are added assuming identical clustering assignments (Fukunaga, 1990; Narendra and Fukunaga, 1977). However, the separability criterion may not be monotonically increasing with respect to dimension when the clustering assignments change.

Scatter separability or the *trace* criterion prefers higher dimensions, intuitively, because data is more scattered in higher dimensions, and mathematically, because more features mean adding more terms in the *trace* function. Observe that in Figure 6, feature y does not provide additional discrimination to the two-cluster data set. Yet, the *trace* criterion prefers feature subset $\{x, y\}$ over feature subset $\{x\}$. Ideally, we would like the criterion value to remain the same if the discrimination information is the same.

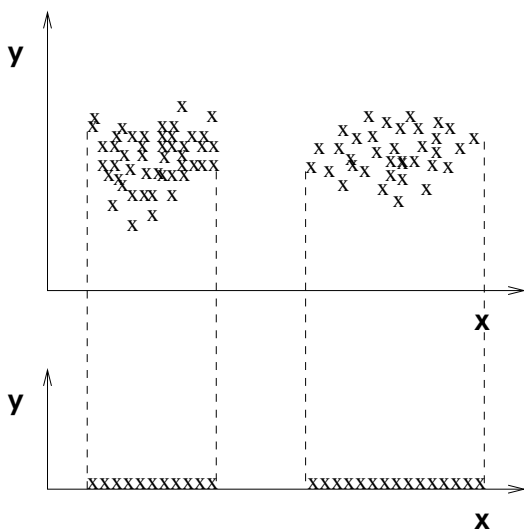


Figure 6: An illustration of scatter separability’s bias with dimension.

The following simple example provides us with an intuitive understanding of this bias. Assume that feature subset S_1 and feature subset S_2 produce identical clustering assignments, $S_1 \subset S_2$ where S_1 and S_2 have d and $d + 1$ features respectively. Assume also that the features are uncorrelated within each cluster. Let S_{w_d} and S_{b_d} be the within-class scatter and between-class scatter in dimension d respectively. To compute $trace(S_{w_{d+1}}^{-1} S_{b_{d+1}})$ for $d + 1$ dimensions, we simply add a positive term to the $trace(S_{w_d}^{-1} S_{b_d})$ value for d dimensions. $S_{w_{d+1}}$ and $S_{b_{d+1}}$ in the $d + 1$ dimensional space are computed as

$$S_{w_{d+1}} = \begin{bmatrix} S_{w_d} & 0 \\ 0 & \sigma_{w_{d+1}}^2 \end{bmatrix}$$

and

$$S_{b_{d+1}} = \begin{bmatrix} S_{b_d} & 0 \\ 0 & \sigma_{b_{d+1}}^2 \end{bmatrix}.$$

Since

$$S_{w_{d+1}}^{-1} = \begin{bmatrix} S_{w_d}^{-1} & 0 \\ 0 & \frac{1}{\sigma_{w_{d+1}}^2} \end{bmatrix},$$

$\text{trace}(S_{w_{d+1}}^{-1} S_{b_{d+1}})$ would be $\text{trace}(S_{w_d}^{-1} S_{b_d}) + \frac{\sigma_{b_{d+1}}^2}{\sigma_{w_{d+1}}^2}$. Since $\sigma_{b_{d+1}}^2 \geq 0$ and $\sigma_{w_{d+1}}^2 > 0$, the trace of the $d + 1$ clustering will always be greater than or equal to trace of the d clustering under the stated assumptions.

The separability criterion monotonically increases with dimension even when the features are correlated as long as the clustering assignments remain the same. Narendra and Fukunaga (1977) proved that a criterion of the form $X_d^T S_d^{-1} X_d$, where X_d is a d -column vector and S_d is a $d \times d$ positive definite matrix, monotonically increases with dimension. They showed that

$$X_{d-1}^T S_{d-1}^{-1} X_{d-1} = X_d^T S_d^{-1} X_d - \frac{1}{b} [(C^T : b) X_d]^2, \tag{4}$$

where

$$X_d = \begin{bmatrix} X_{d-1} \\ x_d \end{bmatrix},$$

$$S_d^{-1} = \begin{bmatrix} A & C \\ C^T & b \end{bmatrix},$$

X_{d-1} and C are $d - 1$ column vectors, x_d and b are scalars, A is a $(d - 1) \times (d - 1)$ matrix, and the symbol $:$ means matrix augmentation. We can show that $\text{trace}(S_{w_d}^{-1} S_{b_d})$ can be expressed as a criterion of the form $\sum_{j=1}^k X_{j d}^T S_d^{-1} X_{j d}$. S_{b_d} can be expressed as $\sum_{j=1}^k Z_{j b_d} Z_{j b_d}^T$ where $Z_{j b_d}$ is a d -column vector:

$$\begin{aligned} \text{trace}(S_{w_d}^{-1} S_{b_d}) &= \text{trace}(S_{w_d}^{-1} \sum_{j=1}^k Z_{j b_d} Z_{j b_d}^T) \\ &= \text{trace}(\sum_{j=1}^k S_{w_d}^{-1} Z_{j b_d} Z_{j b_d}^T) \\ &= \sum_{j=1}^k \text{trace}(S_{w_d}^{-1} Z_{j b_d} Z_{j b_d}^T) \\ &= \sum_{j=1}^k \text{trace}(Z_{j b_d}^T S_{w_d}^{-1} Z_{j b_d}), \end{aligned}$$

since $\text{trace}(A_{p \times q} B_{q \times p}) = \text{trace}(B_{q \times p} A_{p \times q})$ for any rectangular matrices $A_{p \times q}$ and $B_{q \times p}$.

Because $Z_{j b_d}^T S_{w_d}^{-1} Z_{j b_d}$ is scalar,

$$\sum_{j=1}^k \text{trace}(Z_{j b_d}^T S_{w_d}^{-1} Z_{j b_d}) = \sum_{j=1}^k Z_{j b_d}^T S_{w_d}^{-1} Z_{j b_d}.$$

Since each term monotonically increases with dimension, the summation also monotonically increases with dimension. Thus, the scatter separability criterion increases with dimension assuming the clustering assignments remain the same. This means that even if the new feature does not facilitate finding new clusters, the criterion function increases.

4.2 Bias of the Maximum Likelihood (ML) Criterion

Contrary to finding the number of clusters problem, wherein ML increases as the number of model parameters (k) is increased, in feature subset selection, ML prefers lower dimensions. In finding the number of clusters, we try to fit the best Gaussian mixture to the data. The data is fixed and we try to fit our model as best as we can. In feature selection, given different feature spaces, we select the feature subset that is best modeled by a Gaussian mixture.

This bias problem occurs because we define likelihood as the likelihood of the data corresponding to the candidate feature subset (see Equation 10 in Appendix B). To avoid this bias, the comparison can be between two complete (relevant and irrelevant features included) models of the data. In this case, likelihood is defined such that the candidate relevant features are modeled as dependent on the clusters, and the irrelevant features are modeled as having no dependence on the cluster variable. The problem with this approach is the need to define a model for the irrelevant features. Vaithyanathan and Dom uses this for document clustering (Vaithyanathan and Dom, 1999). The multinomial distribution for the relevant and irrelevant features is an appropriate model for text features in document clustering. In other domains, defining models for the irrelevant features may be difficult. Moreover, modeling irrelevant features means more parameters to predict. This implies that we still work with all the features, and as we mentioned earlier, algorithms may break down with high dimensions; we may not have enough data to predict all model parameters. One may avoid this problem by adding the assumption of independence among irrelevant features which may not be true. A poorly-fitting irrelevant feature distribution may cause the algorithm to select too many features. Throughout this paper, we use the maximum likelihood definition only for the relevant features.

For a fixed number of samples, ML prefers lower dimensions. The problem occurs when we compare feature set A with feature set B wherein set A is a subset of set B , and the joint probability of a single point (x, y) is less than or equal to its marginal probability (x) . For sequential searches, this can lead to the trivial result of selecting only a single feature.

ML prefers lower dimensions for discrete random features. The joint probability mass function of discrete random vectors X and Y is $p(X, Y) = p(Y|X)p(X)$. Since $0 \leq p(Y|X) \leq 1$, $p(X, Y) = p(Y|X)p(X) \leq p(X)$. Thus, $p(X)$ is always greater than or equal to $p(X, Y)$ for any X . When we deal with continuous random variables, as in this paper, the definition, $f(X, Y) = f(Y|X)f(X)$ still holds, where $f(\cdot)$ is now the probability density function. $f(Y|X)$ is always greater than or equal to zero. However, $f(Y|X)$ can be greater than one. The marginal density $f(X)$ is greater than or equal to the joint probability $f(X, Y)$ iff $f(Y|X) \leq 1$.

Theorem 4.1 *For a finite multivariate Gaussian mixture, assuming identical clustering assignments for feature subsets A and B with dimensions $d_B \geq d_A$, $ML(\Phi_A) \geq ML(\Phi_B)$ iff*

$$\prod_{j=1}^k \left(\frac{|\Sigma_B|_j}{|\Sigma_A|_j} \right)^{\pi_j} \geq \frac{1}{(2\pi e)^{(d_B - d_A)}},$$

where Φ_A represents the parameters and Σ_{A_j} is the covariance matrix modelling cluster j in feature subset A , π_j is the mixture proportion of cluster j , and k is the number of clusters.

Corollary 4.1 For a finite multivariate Gaussian mixture, assuming identical clustering assignments for feature subsets X and (X, Y) , where X and Y are disjoint, $ML(\Phi_X) \geq ML(\Phi_{XY})$ iff

$$\prod_{j=1}^k |\Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}|_j^{\pi_j} \geq \frac{1}{(2\pi e)^{d_Y}},$$

where the covariance matrix in feature subset (X, Y) is $\begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix}$, and d_Y is the dimension in Y .

We prove Theorem 4.1 and Corollary 4.1 in Appendix B. Theorem 4.1 and Corollary 4.1 reveal the dependencies of comparing the ML criterion for different dimensions. Note that each j th component of the left hand side term of Corollary 4.1 is the determinant of the conditional covariance of $f(Y|X)$. This covariance term is the covariance of Y eliminating the effects of the conditioning variable X , i.e., the conditional covariance does not depend on X . The right hand side is approximately equal to $(0.06)^{d_Y}$. This means that the ML criterion increases when the feature or feature subset to be added (Y) has a generalized variance (determinant of the covariance matrix) smaller than $(0.06)^{d_Y}$. Ideally, we would like our criterion measure to remain the same when the subsets reveal the same clusters. Even when the feature subsets reveal the same cluster, Corollary 4.1 informs us that ML decreases or increases depending on whether or not the generalized variance of the new features is greater than or less than a constant respectively.

5. Normalizing the Criterion Values: Cross-Projection Method

The arguments from the previous section illustrate that to apply the ML and *trace* criteria to feature selection, we need to normalize their values with respect to dimension. A typical approach to normalization is to divide by a penalty factor. For example, for the scatter criterion, we could divide by the dimension, d . Similarly for the ML criterion, we could divide by $\frac{1}{(2\pi e)^d}$. But, $\frac{1}{(2\pi e)^d}$ would not remove the covariance terms due to the increase in dimension. We could also divide $\log ML$ by d , or divide only the portions of the criterion affected by d . The problem with dividing by a penalty is that it requires specification of a different magic function for each criterion.

The approach we take is to project our clusters to the subspaces that we are comparing. Given two feature subsets, S_1 and S_2 , of different dimension, clustering our data using subset S_1 produces cluster C_1 . In the same way, we obtain the clustering C_2 using the features in subset S_2 . Which feature subset, S_1 or S_2 , enables us to discover better clusters? Let $CRIT(S_i, C_j)$ be the feature selection criterion value using feature subset S_i to represent the data and C_j as the clustering assignment. $CRIT(\cdot)$ represents either of the criteria presented in Section 2.3. We normalize the criterion value for S_1, C_1 as

$$normalizedValue(S_1, C_1) = CRIT(S_1, C_1) \cdot CRIT(S_2, C_1),$$

and, the criterion value for S_2, C_2 as

$$normalizedValue(S_2, C_2) = CRIT(S_2, C_2) \cdot CRIT(S_1, C_2).$$

If $normalizedValue(S_i, C_i) > normalizedValue(S_j, C_j)$, we choose feature subset S_i . When the normalized criterion values are equal for S_i and S_j , we favor the lower dimensional feature subset. The

choice of a product or sum operation is arbitrary. Taking the product will be similar to obtaining the geometric mean, and a sum with an arithmetic mean. In general, one should perform normalization based on the semantics of the criterion function. For example, geometric mean would be appropriate for likelihood functions, and an arithmetic mean for the log-likelihood.

When the clustering assignments resulting from different feature subsets, S_1 and S_2 , are identical (i.e., $C_1 = C_2$), the $normalizedValue(S_1, C_1)$ would be equal to the $normalizedValue(S_2, C_2)$, which is what we want. More formally:

Proposition 1 *Given that $C_1 = C_2$, equal clustering assignments, for two different feature subsets, S_1 and S_2 , then $normalizedValue(S_1, C_1) = normalizedValue(S_2, C_2)$.*

Proof: From the definition of $normalizedValue(\cdot)$ we have

$$normalizedValue(S_1, C_1) = CRIT(S_1, C_1) \cdot CRIT(S_2, C_1).$$

Substituting $C_1 = C_2$,

$$\begin{aligned} normalizedValue(S_1, C_1) &= CRIT(S_1, C_2) \cdot CRIT(S_2, C_2). \\ &= normalizedValue(S_2, C_2). \quad \square \end{aligned}$$

To understand why cross-projection normalization removes some of the bias introduced by the difference in dimension, we focus on $normalizedValue(S_1, C_1)$. The common factor is C_1 (the clusters found using feature subset S_1). We measure the criterion values on both feature subsets to evaluate the clusters C_1 . Since the clusters are projected on both feature subsets, the bias due to data representation and dimension is diminished. The normalized value focuses on the quality of the clusters obtained.

For example, in Figure 7, we would like to see whether subset S_1 leads to better clusters than subset S_2 . $CRIT(S_1, C_1)$ and $CRIT(S_2, C_2)$ give the criterion values of S_1 and S_2 for the clusters found in those feature subspaces (see Figures 7a and 7b). We project clustering C_1 to S_2 in Figure 7c and apply the criterion to obtain $CRIT(S_2, C_1)$. Similarly, we project C_2 to feature space S_1 to obtain the result shown in Figure 7d. We measure the result as $CRIT(S_1, C_2)$. For example, if $ML(S_1, C_1)$ is the maximum likelihood of the clusters found in subset S_1 (using Equation 10, Appendix B),² then to compute $ML(S_2, C_1)$, we use the same cluster assignments, C_1 , i.e., the $E[z_{ij}]$'s (the membership probabilities) for each data point x_i remain the same. To compute $ML(S_2, C_1)$, we apply the maximization-step EM clustering update equations (Equations 7-9 in Appendix A to compute the model parameters in the increased feature space, $S_2 = \{F_2, F_3\}$).

Since we project data in both subsets, we are essentially comparing criteria in the same number of dimensions. We are comparing $CRIT(S_1, C_1)$ (Figure 7a) with $CRIT(S_1, C_2)$ (Figure 7d) and $CRIT(S_2, C_1)$ (Figure 7c) with $CRIT(S_2, C_2)$ (Figure 7b). In this example, normalized *trace* chooses subset S_2 , because there exists a better cluster separation in both subspaces using C_2 rather than C_1 . Normalized ML also chooses subset S_2 . C_2 has a better Gaussian mixture fit (smaller variance clusters) in both subspaces (Figures 7b and d) than C_1 (Figures 7a and c). Note that the underlying

2. One can compute the maximum log-likelihood, log ML, efficiently as $Q(\Phi, \Phi) + H(\Phi, \Phi)$ by applying Lemma B.1 and Equation 16 in Appendix B. Lemma B.1 expresses the $Q(\cdot)$ in terms only of the parameter estimates. Equation 16, $H(\Phi, \Phi)$, is the cluster entropy which requires only the $E[z_{ij}]$ values. In practice, we work with log ML to avoid precision problems. The product $normalizedValue(\cdot)$ function then becomes $\log normalizedValue(S_i, C_i) = \log ML(S_i, C_i) + \log ML(S_j, C_i)$.

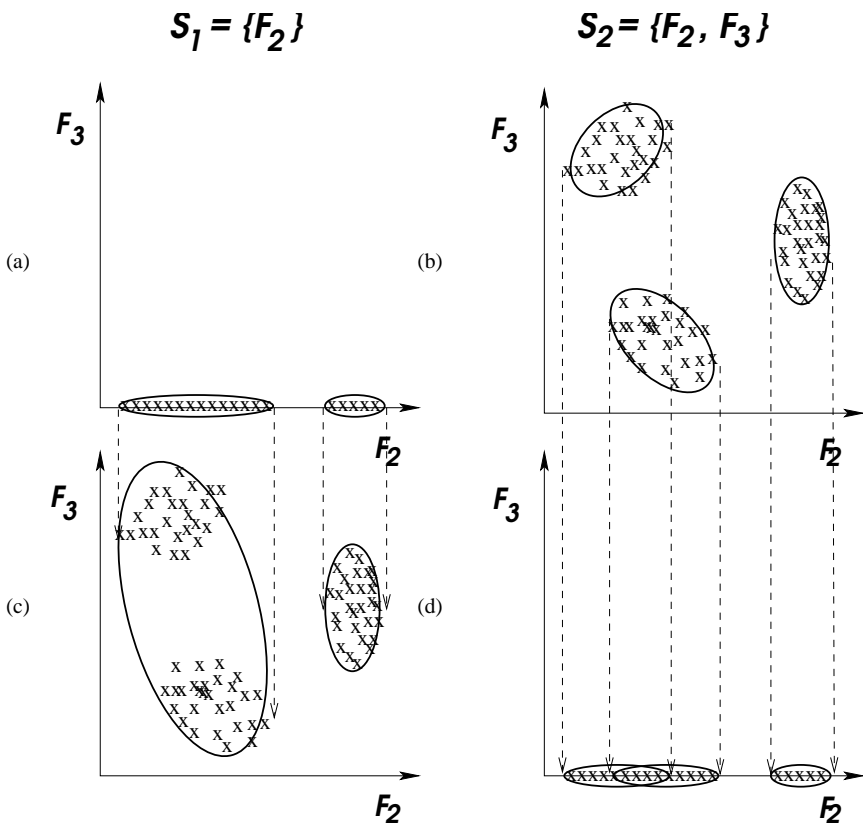


Figure 7: Illustration on normalizing the criterion values. To compare subsets, S_1 and S_2 , we project the clustering results of S_1 , we call C_1 in (a), to feature space S_2 as shown in (c). We also project the clustering results of S_2 , C_2 in (b), onto feature space S_1 as shown in (d). In (a), $tr(S_1, C_1) = 6.094$, $ML(S_1, C_1) = 1.9 \times 10^{-64}$, and $\log ML(S_1, C_1) = -146.7$. In (b), $tr(S_2, C_2) = 9.390$, $ML(S_2, C_2) = 4.5 \times 10^{-122}$, and $\log ML(S_1, C_2) = -279.4$. In (c), $tr(S_2, C_1) = 6.853$, $ML(S_2, C_1) = 3.6 \times 10^{-147}$, and $\log ML(S_2, C_1) = -337.2$. In (d), $tr(S_1, C_2) = 7.358$, $ML(S_1, C_2) = 2.1 \times 10^{-64}$, and $\log ML(S_1, C_2) = -146.6$. We evaluate subset S_1 with *normalized* $tr(S_1, C_1) = 41.76$ and subset S_2 with *normalized* $tr(S_2, C_2) = 69.09$. In the same way, using ML, the normalized values are: 6.9×10^{-211} for subset S_1 and 9.4×10^{-186} for subset S_2 . With log ML, the normalized values are: -483.9 and -426.0 for subsets S_1 and S_2 respectively.

assumption behind this normalization scheme is that the clusters found in the new feature space should be consistent with the structure of the data in the previous feature subset. For the ML criterion, this means that C_i should model S_1 and S_2 well. For the *trace* criterion, this means that the clusters C_i should be well separated in both S_1 and S_2 .

6. Experimental Evaluation

In our experiments, we 1) investigate whether feature selection leads to better clusters than using all the features, 2) examine the results of feature selection with and without criterion normalization, 3) check whether or not finding the number of clusters helps feature selection, and 4) compare the ML and the *trace* criteria. We first present experiments with synthetic data and then a detailed analysis of the FSSEM variants using four real-world data sets. In this section, we first describe our synthetic Gaussian data, our evaluation methods for the synthetic data, and our EM clustering implementation details. We then present the results of our experiments on the synthetic data. Finally, in Section 6.5, we present and discuss experiments with three benchmark machine learning data sets and one new real world data set.

6.1 Synthetic Gaussian Mixture Data

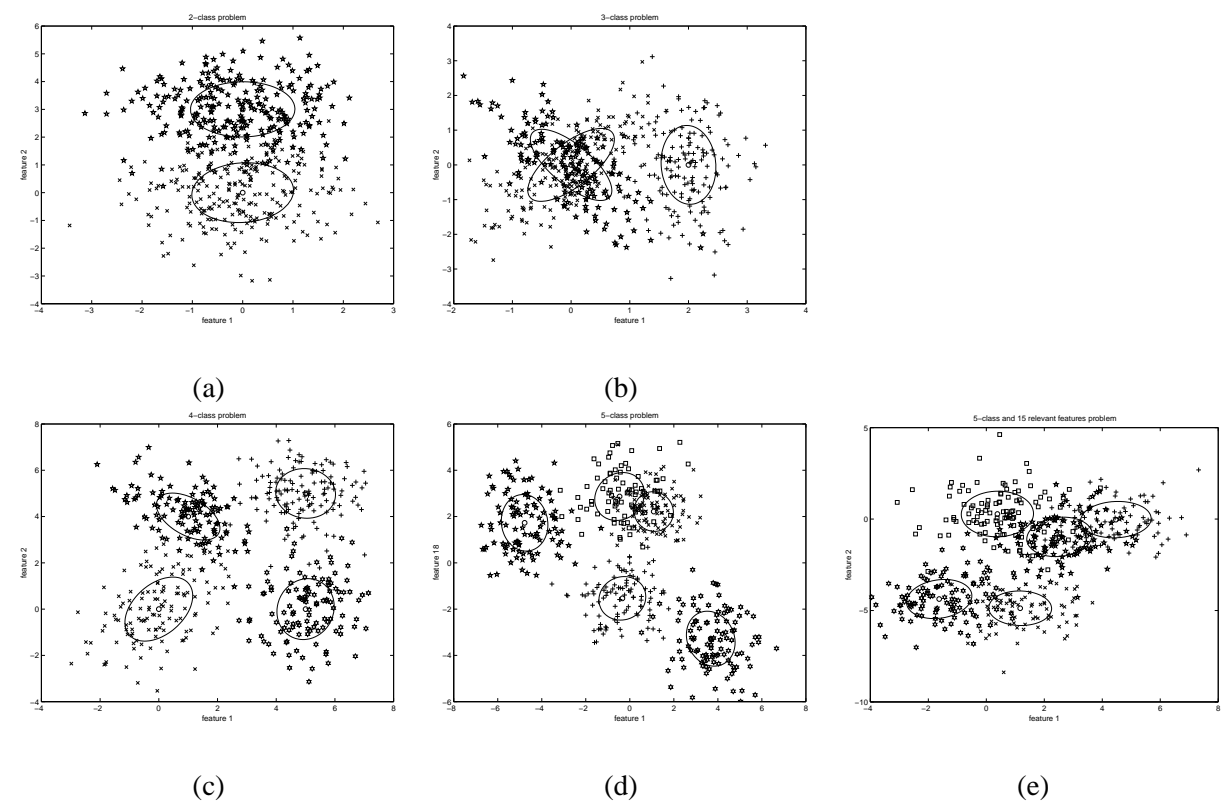


Figure 8: Synthetic Gaussian data.

To understand the performance of our algorithm, we experiment with five sets of synthetic Gaussian mixture data. For each data set we have “relevant” and “irrelevant” features, where relevant means that we created our k component mixture model using these features. Irrelevant features are generated as Gaussian normal random variables. For all five synthetic data sets, we generated $N = 500$ data points and generated clusters that are of equal proportions.

2-class, 2 relevant features and 3 noise features: The first data set (shown in Figure 8a) consists of two Gaussian clusters, both with covariance matrix, $\Sigma_1 = \Sigma_2 = I$ and means $\mu_1 = (0, 0)$ and $\mu_2 = (0, 3)$. This is similar to the two-class data set used by (Smyth, 1996). There is considerable overlap between the two clusters, and the three additional “noise” features increase the difficulty of the problem.

3-class, 2 relevant features and 3 noise features: The second data set consists of three Gaussian clusters and is shown in Figure 8b. Two clusters have means at $(0, 0)$ but the covariance matrices are orthogonal to each other. The third cluster overlaps the tails on the right side of the other two clusters. We add three irrelevant features to the three-class data set used by (Smyth, 1996).

4-class, 2 relevant features and 3 noise features: The third data set (Figure 8c) has four clusters with means at $(0, 0)$, $(1, 4)$, $(5, 5)$ and $(5, 0)$ and covariances equal to I . We add three Gaussian normal random “noise” features.

5-class, 5 relevant features and 15 noise features: For the fourth data set, there are twenty features, but only five are relevant (features $\{1, 10, 18, 19, 20\}$). The true means μ were sampled from a uniform distribution on $[-5, 5]$. The elements of the diagonal covariance matrices σ were sampled from a uniform distribution on $[0.7, 1.5]$ (Fayyad et al., 1998). Figure 8d shows the scatter plot of the data in two of its relevant features.

5-class, 15 relevant features and 5 noise features: The fifth data set (Figure 8e shown in two of its relevant features) has twenty features with fifteen relevant features $\{1, 2, 3, 5, 8, 9, 10, 11, 12, 13, 14, 16, 17, 18, 20\}$. The true means μ were sampled from a uniform distribution on $[-5, 5]$. The elements of the diagonal covariance matrices σ were sampled from a uniform distribution on $[0.7, 1.5]$ (Fayyad et al., 1998).

6.2 Evaluation Measures

We would like to measure our algorithm’s ability to select relevant features, to correctly identify k , and to find structure in the data (clusters). There are no standard measures for evaluating clusters in the clustering literature (Jain and Dubes, 1988). Moreover, no single clustering assignment (or class label) explains every application (Hartigan, 1985). Nevertheless, we need some measure of performance. Fisher (1996) provides and discusses different internal and external criteria for measuring clustering performance.

Since we generated the synthetic data, we know the ‘true’ cluster to which each instance belongs. This ‘true’ cluster is the component that generates that instance. We refer to these ‘true’ clusters as our known ‘class’ labels. Although we used the class labels to measure the performance of FSSEM, we did not use this information during training (i.e., in selecting features and discovering clusters).

Cross-Validated Class Error: We define class error as the number of instances misclassified divided by the total number of instances. We assign each data point to its most likely cluster, and assign each cluster to a class based on examining the class labels of the training data assigned to each cluster and choosing the majority class. Since we have the true cluster labels, we can compute classification error. One should be careful when comparing clusterings with

different number of clusters using training error. Class error based on training decreases with an increase in the number of clusters, k , with the trivial result of 0% error when each data point is a cluster. To ameliorate this problem, we use ten-fold cross-validation error. Ten-fold cross-validation randomly partitions the data set into ten mutually exclusive subsets. We consider each partition (or fold) as the test set and the rest as the training set. We perform feature selection and clustering on the training set, and compute class error on the test set. For each FSSEM variant, the reported error is the average and standard deviation values from the ten-fold cross-validation runs.

Bayes Error: Since we know the true probability distributions for the synthetic data, we provide the Bayes error (Duda et al., 2001) values to give us the lowest average class error rate achievable for these data sets. Instead of a full integration of the error in possibly discontinuous decision regions in multivariate space, we compute the Bayes error experimentally. Using the relevant features and their true distributions, we classify the generated data with an optimal Bayes classifier and calculate the error.

To evaluate the algorithm’s ability to select “relevant” features, we report the average number of features selected, and the average feature recall and precision. Recall and precision are concepts from text retrieval (Salton and McGill, 1983) and are defined here as:

Recall: the number of relevant features in the selected subset divided by the total number of relevant features.

Precision: the number of relevant features in the selected subset divided by the total number of features selected.

These measures give us an indication of the quality of the features selected. High values of precision and recall are desired. Feature precision also serves as a measure of how well our dimension normalization scheme (a.k.a. our stopping criterion) works. Finally, to evaluate the clustering algorithm’s ability to find the “correct” number of clusters, we report the average number of clusters found.

6.3 Initializing EM and Other Implementation Details

In the EM algorithm, we start with an initial estimate of our parameters, $\Phi^{(0)}$, and then iterate using the update equations until convergence. *Note that EM is initialized for each new feature subset.*

The EM algorithm can get stuck at a local maximum, hence the initialization values are important. We used the sub-sampling initialization algorithm proposed by Fayyad et al. (1998) with 10% sub-sampling and $J = 10$ sub-sampling iterations. Each sub-sample, S_i ($i = 1, \dots, J$), is randomly initialized. We run k -means (Duda et al., 2001) on these sub-samples not permitting empty clusters (i.e., when an empty cluster exists at the end of k -means, we reset the empty cluster’s mean equal to the data furthest from its cluster centroid, and re-run k -means). Each sub-sample results in a set of cluster centroids CM_i , i, \dots, J . We then cluster the combined set, CM , of all CM_i ’s using k -means initialized by CM_i resulting in new centroids FM_i . We select the FM_i , $i = 1, \dots, J$, that maximizes the likelihood of CM as our initial clusters.

After initializing the parameters, EM clustering iterates until convergence (i.e., the likelihood does not change by 0.0001) or up to n (default 500) iterations whichever comes first. We limit

the number of iterations because EM converges very slowly near a maximum. We avoid problems with handling singular matrices by adding a scalar ($\delta = 0.000001\sigma^2$, where σ^2 is the average of the variances of the unclustered data) multiplied to the identity matrix (δI) to each of the component covariance matrices Σ_j . This makes the final matrix positive definite (i.e., all eigenvalues are greater than zero) and hence nonsingular. We constrain our solution away from spurious clusters by deleting clusters with any diagonal element equal to or less than δ .

6.4 Experiments on Gaussian Mixture Data

We investigate the biases and compare the performance of the different feature selection criteria. We refer to FSSEM using the separability criterion as FSSEM-TR and using ML as FSSEM-ML. Aside from evaluating the performance of these algorithms, we also report the performance of EM (clustering using all the features) to see whether or not feature selection helped in finding more “interesting” structures (i.e., structures that reveal class labels). FSSEM and EM assume a fixed number of clusters, k , equal to the number of classes. We refer to EM clustering and FSSEM with finding the number of clusters as EM- k and FSSEM- k respectively. Due to clarity purposes and space constraints, we only present the relevant tables here. We report the results for all of the evaluation measures presented in Section 6.2 in (Dy and Brodley, 2003).

6.4.1 ML VERSUS TRACE

We compare the performance of the different feature selection criteria (FSSEM- k -TR and FSSEM- k -ML) on our synthetic data. We use FSSEM- k rather than FSSEM, because Section 6.4.3 shows that feature selection with finding k (FSSEM- k) is better than feature selection with fixed k (FSSEM). Table 1 shows the cross-validated (CV) error and average number of clusters results for *trace* and ML on the five data sets.

Percent CV Error					
Method	2-Class	3-Class	4-Class	5-Class, 5-Feat.	5-Class, 15-Feat.
FSSEM- k -TR	4.6 ± 2.0	21.4 ± 06.0	4.2 ± 2.3	3.0 ± 1.8	0.0 ± 0.0
FSSEM- k -ML	55.6 ± 3.9	54.8 ± 17.4	79.4 ± 6.1	84.0 ± 4.1	78.2 ± 6.1
Average Number of Clusters					
Method	2-Class	3-Class	4-Class	5-Class, 5-Feat.	5-Class, 15-Feat.
FSSEM- k -TR	2.0 ± 0.0	3.0 ± 0.0	4.0 ± 0.0	5.0 ± 0.0	5.0 ± 0.0
FSSEM- k -ML	1.0 ± 0.0	1.4 ± 0.8	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0

Table 1: Cross-validated error and average number of clusters for FSSEM- k -TR versus FSSEM- k -ML applied to the simulated Gaussian mixture data.

FSSEM- k -TR performed better than FSSEM- k -ML in terms of CV error. Trace performed better than ML, because it selected the features with high cluster separation. ML preferred features with low variance. When the variance of each cluster is the same, ML prefers the feature subset with fewer clusters (which happens to be our noise features). This bias is reflected by an average feature recall of 0.04. FSSEM- k -TR, on the other hand, was biased toward separable clusters identified by our defined relevant features, reflected by an average feature recall of 0.8.

6.4.2 RAW DATA VERSUS STANDARDIZED DATA

In the previous subsection, ML performed worse than *trace* for our synthetic data, because ML prefers features with low variance and fewer clusters (our noise features have lower variance than the relevant features). In this subsection, we investigate whether standardizing the data in each dimension (i.e., normalizing each dimension to yield a variance equal to one) would eliminate this bias. Standardizing data is sometimes done as a pre-processing step in data analysis algorithms to equalize the weight contributed by each feature. We would also like to know how standardization affects the performance of the other FSSEM variants.

Let X be a random data vector and X_f ($f = 1 \dots d$) be the elements of the vector, where d is the number of features. We standardize X by dividing each element by the corresponding feature standard deviation (X_f/σ_f , where σ_f is the standard deviation for feature f).

Table 2 reports the CV error. Additional experimental results can be found in (Dy and Brodley, 2003). Aside from the FSSEM variants, we examine the effect of standardizing data on EM-k, clustering with finding the number of clusters using all the features. We represent the corresponding variant on standardized data with the suffix “-STD”. The results show that only FSSEM-k-ML is affected by standardizing data. The *trace* criterion computes the between-class scatter normalized by the average within-class scatter and is invariant to any linear transformation. Since standardizing data is a linear transformation, the *trace* criterion results remain unchanged.

Standardizing data improves ML’s performance. It eliminates ML’s bias to lower overall variance features. Assuming equal variance clusters, ML prefers a single Gaussian cluster over two well-separated Gaussian clusters. But, after standardization, the two Gaussian clusters become more favorable because each of the two clusters now has lower variance (i.e., higher probabilities) than the single cluster noise feature. Observe that when we now compare FSSEM-k-TR-STD or FSSEM-k-TR with FSSEM-k-ML-STD, the performance is similar for all our data sets. These results show that scale invariance is an important property for a feature evaluation criterion. If a criterion is not scale invariant such as ML, in this case, pre-processing by standardizing the data in each dimension is necessary. Scale invariance can be incorporated to the ML criterion by modifying the function as presented in (Dy and Brodley, 2003). Throughout the rest of the paper, we standardize the data before feature selection and clustering.

Percent CV Error					
Method	2-Class	3-Class	4-Class	5-Class, 5-Feat.	5-Class, 15-Feat.
FSSEM-k-TR	4.6 ± 2.0	21.4 ± 06.0	4.2 ± 2.3	3.0 ± 1.8	0.0 ± 0.0
FSSEM-k-TR-STD	4.6 ± 2.0	21.6 ± 05.4	4.0 ± 2.0	3.0 ± 1.8	0.0 ± 0.0
FSSEM-k-ML	55.6 ± 3.9	54.8 ± 17.4	79.4 ± 6.1	84.0 ± 4.1	78.2 ± 6.1
FSSEM-k-ML-STD	4.8 ± 1.8	21.4 ± 05.1	4.0 ± 2.2	15.2 ± 7.3	0.0 ± 0.0
EM-k	55.6 ± 3.9	63.6 ± 06.0	48.6 ± 9.5	84.0 ± 4.1	55.4 ± 5.5
EM-k-STD	55.6 ± 3.9	63.6 ± 06.0	48.6 ± 9.5	84.0 ± 4.1	56.2 ± 6.1

Table 2: Percent CV error of FSSEM variants on standardized and raw data.

6.4.3 FEATURE SEARCH WITH FIXED k VERSUS SEARCH FOR k

In Section 3, we illustrated that different feature subsets have different numbers of clusters, and that to model the clusters during feature search correctly, we need to incorporate finding the number

of clusters, k , in our approach. In this section, we investigate whether finding k yields better performance than using a fixed number of clusters. We represent the FSSEM and EM variants using a fixed number of clusters (equal to the known classes) as FSSEM and EM. FSSEM- k and EM- k stand for FSSEM and EM with searching for k . Tables 3 and 4 summarize the CV error, average number of cluster, feature precision and recall results of the different algorithms on our five synthetic data sets.

Percent CV Error					
Method	2-Class	3-Class	4-Class	5-Class, 5-Feat.	5-Class, 15-Feat.
FSSEM-TR-STD	4.4 ± 02.0	37.6 ± 05.6	7.4 ± 11.0	21.2 ± 20.7	14.4 ± 22.2
FSSEM-k-TR-STD	4.6 ± 02.0	21.6 ± 05.4	4.0 ± 02.0	3.0 ± 01.8	0.0 ± 00.0
FSSEM-ML-STD	7.8 ± 05.5	22.8 ± 06.6	3.6 ± 01.7	15.4 ± 09.5	4.8 ± 07.5
FSSEM-k-ML-STD	4.8 ± 01.8	21.4 ± 05.1	4.0 ± 02.2	15.2 ± 07.3	0.0 ± 00.0
EM-STD	22.4 ± 15.1	30.8 ± 13.1	23.2 ± 10.1	48.2 ± 07.5	10.2 ± 11.0
EM-k-STD	55.6 ± 03.9	63.6 ± 06.0	48.6 ± 09.5	84.0 ± 04.1	56.2 ± 06.1
Bayes	5.4 ± 00.0	20.4 ± 00.0	3.4 ± 00.0	0.8 ± 00.0	0.0 ± 00.0
Average Number of Clusters					
Method	2-Class	3-Class	4-Class	5-Class, 5-Feat.	5-Class, 15-Feat.
FSSEM-TR-STD	fixed at 2	fixed at 3	fixed at 4	fixed at 5	fixed at 5
FSSEM-k-TR-STD	2.0 ± 0.0	3.0 ± 0.0	4.0 ± 0.0	5.0 ± 0.0	5.0 ± 0.0
FSSEM-ML-STD	fixed at 2	fixed at 3	fixed at 4	fixed at 5	fixed at 5
FSSEM-k-ML-STD	2.0 ± 0.0	3.0 ± 0.0	4.0 ± 0.0	4.2 ± 0.4	5.0 ± 0.0
EM-STD	fixed at 2	fixed at 3	fixed at 4	fixed at 5	fixed at 5
EM-k-STD	1.0 ± 0.0	1.0 ± 0.0	2.0 ± 0.0	1.0 ± 0.0	2.1 ± 0.3

Table 3: Percent CV error and average number of cluster results on FSSEM and EM with fixed number of clusters versus finding the number of clusters.

Looking first at FSSEM-k-TR-STD compared to FSSEM-TR-STD, we see that including order identification (FSSEM-k-TR-STD) with feature selection results in lower CV error for the *trace* criterion. For all data sets except the two-class data, FSSEM-k-TR-STD had significantly lower CV error than FSSEM-TR-STD. Adding the search for k within the feature subset selection search allows the algorithm to find the relevant features (an average of 0.796 feature recall for FSSEM-k-TR-STD versus 0.656 for FSSEM-TR-STD).³ This is because the best number of clusters depends on the chosen feature subset. For example, on closer examination, we noted that on the three-class problem when k is fixed at three, the clusters formed by feature 1 are better separated than clusters that are formed by features 1 and 2 together. As a consequence, FSSEM-TR-STD did not select feature 2. When k is made variable during the feature search, FSSEM-k-TR-STD finds two clusters in feature 1. When feature 2 is considered with feature 1, three or more clusters are found resulting in higher separability.

In the same way, FSSEM-k-ML-STD was better than fixing k , FSSEM-ML-STD, for all data sets in terms of CV error except for the four-class data. FSSEM-k-ML-STD performed slightly better than FSSEM-ML-STD for all the data sets in terms of feature precision and recall. This

3. Note that the recall value is low for the five-class fifteen-features data. This is because some of the “relevant” features are redundant as reflected by the 0.0% CV error obtained by our feature selection algorithms.

Average Feature Precision					
Method	2-Class	3-Class	4-Class	5-Class, 5-Feat.	5-Class, 15-Feat.
FSSEM-TR-STD	0.62 ± 0.26	0.56 ± 0.24	0.68 ± 0.17	0.95 ± 0.15	1.00 ± 0.00
FSSEM-k-TR-STD	0.57 ± 0.23	0.65 ± 0.05	0.53 ± 0.07	1.00 ± 0.00	1.00 ± 0.00
FSSEM-ML-STD	0.24 ± 0.05	0.52 ± 0.17	0.53 ± 0.10	0.98 ± 0.05	1.00 ± 0.00
FSSEM-k-ML-STD	0.33 ± 0.00	0.67 ± 0.13	0.50 ± 0.00	1.00 ± 0.00	1.00 ± 0.00
EM-k	0.20 ± 0.00	0.20 ± 0.00	0.20 ± 0.00	0.25 ± 0.00	0.75 ± 0.00
EM-k-STD	0.20 ± 0.00	0.20 ± 0.00	0.20 ± 0.00	0.25 ± 0.00	0.75 ± 0.00
Average Feature Recall					
Method	2-Class	3-Class	4-Class	5-Class, 5-Feat.	5-Class, 15-Feat.
FSSEM-TR-STD	1.00 ± 0.00	0.55 ± 0.15	0.95 ± 0.15	0.46 ± 0.20	0.32 ± 0.19
FSSEM-k-TR-STD	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.62 ± 0.06	0.36 ± 0.13
FSSEM-ML-STD	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.74 ± 0.13	0.41 ± 0.20
FSSEM-k-ML-STD	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.72 ± 0.16	0.51 ± 0.14
EM-k	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00
EM-k-STD	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00

Table 4: Average feature precision and recall obtained by FSSEM with a fixed number of clusters versus FSSEM with finding the number of clusters.

shows that incorporating finding k helps in selecting the “relevant” features. EM-STD had lower CV error than EM-k-STD due to prior knowledge about the correct number of clusters. Both EM-STD and EM-k-STD had poorer performance than FSSEM-k-TR/ML-STD, because of the retained noisy features.

6.4.4 FEATURE CRITERION NORMALIZATION VERSUS WITHOUT NORMALIZATION

Percent CV Error					
Method	2-Class	3-Class	4-Class	5-Class, 5-Feat.	5-Class, 15-Feat.
FSSEM-k-TR-STD-notnorm	4.6 ± 2.0	23.4 ± 6.5	4.2 ± 2.3	2.6 ± 1.3	0.0 ± 0.0
FSSEM-k-TR-STD	4.6 ± 2.0	21.6 ± 5.4	4.0 ± 2.0	3.0 ± 1.8	0.0 ± 0.0
FSSEM-k-ML-STD-notnorm	4.6 ± 2.2	36.2 ± 4.2	48.2 ± 9.4	63.6 ± 4.9	46.8 ± 6.2
FSSEM-k-ML-STD	4.8 ± 1.8	21.4 ± 5.1	4.0 ± 2.2	15.2 ± 7.3	0.0 ± 0.0
Bayes	5.4 ± 0.0	20.4 ± 0.0	3.4 ± 0.0	0.8 ± 0.0	0.0 ± 0.0
Average Number of Features Selected					
Method	2-Class	3-Class	4-Class	5-Class, 5-Feat.	5-Class, 15-Feat.
FSSEM-k-TR-STD-notnorm	2.30 ± 0.46	3.00 ± 0.00	3.90 ± 0.30	3.30 ± 0.46	9.70 ± 0.46
FSSEM-k-TR-STD	2.00 ± 0.63	3.10 ± 0.30	3.80 ± 0.40	3.10 ± 0.30	5.40 ± 1.96
FSSEM-k-ML-STD-notnorm	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00
FSSEM-k-ML-STD	3.00 ± 0.00	3.10 ± 0.54	4.00 ± 0.00	3.60 ± 0.80	7.70 ± 2.10

Table 5: Percent CV error and average number of features selected by FSSEM with criterion normalization versus without.

Table 5 presents the CV error and average number of features selected by feature selection with cross-projection criterion normalization versus without (those with suffix “notnorm”). Here and throughout the paper, we refer to normalization as the feature normalization scheme (cross-projection method) described in Section 5. For the *trace* criterion, without normalization did not affect the CV error. However, normalization achieved similar CV error performance using fewer features than without normalization. For the ML criterion, criterion normalization is definitely needed. Note that without, FSSEM-k-ML-STD-notnorm selected only a single feature for each data set resulting in worse CV error performance than with normalization (except for the two-class data which has only one relevant feature).

6.4.5 FEATURE SELECTION VERSUS WITHOUT FEATURE SELECTION

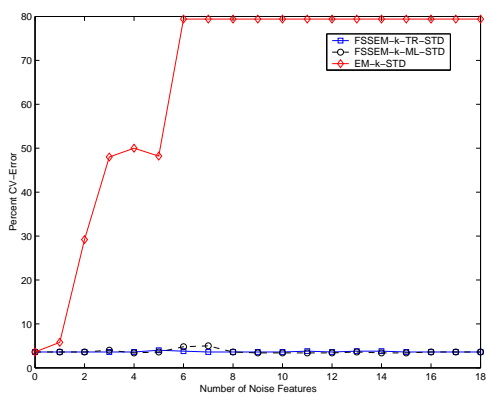
In all cases, feature selection (FSSEM, FSSEM-k) obtained better results than without feature selection (EM, EM-k) as reported in Table 3. Note that for our data sets, the noise features misled EM-k-STD, leading to fewer clusters than the “true” k . Observe too that FSSEM-k was able to find approximately the true number of clusters for the different data sets.

In this subsection, we experiment on the sensitivity of the FSSEM variants to the number of noise features. Figures 9a-e plot the cross-validation error, average number of clusters, average number of noise features, feature precision and recall respectively of feature selection (FSSEM-k-TR-STD and FSSEM-k-ML-STD) and without feature selection (EM-k-STD) as more and more noise features are added to the four-class data. Note that the CV error performance, average number of clusters, average number of selected features and feature recall for the feature selection algorithms are more or less constant throughout and are approximately equal to clustering with no noise. The feature precision and recall plots reveal that the CV error performance of feature selection was not affected by noise, because the FSSEM-k variants were able to select the relevant features (recall = 1) and discard the noisy features (high precision). Figure 9 demonstrates the need for feature selection as irrelevant features can mislead clustering results (reflected by EM-k-STD’s performance as more and more noise features are added).

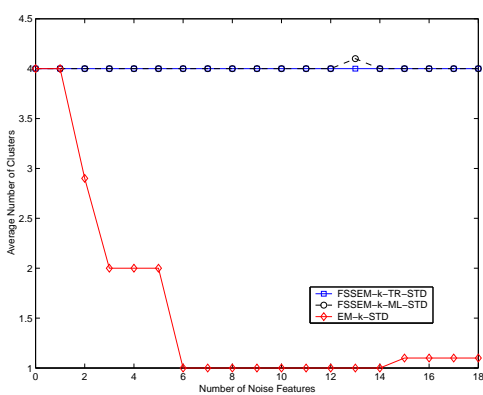
6.4.6 CONCLUSIONS ON EXPERIMENTS WITH SYNTHETIC DATA

Experiments on simulated Gaussian mixture data reveal that:

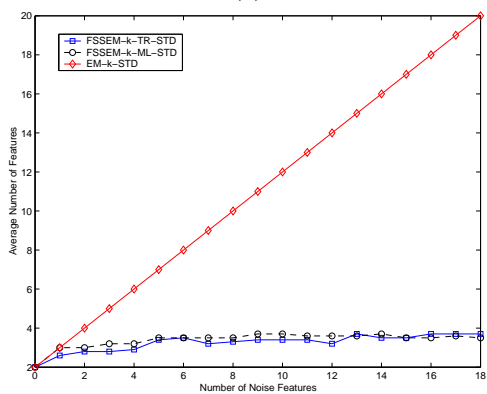
- Standardizing the data before feature subset selection in conjunction with the ML criterion is needed to remove ML’s preference for low variance features.
- Order identification led to better results than fixing k , because different feature subsets have different number of clusters as illustrated in Section 3.
- The criterion normalization scheme (cross-projection) introduced in Section 5 removed the biases of *trace* and ML with respect to dimension. The normalization scheme enabled feature selection with *trace* to remove “redundant” features and prevented feature selection with ML from selecting only a single feature (a trivial result).
- Both ML and *trace* with feature selection performed equally well for our five data sets. Both criteria were able to find the “relevant” features.
- Feature selection obtained better results than without feature selection.



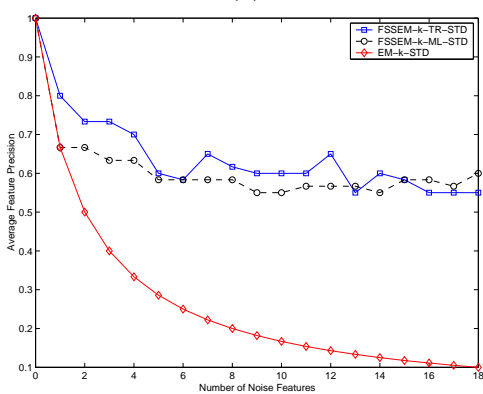
(a)



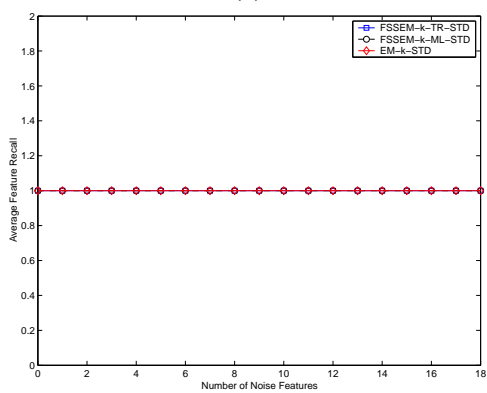
(b)



(c)



(d)



(e)

Figure 9: Feature selection versus without feature selection on the four-class data.

6.5 Experiments on Real Data

We examine the FSSEM variants on the iris, wine, and ionosphere data set from the UCI learning repository (Blake and Merz, 1998), and on a high resolution computed tomography (HRCT) lung

image data which we collected from IUPUI medical center (Dy et al., 2003; Dy et al., 1999). Although for each data set the class information is known, we remove the class labels during training.

Unlike synthetic data, we do not know the “true” number of (Gaussian) clusters for real-world data sets. Each class may be composed of many Gaussian clusters. Moreover, the clusters may not even have a Gaussian distribution. To see whether the clustering algorithms found clusters that correspond to classes (wherein a class can be multi-modal), we compute the cross-validated class error in the same way as for the synthetic Gaussian data. On real data sets, we do not know the “relevant” features. Hence, we cannot compute precision and recall and therefore report only the average number of features selected and the average number of clusters found.

Although we use class error as a measure of cluster performance, we should not let it misguide us in its interpretation. Cluster quality or interestingness is difficult to measure because it depends on the particular application. This is a major distinction between unsupervised clustering and supervised learning. Here, class error is just *one interpretation of the data*. We can also measure cluster performance in terms of the *trace* criterion and the ML criterion. Naturally, FSSEM-k-TR and FSSEM-TR performed best in terms of *trace*; and, FSSEM-k-ML and FSSEM-ML were best in terms of maximum likelihood. Choosing either TR or ML depends on your application goals. If you are interested in finding the features that best separate the data, use FSSEM-k-TR. If you are interested in finding features that model Gaussian clusters best, use FSSEM-k-ML.

To illustrate the generality and ease of applying other clustering methods in the wrapper framework, we also show the results for different variants of feature selection wrapped around the k -means clustering algorithm (Forgy, 1965; Duda et al., 2001) coupled with the TR and ML criteria. We use sequential forward search for feature search. To find the number of clusters, we apply the BIC penalty criterion (Pelleg and Moore, 2000). We use the following acronyms throughout the rest of the paper: Kmeans stands for the k -means algorithm, FSS-Kmeans stands for feature selection wrapped around k -means, TR represents the *trace* criterion for feature evaluation, ML represents ML criterion for evaluating features, “- k ” represents that the variant finds the number of clusters, and “-STD” shows that the data was standardized such that each feature has variance equal to one.

Since cluster quality depends on the initialization method used for clustering, we performed EM clustering using three different initialization methods:

1. Initialize using ten k -means starts with each k -means initialized by a random seed, then pick the final clustering corresponding to the highest likelihood.
2. Ten random re-starts.
3. Fayyad et al.’s method as described earlier in Section 6.3 (Fayyad et al., 1998).

Items one and two are similar for the k -means clustering. Hence, for k -means, we initialize with items two and three (with item three performed using Fayyad et al.’s method for k -means (Bradley and Fayyad, 1998) which applies k -means to the sub-sampled data instead of EM and distortion to pick the best clustering instead of the ML criterion). In the discussion section as follows, we show the results for FSSEM and FSS-Kmeans variants using the initialization which provides consistently good CV-error across all methods. We present the results using each initialization method on all the FSSEM and FSS-Kmeans variants in (Dy and Brodley, 2003) Appendix E. On the tables, “-1”, “-2”, and “-3” represent the initialization methods 1, 2, and 3 respectively.

Iris Data and FSSEM Variants			
Method	%CV Error	Ave. No. of Clusters	Ave. No. of Features
FSSEM-TR-STD-1	2.7 ± 04.4	fixed at 3	3.5 ± 0.7
FSSEM-k-TR-STD-1	4.7 ± 05.2	3.1 ± 0.3	2.7 ± 0.5
FSSEM-ML-STD-1	7.3 ± 12.1	fixed at 3	3.6 ± 0.9
FSSEM-k-ML-STD-1	3.3 ± 04.5	3.0 ± 0.0	2.5 ± 0.5
EM-STD-1	3.3 ± 05.4	fixed at 3	fixed at 4
EM-k-STD-1	42.0 ± 14.3	2.2 ± 0.6	fixed at 4
Iris Data and FSS-Kmeans Variants			
Method	%CV Error	Ave. No. of Clusters	Ave. No. of Features
FSS-Kmeans-TR-STD-2	2.7 ± 03.3	fixed at 3	1.9 ± 0.3
FSS-Kmeans-k-TR-STD-2	13.3 ± 09.4	4.5 ± 0.7	2.3 ± 0.5
FSS-Kmeans-ML-STD-2	2.0 ± 03.1	fixed at 3	2.0 ± 0.0
FSS-Kmeans-k-ML-STD-2	4.7 ± 04.3	3.4 ± 0.5	2.4 ± 0.5
Kmeans-STD-2	17.3 ± 10.8	fixed at 3	fixed at 4
Kmeans-k-STD-2	44.0 ± 11.2	2.0 ± 0.0	fixed at 4

Table 6: Results for the different variants on the iris data.

6.5.1 IRIS DATA

We first look at the simplest case, the Iris data. This data has three classes, four features, and 150 instances. Fayyad et. al’s method of initialization works best for large data sets. Since the Iris data only has a few number of instances and classes that are well-separated, ten k -means starts provided the consistently best result for initializing EM clustering across the different methods. Table 6 summarizes the results for the different variants of FSSEM compared to EM clustering without feature selection. For the iris data, we set K_{max} in FSSEM-k equal to six, and for FSSEM we fixed k at three (equal to the number of labeled classes). The CV error for FSSEM-k-TR-STD and FSSEM-k-ML-STD are much better than EM-k-STD. This means that when you do not know the “true” number of clusters, feature selection helps find good clusters. FSSEM-k even found the “correct” number of clusters. EM clustering with the “true” number of clusters (EM-STD) gave good results. Feature selection, in this case, did not improve the CV-error of EM-STD, however, they produced similar error rates with fewer features. FSSEM with the different variants consistently chose feature 3 (petal-length), and feature 4 (petal-width). In fact, we learned from this experiment that only these two features are needed to correctly cluster the iris data to three groups corresponding to iris-setosa, iris-versicolor and iris-viginica. Figures 10 (a) and (b) show the clustering results as a scatterplot on the first two features chosen by FSSEM-k-TR and FSSEM-k-ML respectively. The results for feature selection wrapped around k -means are also shown in Table 6. We can infer similar conclusions from the results on FSS-Kmeans variants as with the FSSEM variants for this data set.

6.5.2 WINE DATA

The wine data has three classes, thirteen features and 178 instances. For this data, we set K_{max} in FSSEM-k equal to six, and for FSSEM we fixed k at three (equal to the number of labeled classes). Table 7 summarizes the results when FSSEM and the FSS-Kmeans variants are initialized with ten k -means starts and ten random re-starts respectively. These are the initialization methods which led to the best performance for EM and k -means without feature selection. When “k” is

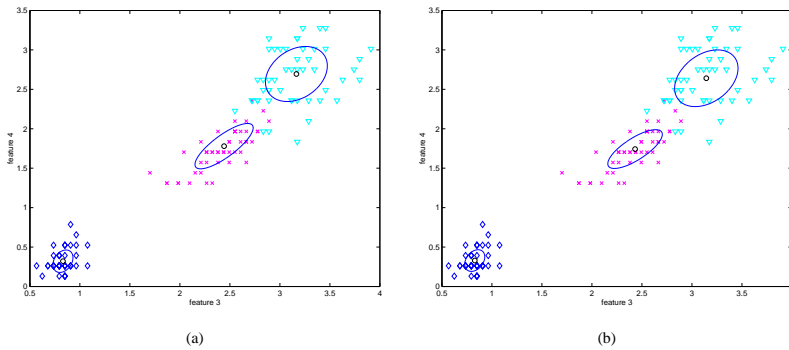


Figure 10: The scatter plots on iris data using the first two features chosen by FSSEM-k-TR (a) and FSSEM-k-ML (b). \diamond , \times and ∇ represent the different class assignments. \circ are the cluster means, and the ellipses are the covariances corresponding to the clusters discovered by FSSEM-k-TR and FSSEM-k-ML.

Wine Data and FSSEM Variants			
Method	%CV Error	Ave. No. of Clusters	Ave. No. of Features
FSSEM-TR-STD-1	44.0 ± 08.1	fixed at 3	1.4 ± 0.5
FSSEM-k-TR-STD-1	12.4 ± 13.0	3.6 ± 0.8	3.8 ± 1.8
FSSEM-ML-STD-1	30.6 ± 21.8	fixed at 3	2.9 ± 0.8
FSSEM-k-ML-STD-1	23.6 ± 14.4	3.9 ± 0.8	3.0 ± 0.8
EM-STD-1	10.0 ± 17.3	fixed at 3	fixed at 13
EM-k-STD-1	37.1 ± 12.6	3.2 ± 0.4	fixed at 13
Wine Data and FSS-Kmeans Variants			
Method	%CV Error	Ave. No. of Clusters	Ave. No. of Features
FSS-Kmeans-TR-STD-2	37.3 ± 14.0	fixed at 3	1.0 ± 0.0
FSS-Kmeans-k-TR-STD-2	28.1 ± 09.6	3.6 ± 0.5	2.5 ± 0.9
FSS-Kmeans-ML-STD-2	16.1 ± 09.9	fixed at 3	3.1 ± 0.3
FSS-Kmeans-k-ML-STD-2	18.5 ± 07.2	4.2 ± 0.6	3.1 ± 0.7
Kmeans-STD-2	0.0 ± 00.0	fixed at 3	fixed at 13
Kmeans-k-STD-2	33.4 ± 21.3	2.6 ± 0.8	fixed at 13

Table 7: Results for the different variants on the wine data set.

known, k -means was able to find the clusters corresponding to the “true” classes correctly. EM clustering also performed well when “ k ” is given. EM and k -means clustering performed poorly in terms of CV error when “ k ” is unknown. It is in this situation where feature selection, FSSEM-k and FSS-Kmeans-k, helped the base clustering methods find good groupings. Interestingly, for the wine data, FSSEM-k-TR performed better than FSSEM-k-ML, and FSS-Kmeans-ML had better CV-error than FSS-Kmeans-TR. This is an example on where using different criteria for feature selection and clustering improved the results through their interaction. Figures 11 (a) and (b) show the scatterplots and clusters discovered projected on the first two features chosen by FSSEM-k-TR and FSS-Kmeans-k-ML respectively. FSSEM-k-TR picked features $\{12, 13, 7, 5, 10, 1, 4\}$ and

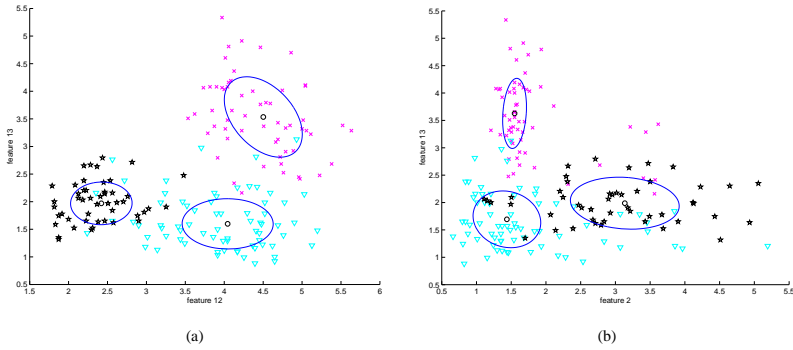


Figure 11: The scatter plots on the wine data using the first two features chosen by FSSEM-k-TR (a) and FSSEM-k-ML (b). \star , \times and ∇ represent the different class assignments. \circ are the cluster means, and the ellipses are the covariances corresponding to the clusters discovered by FSSEM-k-TR and FSS-Kmeans-k-ML.

FSS-Kmeans-k-ML selected features $\{2, 13, 12\}$.⁴ Features 12 and 13 stand for “OD280-OD315 of diluted wines” and “proline.”

6.5.3 IONOSPHERE DATA

The radar data is collected from a phased array of sixteen high-frequency antennas. The targets are free electrons in the atmosphere. Classes label the data as either good (radar returns showing structure in the ionosphere) or bad returns. There are 351 instances with 34 continuous attributes (measuring time of pulse and pulse number). Features 1 and 2 are discarded, because their values are constant or discrete for all instances. Constant feature values produce an infinite likelihood value for a Gaussian mixture model. Discrete feature values with discrete levels less than or equal to the number of clusters also produce an infinite likelihood value for a finite Gaussian mixture model.

Table 8 reports the ten-fold cross-validation error and the number of clusters found by the different EM and FSSEM algorithms. For the ionosphere data, we set K_{max} in FSSEM-k equal to ten, and fixed k at two (equal to the number of labeled classes) in FSSEM. FSSEM-k-ML and EM clustering with “k” known performed better in terms of CV error compared to the rest of the EM variants. Note that FSSEM-k-ML gave comparable performance with EM using fewer features and with no knowledge of the “true” number of clusters. Table 8 also shows the results for the different k -means variants. FSS-Kmeans-k-ML-STD obtains the best CV error followed closely by FSS-Kmeans-ML-STD. Interestingly, these two methods and FSSEM-k-ML all chose features 5 and 3 (based on the original 34 features) as their first two features.

Figures 12a and b present scatterplots of the ionosphere data on the first two features chosen by FSSEM-k-TR and FSSEM-k-ML together with their corresponding means (in \circ 's) and covariances (in ellipses) discovered. Observe that FSSEM-k-TR favored the clusters and features in Figure 12a because the clusters are well separated. On the other hand, FSSEM-k-ML favored the clusters in Figure 12b, which have small generalized variances. Since the ML criterion matches the ionosphere

⁴. These feature subsets are the features which provided the best CV-error performance among the ten-fold runs.

Ionosphere Data and FSSEM Variants			
Method	%CV Error	Ave. No. of Clusters	Ave. No. of Features
FSSEM-TR-STD-2	38.5 ± 06.5	fixed at 2	2.3 ± 0.9
FSSEM-k-TR-STD-2	23.1 ± 05.8	6.6 ± 1.3	1.1 ± 0.3
FSSEM-ML-STD-2	37.9 ± 07.5	fixed at 2	2.7 ± 2.1
FSSEM-k-ML-STD-2	18.8 ± 06.9	7.6 ± 1.0	2.9 ± 1.1
EM-STD-2	16.8 ± 07.3	fixed at 2	fixed at 32
EM-k-STD-2	35.3 ± 10.3	8.4 ± 1.0	fixed at 32
Ionosphere Data and FSS-Kmeans Variants			
Method	%CV Error	Ave. No. of Clusters	Ave. No. of Features
FSS-Kmeans-TR-STD-2	35.3 ± 06.5	fixed at 2	1.0 ± 0.0
FSS-Kmeans-k-TR-STD-2	22.8 ± 08.5	9.8 ± 0.4	1.0 ± 0.0
FSS-Kmeans-ML-STD-2	17.7 ± 04.9	fixed at 2	3.5 ± 0.8
FSS-Kmeans-k-ML-STD-2	16.2 ± 04.8	9.3 ± 0.8	1.7 ± 0.8
Kmeans-STD-2	23.4 ± 10.1	fixed at 2	fixed at 32
Kmeans-k-STD-2	28.8 ± 10.8	7.7 ± 0.6	fixed at 32

Table 8: Results for the different variants on the ionosphere data set.

class labels more closely, FSSEM-k-ML performed better with respect to CV error. FSSEM-k-ML obtained better CV error than EM-k; FSS-Kmeans-ML and FSS-Kmeans-k-ML also performed better than Kmeans and Kmeans-k in terms of CV error. The feature selection variants performed better using fewer features compared to the 32 features used by EM-k, Kmeans, and Kmeans-k. It is interesting to note that for this data, random re-start initialization obtained significantly better CV error for EM clustering (16.8%) compared to the other initialization methods (20.5% and 24.8% for ten k -means starts and Fayyad et al.’s method respectively). This is because the two “true” classes are highly overlapped. Ten k -means starts tend to start-off with well-separated clusters.

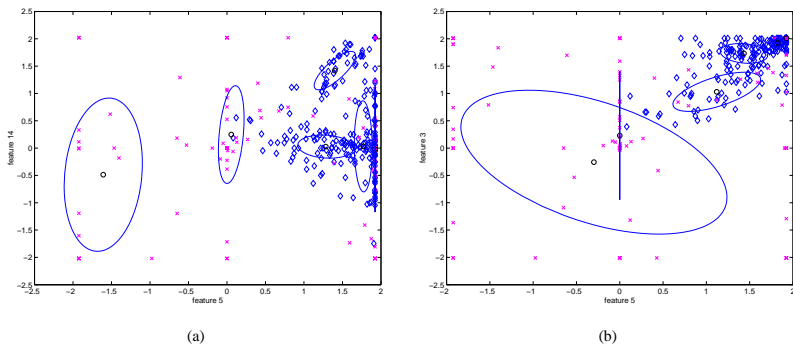


Figure 12: The scatter plots on the ionosphere data using the first two features chosen by FSSEM-k-TR (a) and FSSEM-k-ML (b). \times and \diamond represent the different class assignments. \circ are the cluster means, and the ellipses are the covariances corresponding to the clusters discovered by FSSEM-k-TR and FSSEM-k-ML.

6.5.4 HRCT-LUNG DATA



(a) Centrilobular Emphysema

(b) Paraseptal Emphysema

(c) IPF

Figure 13: HRCT-lung images.

HRCT-lung consists of 500 instances. Each of these instances are represented by 110 low-level continuous features measuring geometric, gray level and texture features (Dy et al., 1999). We actually used only 108 features because two of the features are constant or discrete. We also log-transformed the data to make our features which are mostly positive real-valued numbers more Gaussian. For features with negative values (like the feature, local mean minus global mean), we add an offset making the minimum value equal to zero. We assign $\log(0)$ to be $\log(0.00000000000001)$. The data is classified into five disease classes (Centrilobular Emphysema, Paraseptal Emphysema, EG, IPF, and Panacinar). Figure 13 shows three HRCT-lung images from three of the disease classes. The white marking is the pathology bearing region (PBR) marked by a radiologist. An instance represents a PBR. An image may contain more than one PBR and more than one disease classification. Note that Centrilobular Emphysema (CE) is characterized by a large number of low intensity (darker) regions which may occupy the entire lung as in Figure 13a. Paraseptal Emphysema (PE) is also characterized by low intensity regions (see Figure 13b). Unlike CE, these regions occur near the boundaries or near fissures. The dark regions are usually separated by thin walls from their adjacent boundary or fissure. CE and PE can be further grouped according to disease severity characterized by the intensity of the regions. Lower intensities indicate more severe cases. The lung image of IPF is characterized by high intensities forming a “glass-like” structure as shown in Figure 13c. Feature selection is important for this data set, because EM clustering using all the features results in just one cluster.

Table 9 presents the results on the HRCT-lung data set. For the HRCT lung data, FSSEM-k-TR and FSS-Kmeans-k-TR performed better than FSSEM-k-ML and FSS-Kmeans-k-ML respectively in terms of CV error. Figures 14 (a) and (b) present scatterplots of the HRCT-lung data on the first two features chosen by FSSEM-k-TR and FSSEM-k-ML. Observe that the clusters found by FSSEM-k-TR are well separated and match the class labels well. FSSEM-k-ML, on the other hand, selects features that result in high-density clusters. Figure 14 (b) demonstrates this clearly. Note also that the “true” number of clusters for this data is more than five (the number of labeled classes). This helped FSSEM-k-TR and FSS-Kmeans-k-TR obtained better results than their fixed-k variants.

HRCT-Lung Data and FSSEM Variants			
Method	%CV Error	Ave. No. of Clusters	Ave. No. of Features
FSSEM-TR-STD-3	36.8 ± 6.6	fixed at 5	1.3 ± 0.5
FSSEM-k-TR-STD-3	26.6 ± 7.7	6.0 ± 2.7	1.7 ± 0.9
FSSEM-ML-STD-3	37.2 ± 5.5	fixed at 5	3.3 ± 0.6
FSSEM-k-ML-STD-3	37.0 ± 5.7	5.2 ± 1.7	6.6 ± 2.8
EM-STD-3	37.2 ± 5.5	fixed at 5	fixed at 108
EM-k-STD-3	37.2 ± 5.5	1.1 ± 0.3	fixed at 108
HRCT-Lung Data and FSS-Kmeans Variants			
Method	%CV Error	Ave. No. of Clusters	Ave. No. of Features
FSS-Kmeans-TR-STD-3	37.2 ± 05.5	fixed at 5	1.0 ± 0.0
FSS-Kmeans-k-TR-STD-3	28.0 ± 10.7	7.5 ± 1.9	2.9 ± 2.3
FSS-Kmeans-ML-STD-3	36.8 ± 05.9	fixed at 5	3.4 ± 0.7
FSS-Kmeans-k-ML-STD-3	35.6 ± 06.7	4.3 ± 0.9	5.8 ± 3.1
Kmeans-STD-3	36.6 ± 04.9	fixed at 5	fixed at 108
Kmeans-k-STD-3	37.0 ± 05.3	3.4 ± 0.5	fixed at 108

Table 9: Results on the HRCT-lung image data set for the different variants.

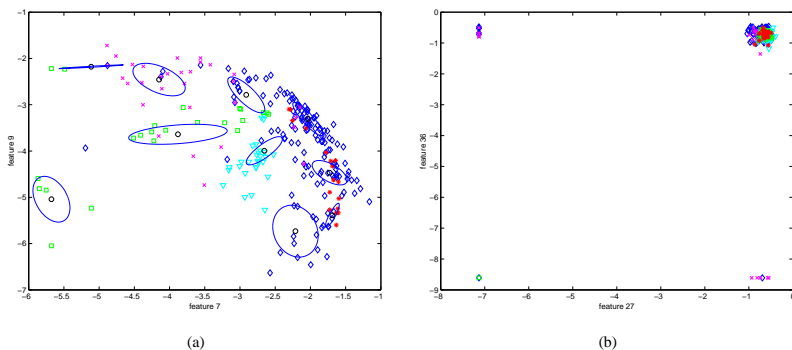


Figure 14: The scatter plots on the HRCT-lung data using the first two features chosen by FSSEM-k-TR (a) and FSSEM-k-ML (b). \times , \diamond , \square , $*$, and ∇ represent the different class assignments. \circ are the cluster means, and the ellipses are the covariances corresponding to the clusters discovered by FSSEM-k-TR.

HRCT-lung is a difficult data set due to its skewed class distribution (approximately 62.8% of the data is from the disease Centrilobular Emphysema). Because of this, even though EM-k discovered approximately only one cluster, its class error (which is equal to the error using a majority classification rule) is close to the values obtained by the other methods. The high dimensions obscure the HRCT-lung's classes and result in EM-k finding only one cluster. Even with a difficult problem such as this, feature selection obtained better CV-error than without feature selection using much fewer features (an average of 1.7 for FSSEM-k-TR and 2.9 for FSS-Kmeans-k-TR) compared to the original 108 features. FSSEM-k-TR picked features $\{7, 9\}$ and FSS-Kmeans-k-TR chose features $\{8, 6, 59\}$. Features 6, 7, 8, and 9 are gray level histogram values of the lung region, and feature 59 is a histogram value at a local pathology bearing region. These features make sense in discriminating

between Centrilobular Emphysema (the largest class) from the rest, as this class is characterized by low gray level values.

6.5.5 CONCLUSIONS ON EXPERIMENTS WITH REAL DATA

Our results on real data show that feature selection improved the performance of clustering algorithms in finding “interesting” patterns. We measure “interestingness” performance here by how well the discovered clusters match labeled classes (CV-error). FSSEM-k and FSS-Kmeans-k obtained better CV-error than EM-k and k -means using fewer features. Moreover, our experiments reveal that no one feature selection criterion (ML or TR) is better than the other. They have different biases. ML selects features that results in high-density clusters, and performed better than TR on the ionosphere data. Scatter separability (TR) prefers features that reveal well-separated clusters, and performed better than ML on the HRCT-lung data. They both did well on the iris and wine data.

7. Related Work: A Review of Feature Selection Algorithms for Unsupervised Learning

There are three different ways to select features from unsupervised data: 1) after clustering, 2) before clustering, and 3) during clustering. An example algorithm that performs feature selection after clustering is (Mirkin, 1999). The method first applies a new separate-and-conquer version of k -means clustering. Then, it computes the contribution weight of each variable in proportion to the squared deviation of each variable’s within-cluster mean from the total mean. It represents clusters by conjunctive concepts starting from the variable with the highest weight, until adding variables (with its conceptual description) does not improve the cluster “precision error”. Feature selection after clustering is important for conceptual learning, for describing and summarizing structure from data. This type of selecting features can remove redundancy but not feature irrelevance because the initial clustering is performed using all the features. As pointed out earlier, the existence of irrelevant features can misguide clustering results. Using all the features for clustering also assumes that our clustering algorithm does not break down with high dimensional data. In this paper, we only examine feature selection algorithms that affect (can change) the clustering outcomes; i.e., before or during clustering.

A significant body of research exists on methods for feature subset selection for supervised data. These methods can be grouped as *filter* (Marill and Green, 1963; Narendra and Fukunaga, 1977; Almuallim and Dietterich, 1991; Kira and Rendell, 1992; Kononenko, 1994; Liu and Setiono, 1996; Cardie, 1993; Singh and Provan, 1995) or *wrapper* (John et al., 1994; Doak, 1992; Caruana and Freitag, 1994; Aha and Bankert, 1994; Langley and Sage, 1994; Pazzani, 1995) approaches. To maintain the filter/wrapper model distinction used in supervised learning, we define *filter* methods in unsupervised learning as using some intrinsic property of the data to select features without utilizing the clustering algorithm that will ultimately be applied. *Wrapper* approaches, on the other hand, apply the unsupervised learning algorithm to each candidate feature subset and then evaluate the feature subset by criterion functions that utilize the clustering result.

When we first started this research, not much work has been done in feature subset selection for unsupervised learning in the context of machine learning, although research in the form of principal components analysis (PCA) (Chang, 1983), factor analysis (Johnson and Wichern, 1998) and projection pursuit (Friedman, 1987; Huber, 1985) existed. These early works in data reduction for unsupervised data can be thought of as filter methods, because they select the features prior to apply-

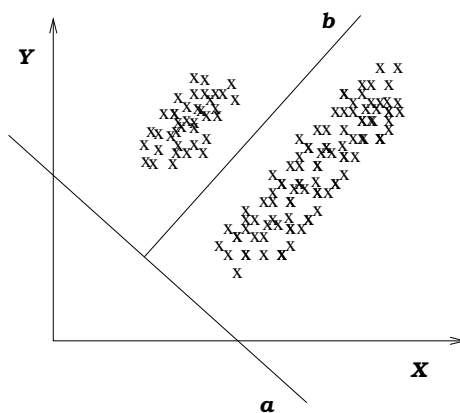


Figure 15: Illustration on when PCA is a poor discriminator.

ing clustering. But rather than selecting a subset of the features, they involve some type of feature transformation. PCA and factor analysis aim to reduce the dimension such that the representation is as faithful as possible to the original data. Note that data reduction techniques based on representation (like PCA) are better suited for compression applications rather than classification (Fukunaga (1990) provides an illustrative example on this). Figure 15 recreates this example. PCA chooses the projection with the highest variance. Projecting two dimensions to one dimension in this example, PCA would project the data to axis b , which is clearly inferior to axis a for discriminating the two clusters. Contrary to PCA and factor analysis, projection pursuit aims to find “interesting” projections from multi-dimensional data for visualizing structure in the data. A recent method for finding transformations called independent components analysis (ICA) (Hyvärinen, 1999) has gained widespread attention in signal processing. ICA tries to find a transformation such that the transformed variables are statistically independent.

The filter methods described in the previous paragraph all involve transformations of the original variable space. In this paper, we are interested in subsets of the original space, because some domains prefer the original variables in order to maintain the physical interpretation of these features. Moreover, transformations of the variable space require computation or collection of all the features before dimension reduction can be achieved, whereas subsets of the original space require computation or collection of only the selected feature subsets after feature selection is determined. If some features cost more than others, one can consider these costs in selecting features. In this paper, we assume each feature has equal cost. Other interesting and current directions in feature selection involving feature transformations are mixtures of principal component analyzers (Kambhatla and Leen, 1997; Tipping and Bishop, 1999) and mixtures of factor analyzers (Ghahramani and Beal, 2000; Ghahramani and Hinton, 1996; Ueda et al., 1999). We consider these mixture algorithms as wrapper approaches.

In recent years, more attention has been paid to unsupervised feature subset selection. Most of these methods are wrapper approaches. Gennari (1991) incorporates feature selection (they call “attention”) to CLASSIT (an incremental concept formation hierarchical clustering algorithm introduced in (Gennari et al., 1989)). The attention algorithm inspects the features starting with the most salient (“per-attribute contribution to category utility”) attribute to the least salient attribute,

and stop inspecting features if the remaining features do not change the current clustering decision. The purpose of this attention mechanism is to increase efficiency without loss of prediction accuracy. Devaney and Ram (1997) applied sequential forward and backward search. To evaluate each candidate subset, they measured the category utility of the clusters found by applying COBWEB (Fisher, 1987) in conjunction with the feature subset. Talavera (1999) applied “blind” (similar to the filter) and “feedback” (analogous to the wrapper) approaches to COBWEB, and used a feature dependence measure to select features. Vaithyanathan and Dom (1999) formulated an objective function for choosing the feature subset and finding the optimal number of clusters for a document clustering problem using a Bayesian statistical estimation framework. They modeled each cluster as a multinomial. They extended this concept to create hierarchical clusters (Vaithyanathan and Dom, 2000). Agrawal, et al. (1998) introduced a clustering algorithm (CLIQUE) which proceeds level-by-level from one feature to the highest dimension or until no more feature subspaces with clusters (regions with high density points) are generated. CLIQUE is a density based clustering algorithm which does not assume any density model. However, CLIQUE needs to specify parameters τ (the density threshold) and v (the equal length interval partitioning for each dimension). In contrast, our method makes assumptions about distributions to avoid specifying parameters. Kim, Street and Menczer (2002) apply an evolutionary local selection algorithm (ELSA) to search the feature subset and number of clusters on two clustering algorithms: K-means and EM clustering (with diagonal covariances), and a Pareto front to combine multiple objective evaluation functions. Law, Figueiredo and Jain (2002) estimate feature saliency using EM by modeling relevant features as conditionally independent given the component label, and irrelevant features with a probability density identical for all components. They also developed a wrapper approach that selects features using Kullback-Leibler divergence and entropy. Friedman and Meulman (2003) designed a distance measure for attribute-value data for clustering on subsets of attributes, and allow feature subsets for each cluster to be different.

8. Summary

In this paper, we introduced a wrapper framework for performing feature subset selection for unsupervised learning. We explored the issues involved in developing algorithms under this framework. We identified the need for finding the number of clusters in feature search and provided proofs for the biases of ML and scatter separability with respect to dimension. We, then, presented methods to ameliorate these problems.

Our experimental results showed that incorporating finding the number of clusters k into the feature subset selection process led to better results than fixing k to be the true number of classes. There are two reasons: 1) the number of classes is not necessarily equal to the number of Gaussian clusters, and 2) different feature subsets have different number of clusters. Supporting theory, our experiments on simulated data showed that ML and scatter separability are in some ways biased with respect to dimension. Thus, a normalization scheme is needed for the chosen feature selection criterion. Our proposed cross-projection criterion normalization scheme was able to eliminate these biases.

Although we examined the wrapper framework using FSSEM, the search method, feature selection criteria (especially the *trace* criterion), and the feature normalization scheme can be easily applied to any clustering method. The issues we have encountered and solutions presented are applicable to any feature subset wrapper approach. FSSEM serves as an example. Depending on one’s

application, one may choose to apply a more appropriate search method, clustering and feature selection criteria.

9. Future Directions

Research in feature subset selection for unsupervised learning is quite young. Even though we have addressed some issues, the paper opens up more questions that need to be answered.

Hartigan (1985) pointed out that no single criterion is best for all applications. This is reiterated by our results on the HRCT and Ionosphere data. This led us to work in visualization and user interaction to guide the feature search (Dy and Brodley, 2000b). Another interesting direction is to look at feature selection with hierarchical clustering (Gennari, 1991; Fisher, 1996; Devaney and Ram, 1997; Talavera, 1999; Vaithyanathan and Dom, 2000), since hierarchical clustering provides groupings at various perceptual levels. In addition, a cluster may be modeled better by a different feature subset from other clusters. One may wish to develop algorithms that select a different feature subset for each cluster component.

We explored unsupervised feature selection through the wrapper framework. It would be interesting to do a rigorous investigation of filter versus wrapper approach for unsupervised learning. One may also wish to venture in transformations of the original variable space. In particular, investigate on mixtures of principal component analyzers (Kambhatla and Leen, 1997; Tipping and Bishop, 1999), mixtures of factor analyzers (Ghahramani and Beal, 2000; Ghahramani and Hinton, 1996; Ueda et al., 1999) and mixtures of independent component analyzers (Hyvärinen, 1999).

The difficulty with unsupervised learning is the absence of labeled examples to guide the search. Breiman (Breiman, 2002) suggests transforming the clustering problem into a classification problem by assigning the unlabeled data to class one, and adding the same amount of random vectors into another class two. The second set is generated by independent sampling from the one-dimensional marginal distributions of class one. Understanding and developing tricks such as these to uncover structure from unlabeled data remains as topics that need further investigation. Another avenue for future work is to explore semi-supervised (few labeled examples and large amounts of unlabeled data) methods for feature selection.

Finally, in feature selection for unsupervised learning, several fundamental questions are still unanswered:

1. How do you define what “interestingness” means?
2. Should the criterion for “interestingness” (feature selection criterion) be the same as the criterion for “natural” grouping (clustering criterion)? Most of the literature uses the same criterion for feature selection and clustering as this leads to a clean optimization formulation. However, defining “interestingness” into a mathematical criterion is a difficult problem. Allowing different criteria to interact may provide a better model. Our experimental results on the wine data suggest this direction.
3. Our experiments on synthetic data indicate the need to standardize features. Mirkin, 1999, also standardized his features. Should features always be standardized before feature selection? If so, how do you standardize data containing different feature types (real-valued, nominal, and discrete)?

4. What is the best way to evaluate the results? In this paper, we evaluate performance using an external criterion (cross-validated class error). This is a standard measure used by most papers in the feature selection for unsupervised learning literature. Class error is task specific and measures the performance for one labeling solution. Is this the best way to compare different clustering algorithms?

Acknowledgments

The authors wish to thank Dr. Craig Codrington for discussions on criterion normalization, the ML-lunch group at Purdue University for helpful comments, and the reviewers for their constructive remarks. This research is supported by NSF Grant No. IRI9711535, and NIH Grant No. 1 R01 LM06543-01A1.

Appendix A. EM Clustering

Clustering using finite mixture models is a well-known method and has been used for a long time in pattern recognition Duda and Hart (1973); Fukunaga (1990); Jain and Dubes (1988) and statistics McLachlan and Basford (1988); Titterington et al. (1985); Fraley and Raftery (2000). In this model, one assumes that the data is generated from a mixture of component density functions, in which each component density function represents a cluster. The probability distribution function of the data has the following form:

$$f(X_i|\Phi) = \sum_{j=1}^k \pi_j f_j(X_i|\theta_j) \quad (5)$$

where $f_j(X_i|\theta_j)$ is the probability density function for class j , π_j is the mixing proportion of class j (prior probability of class j), k is the number of clusters, X_i is a d -dimensional random data vector, θ_j is the set of parameters for class j , $\Phi = (\pi, \theta)$ is the set of all parameters and $f(X_i|\Phi)$ is the probability density function of our observed data point X_i given the parameters Φ . Since the π_j 's are prior probabilities, they are subject to the following constraints: $\pi_j \geq 0$ and $\sum_{j=1}^k \pi_j = 1$.

The X_i 's, where $i = 1 \dots N$, are the data vectors we are trying to cluster, and N is the number of samples. To cluster X_i , we need to estimate the parameters, Φ . One method for estimating Φ is to find Φ that maximizes the log-likelihood, $\log f(X|\Phi) = \sum_{i=1}^N \log f(X_i|\Phi)$. To compute $f(X_i|\Phi)$, we need to know the cluster (the missing data) to which X_i (the observed data) belongs. We apply the EM algorithm, which provides us with "soft-clustering" information; i.e., a data point X_i can belong to more than one cluster (weighted by its probability to belong to each cluster). The expectation-maximization (EM) algorithm, introduced in some generality by Dempster, Laird and Rubin in 1977, is an iterative approximation algorithm for computing the maximum likelihood (ML) estimate of missing data problems.

Going through the derivation of applying EM on our Gaussian mixture model, we obtain the following EM update equations (Wolfe, 1970):

$$E[z_{ij}]^{(t)} = p(z_{ij} = 1 | X, \Phi^{(t)}) = \frac{f_j(X_i|\Phi_j^{(t)})\pi_j^{(t)}}{\sum_{s=1}^k f_s(X_i|\Phi_s^{(t)})\pi_s^{(t)}}; \quad (6)$$

$$\pi_j^{(t+1)} = \frac{1}{N} \sum_{i=1}^N E[z_{ij}]^{(t)}; \quad (7)$$

$$\mu_j^{(t+1)} = \frac{1}{N\pi_j^{(t+1)}} \sum_{i=1}^N E[z_{ij}]^{(t)} X_i; \quad (8)$$

$$\Sigma_j^{(t+1)} = \frac{1}{N\pi_j^{(t+1)}} \sum_{i=1}^N E[z_{ij}]^{(t)} (X_i - \mu_j^{(t+1)})(X_i - \mu_j^{(t+1)})^T; \quad (9)$$

where $E[z_{ij}]$ is the probability that X_i belongs to cluster j given our current parameters and X_i , $\sum_{i=1}^N E[z_{ij}]$ is the estimated number of data points in class j , and the superscript t refers to the iteration.

Appendix B. Additional Proofs on ML's Bias with Dimension

In this appendix, we prove Theorem 4.1 and Corollary 4.1 which state the condition that needs to be satisfied for the maximum likelihood of feature subset A , $ML(\Phi_A)$, to be greater than or equal to the maximum likelihood of feature subset B , $ML(\Phi_B)$. To prove these results, we first define the maximum likelihood criterion for a mixture of Gaussians, prove Lemma B.1 which derives a simplified form of $\exp(Q(\Phi, \Phi))$ for a finite Gaussian mixture, and Lemma B.2 which states the condition that needs to be satisfied for the complete expected data log-likelihood $Q(\cdot)$ function given the observed data and the parameter estimates in feature subset A , $Q(\Phi_A, \Phi_A)$, to be greater than or equal to the $Q(\cdot)$ function of feature subset B , $Q(\Phi_B, \Phi_B)$.

The maximum likelihood of our data, X , is

$$ML = \max_{\Phi} (f(X|\Phi)) = \max_{\Phi} \prod_{i=1}^N \left(\sum_{j=1}^k \pi_j f_j(X_i|\theta_j) \right), \quad (10)$$

where $f_j(X_i|\theta_j)$ is the probability density function for class j , π_j is the mixing proportion of class j (prior probability of class j), N is the number of data points, k is the number of clusters, X_i is a d -dimensional random data vector, θ_j is the set of parameters for class j , $\Phi = (\pi, \theta)$ is the set of all parameters and $f(X|\Phi)$ is the probability density function of our observed data $X = X_1, X_2, \dots, X_N$ given the parameters Φ . We choose the feature subset that maximizes this criterion.

Lemma B.1 *For a finite mixture of Gaussians,*

$$\exp(Q(\Phi, \Phi)) = \prod_{j=1}^K \pi_j^{N\pi_j} \frac{1}{(2\pi)^{\frac{dN\pi_j}{2}} |\Sigma_j|^{\frac{N\pi_j}{2}}} e^{-\frac{1}{2}dN\pi_j},$$

where x_i , $i = 1 \dots N$, are the N observed data points, z_{ij} is the missing variable equal to one if x_i belongs to cluster j and zero otherwise, π_j is the mixture proportion, μ_j is the mean and Σ_j is the covariance matrix of each Gaussian cluster respectively, and $\Phi = (\pi, \mu, \Sigma)$ is the set of all estimated parameters.

Proof:

$$\begin{aligned}
 Q(\Phi, \Phi) &\triangleq E_{z|x}[\log f(x, z|\Phi)|x, \Phi] \\
 &= E_{z|x}[\log f(x|z, \Phi)|x, \Phi] + E_{z|x}[\log f(z|\Phi)|x, \Phi] \\
 &= \sum_{i=1}^N \sum_{j=1}^K p(z_{ij} = 1|x, \Phi) \log f_j(x_i|\phi_j) + \sum_{i=1}^N \sum_{j=1}^K p(z_{ij} = 1|x, \Phi) \log \pi_j \\
 &= \sum_{i=1}^N \sum_{j=1}^K p(z_{ij} = 1|x, \Phi) \log(\pi_j f_j(x_i|\phi_j)) \tag{11}
 \end{aligned}$$

$$\begin{aligned}
 &= \sum_{i=1}^N \sum_{j=1}^K E[z_{ij}] \log(\pi_j f_j(x_i|\phi_j)) \\
 \exp(Q(\Phi, \Phi)) &= \prod_{i=1}^N \prod_{j=1}^K (\pi_j f_j(x_i|\phi_j))^{E[z_{ij}]} \tag{12}
 \end{aligned}$$

Substituting our parameter estimates to Equation 12 and sample data x_i 's,

$$\begin{aligned}
 \exp(Q(\Phi, \Phi)) &= \prod_{j=1}^K \pi_j^{\sum_{i=1}^N E[z_{ij}]} \prod_{i=1}^N \left(\frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_j|^{\frac{1}{2}}} e^{-\frac{1}{2}(x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j)} \right)^{E[z_{ij}]} \\
 &= \prod_{j=1}^K \pi_j^{N\pi_j} \frac{1}{(2\pi)^{\frac{dN\pi_j}{2}} |\Sigma_j|^{\frac{N\pi_j}{2}}} e^{-\frac{1}{2} \sum_{i=1}^N E[z_{ij}] (x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j)} \tag{13}
 \end{aligned}$$

Simplifying the exponent of e we obtain

$$\begin{aligned}
 &-\frac{1}{2} \sum_{i=1}^N E[z_{ij}] (x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j) \\
 &= -\frac{1}{2} \sum_{i=1}^N E[z_{ij}] \text{tr}((x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j)) \\
 &= -\frac{1}{2} \sum_{i=1}^N E[z_{ij}] \text{tr}(\Sigma_j^{-1} (x_i - \mu_j) (x_i - \mu_j)^T) \\
 &= -\frac{1}{2} \text{tr}(\Sigma_j^{-1} (\sum_{i=1}^N E[z_{ij}] (x_i - \mu_j) (x_i - \mu_j)^T)).
 \end{aligned}$$

Adding and subtracting \bar{x}_j , where $\bar{x}_j = \frac{1}{N\pi_j} \sum_{i=1}^N E[z_{ij}] x_i$, this last expression becomes

$$-\frac{1}{2} \text{tr}(\Sigma_j^{-1} (\sum_{i=1}^N E[z_{ij}] (x_i - \bar{x}_j + \bar{x}_j - \mu_j) (x_i - \bar{x}_j + \bar{x}_j - \mu_j)^T)).$$

Cancelling cross-product terms yields

$$-\frac{1}{2} \text{tr}(\Sigma_j^{-1} (\sum_{i=1}^N E[z_{ij}] (x_i - \bar{x}_j) (x_i - \bar{x}_j)^T + \sum_{i=1}^N E[z_{ij}] (\bar{x}_j - \mu_j) (\bar{x}_j - \mu_j)^T)),$$

and finally substituting the parameter estimates (Equations 7-9) gives the expression

$$\begin{aligned} & -\frac{1}{2}tr(\Sigma_j^{-1}\Sigma_jN\pi_j) \\ & = -\frac{1}{2}dN\pi_j. \end{aligned} \tag{14}$$

Thus, $\exp(Q(\Phi, \Phi))$ can be expressed as

$$\exp(Q(\Phi, \Phi)) = \prod_{j=1}^K \pi_j^{N\pi_j} \frac{1}{(2\pi)^{\frac{dN\pi_j}{2}} |\Sigma_j|^{\frac{N\pi_j}{2}}} e^{-\frac{1}{2}dN\pi_j} \quad \square$$

Lemma B.2 *Assuming identical clustering assignments for feature subsets A and B with dimensions $d_B \geq d_A$, $Q(\Phi_A, \Phi_A) \geq Q(\Phi_B, \Phi_B)$ iff*

$$\prod_{j=1}^k \left(\frac{|\Sigma_B|_j}{|\Sigma_A|_j} \right)^{\pi_j} \geq \frac{1}{(2\pi e)^{(d_B-d_A)}}.$$

Proof:

Applying Lemma B.1, and assuming subsets A and B have equal clustering assignments,

$$\begin{aligned} & \frac{\exp(Q(\Phi_A, \Phi_A))}{\exp(Q(\Phi_B, \Phi_B))} \geq 1; \\ & \frac{\prod_{j=1}^k \pi_j^{N\pi_j} \frac{1}{(2\pi e)^{\frac{d_A N\pi_j}{2}}} \frac{1}{|\Sigma_A|_j^{\frac{N\pi_j}{2}}}}{\prod_{j=1}^k \pi_j^{N\pi_j} \frac{1}{(2\pi e)^{\frac{(d_B)N\pi_j}{2}}} \frac{1}{|\Sigma_B|_j^{\frac{N\pi_j}{2}}}} \geq 1. \end{aligned} \tag{15}$$

Given $d_B \geq d_A$, without loss of generality and cancelling common terms,

$$\begin{aligned} & \prod_{j=1}^k \left(\frac{|\Sigma_B|_j}{|\Sigma_A|_j} \right)^{\frac{N\pi_j}{2}} (2\pi e)^{\frac{(d_B-d_A)N\pi_j}{2}} \geq 1; \\ & \prod_{j=1}^k \left(\frac{|\Sigma_B|_j}{|\Sigma_A|_j} \right)^{\pi_j} (2\pi e)^{(d_B-d_A)\pi_j} \geq 1; \\ & (2\pi e)^{(d_B-d_A)\sum_{j=1}^k \pi_j} \prod_{j=1}^k \left(\frac{|\Sigma_B|_j}{|\Sigma_A|_j} \right)^{\pi_j} \geq 1; \\ & \prod_{j=1}^k \left(\frac{|\Sigma_B|_j}{|\Sigma_A|_j} \right)^{\pi_j} \geq \frac{1}{(2\pi e)^{(d_B-d_A)}}. \quad \square \end{aligned}$$

Theorem B.1 (Theorem 4.1 restated) *For a finite multivariate Gaussian mixture, assuming identical clustering assignments for feature subsets A and B with dimensions $d_B \geq d_A$, $ML(\Phi_A) \geq ML(\Phi_B)$ iff*

$$\prod_{j=1}^k \left(\frac{|\Sigma_B|_j}{|\Sigma_A|_j} \right)^{\pi_j} \geq \frac{1}{(2\pi e)^{(d_B-d_A)}}.$$

Proof:

The log-likelihood, $\log L(\Phi') = \log f(x|\Phi')$.

$$\begin{aligned} \log L(\Phi') &= E_{z|x}[\log f(x, z|\Phi')|x, \Phi] - E_{z|x}[\log f(z|x, \Phi')|x, \Phi] \\ &\stackrel{\triangle}{=} Q(\Phi', \Phi) + H(\Phi', \Phi). \end{aligned}$$

$$\begin{aligned} H(\Phi, \Phi) &= -E[\log f(z|x, \Phi)|x, \Phi] \\ &= -\sum_{i=1}^N \sum_{j=1}^k p(z_{ij} = 1|x_i, \phi_j) \log p(z_{ij} = 1|x_i, \phi_j) \\ &= -\sum_{i=1}^N \sum_{j=1}^k E[z_{ij}] \log E[z_{ij}]. \end{aligned} \tag{16}$$

Since the identical clustering assignment assumption means that $E[z_{ij}]$ for feature set A is equal to $E[z_{ij}]$ for feature set B ,

$$H(\Phi_A, \Phi_A) = H(\Phi_B, \Phi_B).$$

Thus,

$$\frac{ML(\Phi_A)}{ML(\Phi_B)} = \frac{\exp(Q(\Phi_A, \Phi_A))}{\exp(Q(\Phi_B, \Phi_B))}.$$

For a finite Gaussian mixture, from Lemma B.2, $\frac{ML(\Phi_A)}{ML(\Phi_B)} \geq 1$ iff

$$\prod_{j=1}^k \left(\frac{|\Sigma_B|_j}{|\Sigma_A|_j} \right)^{\pi_j} \geq \frac{1}{(2\pi e)^{(d_B - d_A)}}. \quad \square$$

Corollary B.1 (Corollary 4.1 restated) *For a finite multivariate Gaussian mixture, assuming identical clustering assignments for feature subsets X and (X, Y) , where X and Y are disjoint, $ML(\Phi_X) \geq ML(\Phi_{XY})$ iff*

$$\prod_{j=1}^k |\Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}|_j^{\pi_j} \geq \frac{1}{(2\pi e)^{d_Y}}.$$

Proof:

Applying Theorem 4.1, and if we let A be the marginal feature vector X with dimension d_X and B be the joint feature vector (X, Y) with dimension $d_X + d_Y$ (where subsets X and Y are disjoint), then the maximum likelihood of X is greater than or equal to the maximum likelihood of (X, Y) iff

$$\begin{aligned} \frac{ML(\Phi_X)}{ML(\Phi_{XY})} &\geq 1 \\ (2\pi e)^{d_Y} \prod_{j=1}^k \left(\frac{\begin{vmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{vmatrix}_j}{|\Sigma_{XX}|_j} \right)^{\pi_j} &\geq 1. \end{aligned}$$

Exercise 4.11 of Johnson and Wichern (1998) shows that for any square matrix A ,

$$\begin{aligned}
A &= \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \\
|A| &= |A_{22}| |A_{11} - A_{12} A_{22}^{-1} A_{21}| \quad \text{for } |A_{22}| \neq 0 \\
&= |A_{11}| |A_{22} - A_{21} A_{11}^{-1} A_{12}| \quad \text{for } |A_{11}| \neq 0.
\end{aligned}$$

Thus,

$$\begin{aligned}
(2\pi e)^{d_Y} \prod_{j=1}^k \left(\frac{1}{|\Sigma_{XX}|_j} |\Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}|_j \right)^{\pi_j} &\geq 1 \\
\prod_{j=1}^k |\Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}|_j^{\pi_j} &\geq \frac{1}{(2\pi e)^{d_Y}}. \quad \square
\end{aligned}$$

Now, what do these results mean? One can compute the maximum log-likelihood, log ML efficiently as $Q(\Phi, \Phi) + H(\Phi, \Phi)$ by applying Lemma B.1 and Equation 16. Lemma B.1 shows that the ML criterion prefers low covariance clusters. Equation 16 shows that the ML criterion penalizes increase in cluster entropy. Theorem 4.1 and Corollary 4.1 reveal the dependencies of comparing the ML criterion for different dimensions. Note that the left hand side term of Corollary 4.1 is the determinant of the covariance of $f(Y|X)$. It is the covariance of Y minus the correlation of Y and X . For a criterion measure to be unbiased with respect to dimension, the criterion value should be the same for the different subsets when the cluster assignments are equal (and should not be dependent on the dimension). But, in this case, ML increases when the additional feature has small variance and decreases when the additional feature has large variance.

References

- R. Agrawal, J. Gehrke, D. Gunopulos, and P. Raghavan. Automatic subspace clustering of high dimensional data for data mining applications. In *Proceedings ACM SIGMOD International Conference on Management of Data*, pages 94–105, Seattle, WA, June 1998. ACM Press.
- D. W. Aha and R. L. Bankert. Feature selection for case-based classification of cloud types: An empirical comparison. In *Proceedings of the 1994 AAAI Workshop on Case-Based Reasoning*, pages 106–112, Seattle, WA, 1994. AAAI Press.
- H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, AC-19(6):716–723, December 1974.
- H. Almuallim and T. G. Dietterich. Learning with many irrelevant features. In *Proceedings of the Ninth National Conference on Artificial Intelligence*, pages 547–552, Anaheim, CA, 1991. AAAI Press.
- C. L. Blake and C. J. Merz. UCI repository of machine learning databases. In <http://www.ics.uci.edu/~mllearn/MLRepository.html>, 1998.

- C. A. Bouman, M. Shapiro, G. W. Cook, C. B. Atkins, and H. Cheng. Cluster: An unsupervised algorithm for modeling gaussian mixtures. In <http://dynamo.ecn.purdue.edu/~bouman/software/cluster>, October 1998.
- P. S. Bradley and U. M. Fayyad. Refining initial points for K-Means clustering. In *Proceedings of the Fifteenth International Conference on Machine Learning*, pages 91–99, San Francisco, CA, 1998. Morgan Kaufmann.
- L. Breiman. Wald lecture II: Looking inside the black box, 2002. 277th meeting of the Institute of Mathematical Statistics, Banff, Alberta, Canada. <http://stat-www.berkeley.edu/users/breiman/wald2002-2.pdf>.
- C. Cardie. Using decision trees to improve case-based learning. In *Machine Learning: Proceedings of the Tenth International Conference*, pages 25–32, Amherst, MA, 1993. Morgan Kaufmann.
- R. Caruana and D. Freitag. Greedy attribute selection. In *Proceedings of the Eleventh International Conference on Machine Learning*, pages 28–36, New Brunswick, NJ, 1994. Morgan Kaufmann.
- W. C. Chang. On using principal components before separating a mixture of two multivariate normal distributions. *Applied Statistics*, 32:267–275, 1983.
- P. Cheeseman and J. Stutz. Bayesian classification (autoclass): theory and results. In *Advances in Knowledge Discovery and Data Mining*, pages 153–180, Cambridge, MA, 1996. AAAI/MIT Press.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal Royal Statistical Society, Series B*, 39(1):1–38, 1977.
- M. Devaney and A. Ram. Efficient feature selection in conceptual clustering. In *Proceedings of the Fourteenth International Conference on Machine Learning*, pages 92–97, Nashville, TN, 1997. Morgan Kaufmann.
- J. Doak. An evaluation of feature selection methods and their application to computer security. Technical report, University of California at Davis, 1992.
- R. O. Duda and P. E. Hart. *Pattern Classification and Scene Analysis*. Wiley & Sons, NY, 1973.
- R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern Classification (Second Edition)*. Wiley & Sons, NY, 2001.
- J. G. Dy and C. E. Brodley. Feature subset selection and order identification for unsupervised learning. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pages 247–254, Stanford University, 2000a. Morgan Kaufmann.
- J. G. Dy and C. E. Brodley. Interactive visualization and feature selection for unsupervised data. In *Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 360–364, Boston, MA, August 2000b. ACM Press.
- J. G. Dy and C. E. Brodley. TR-CDSP-03-53: feature selection for unsupervised learning. Technical report, Northeastern University, December 2003.

- J. G. Dy, C. E. Brodley, A. Kak, L. S. Broderick, and A. M. Aisen. Unsupervised feature selection applied to content-based retrieval of lung images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(3):373–378, March 2003.
- J. G. Dy, C. E. Brodley, A. Kak, C. R. Shyu, and L. S. Broderick. The customized-queries approach to CBIR using EM. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 2, pages 400–406, Fort Collins, CO, June 1999. IEEE Computer Society Press.
- U. Fayyad, C. Reina, and P. S. Bradley. Initialization of iterative refinement clustering algorithms. In *Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining*, pages 194–198, New York, August 1998. AAAI Press.
- M. Figueiredo and A. K. Jain. Unsupervised learning of finite mixture models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24:381–396, 2002.
- D. Fisher. Iterative optimization and simplification of hierarchical clusterings. *Journal of Artificial Intelligence Research*, 4:147–179, 1996.
- D. H. Fisher. Knowledge acquisition via incremental conceptual clustering. *Machine Learning*, 2(2):139–172, 1987.
- E. Forgy. Cluster analysis of multivariate data: Efficiency vs. interpretability of classifications. *Biometrics*, 21:768, 1965.
- C. Fraley and A. E. Raftery. Model-based clustering, discriminant analysis, and density estimation. Technical Report Technical Report No. 380, University of Washington, Seattle, WA, 2000.
- J. H. Friedman. Exploratory projection pursuit. *Journal American Statistical Association*, 82:249–266, 1987.
- J. H. Friedman and J. J. Meulman. Clustering objects on subsets of attributes. *to appear at Journal Royal Statistical Society*, 2003.
- K. Fukunaga. *Statistical Pattern Recognition (second edition)*. Academic Press, San Diego, CA, 1990.
- J. H. Gennari. Concept formation and attention. In *Proceedings of the Thirteenth Annual Conference of the Cognitive Science Society*, pages 724–728, Chicago, IL, 1991. Lawrence Erlbaum.
- J. H. Gennari, P. Langley, and D. Fisher. Models of incremental concept formation. *Artificial Intelligence*, 40:11–61, 1989.
- Z. Ghahramani and M. J. Beal. Variational inference for bayesian mixtures of factor analysers. In S. A. Solla, T. K. Leen, and K. Muller, editors, *Advances in Neural Information Processing Systems*, volume 12, pages 449–455. MIT Press, 2000.
- Z. Ghahramani and G. E. Hinton. The EM algorithm for mixtures of factor analyzers. Technical Report Technical Report CRG-TR-96-1, University of Toronto, Toronto, Canada, 1996.
- J. A. Hartigan. Statistical theory in clustering. *Journal of Classification*, 2:63–76, 1985.

- P. J. Huber. Projection pursuit. *The Annals of Statistics*, 13(2):435–475, 1985.
- A. Hyvärinen. Survey on independent component analysis. *Neural Computing Surveys*, 2:94–128, 1999.
- A. K. Jain and R. C. Dubes. *Algorithms for clustering data*. Prentice Hall, Englewood Cliffs, NJ, 1988.
- G. H. John, R. Kohavi, and K. Pfleger. Irrelevant features and the subset selection problem. In *Machine Learning: Proceedings of the Eleventh International Conference*, pages 121–129, New Brunswick, NJ, 1994. Morgan Kaufmann.
- R. A. Johnson and D. W. Wichern. *Applied Multivariate Statistical Analysis*. Prentice-Hall, 4 edition, 1998.
- N. Kambhatla and T. K. Leen. Dimension reduction by local principal component analysis. *Neural Computation*, 9:1493–1516, 1997.
- Y. S. Kim, N. Street, and F. Menczer. Evolutionary model selection in unsupervised learning. *Intelligent Data Analysis*, 6:531–556, 2002.
- K. Kira and L. A. Rendell. A practical approach to feature selection. In *Proceedings of the Ninth International Workshop on Machine Learning*, pages 249–256, Aberdeen, Scotland, 1992. Morgan Kaufmann.
- J. Kittler. Feature set search algorithms. In *Pattern Recognition and Signal Processing*, pages 41–60, 1978.
- R. Kohavi and G. H. John. Wrappers for feature subset selection. *Artificial Intelligence*, 97(1-2): 273–324, 1997.
- I. Kononenko. Estimating attributes: analysis and extensions of relief. In *Proceedings of the European Conference on Machine Learning*, 1994.
- P. Langley and S. Sage. Oblivious decision trees and abstract cases. In *Working Notes of the AAAI-94 Workshop on Case-Based Reasoning*, pages 113–117, Seattle, 1994. AAAI Press.
- M. H. Law, M. Figueiredo, and A. K. Jain. Feature selection in mixture-based clustering. In *Advances in Neural Information Processing Systems 15*, Vancouver, December 2002.
- H. Liu and R. Setiono. Dimensionality reduction via discretization. *Knowledge-Based Systems*, 9(1):67–72, February 1996.
- T. Marill and D. M. Green. On the effectiveness of receptors in recognition systems. *IEEE Transactions on Information Theory*, 9:11–17, 1963.
- G. McLachlan and T. Krishnan. *The EM Algorithm and Extensions*. John Wiley, New York, 1997.
- G. J. McLachlan and K. E. Basford. *Mixture Models, Inference and Applications to Clustering*. Marcel Dekker, New York, 1988.

- G. Milligan. A monte carlo study of thirty internal criterion measures for cluster analysis. *Psychometrika*, 46(2):187–199, June 1981.
- B. Mirkin. Concept learning and feature selection based on square-error clustering. *Machine Learning*, 35:25–39, 1999.
- T. K. Moon. The expectation-maximization algorithm. In *IEEE Signal Processing Magazine*, pages 47–59, November 1996.
- P. M. Narendra and K. Fukunaga. A branch and bound algorithm for feature subset selection. *IEEE Transactions on Computers*, C-26(9):917–922, September 1977.
- M. J. Pazzani. Searching for dependencies in bayesian classifiers. In *Proceedings of the fifth International Workshop on Artificial Intelligence and Statistics*, Ft. Lauderdale, Florida, 1995.
- D. Pelleg and A. Moore. X-means: Extending K -means with efficient estimation of the number of clusters. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pages 727–734. Morgan Kaufmann, San Francisco, CA, 2000.
- J. Rissanen. A universal prior for integers and estimation by minimum description length. *Annals of Statistics*, 11(2):416–431, 1983.
- S. J. Russell and P. Norvig. *Artificial Intelligence a Modern Approach*. Prentice Hall, Saddle River, NJ, 1995.
- G. Salton and M. J. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill Book Company, 1983.
- G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464, 1978.
- M. Singh and G. M. Provan. A comparison of induction algorithms for selective and non-selective bayesian classifiers. In *Machine Learning: Proceedings of the Twelfth International Conference*, pages 497–505, 1995.
- P. Smyth. Clustering using Monte Carlo cross-validation. In E. Simoudis, J. Han, and U. Fayyad, editors, *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, pages 126–133, Portland, OR, 1996. AAAI Press.
- L. Talavera. Feature selection as a preprocessing step for hierarchical clustering. In *Proceedings of the Sixteenth International Conference on Machine Learning*, pages 389–397, Bled, Slovenia, 1999. Morgan Kaufmann.
- M. E. Tipping and C. M. Bishop. Mixtures of probabilistic principal component analysers. *Neural Computation*, 11(2):443–482, 1999.
- D. M. Titterton, A. F. M. Smith, and U. E. Makov. *Statistical Analysis of Finite Mixture Distributions*. John Wiley and Sons, Chichester, UK, 1985.
- N. Ueda, R. Nakano, Z. Ghahramani, and G. E. Hinton. SMEM algorithm for mixture models. In *Neural Computation*, 1999. To appear.

- S. Vaithyanathan and B. Dom. Model selection in unsupervised learning with applications to document clustering. In *Proceedings of the Sixteenth International Conference on Machine Learning*, pages 433–443, Bled, Slovenia, 1999. Morgan Kaufmann.
- S. Vaithyanathan and B. Dom. Hierarchical unsupervised learning. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pages 1039–1046, Palo Alto, CA, 2000. Morgan Kaufmann.
- J. H. Wolfe. Pattern clustering by multivariate mixture analysis. *Multivariate Behavioral Research*, 5(3):101–116, 1970.
- C. F. J. Wu. On the convergence properties of the em algorithm. *The Annals of Statistics*, 11(1): 95–103, 1983.