

Hunting for Coherent Co-clusters in High Dimensional and Noisy Datasets

Meghana Deodhar, Joydeep Ghosh, Gunjan Gupta
 Department of ECE
 University of Texas at Austin
 Austin, TX, USA
 deodhar,ghosh,ggupta@ece.utexas.edu

Hyuk Cho, Inderjit Dhillon
 Department of CS
 University of Texas at Austin
 Austin, TX, USA
 hyukcho,inderjit@cs.utexas.edu

Abstract

Clustering problems often involve datasets where only a part of the data is relevant to the problem, e.g., in microarray data analysis only a subset of the genes show cohesive expressions within a subset of the conditions/features. The existence of a large number of non-informative data points and features makes it challenging to hunt for coherent and meaningful clusters from such datasets. Additionally, since clusters could exist in different subspaces of the feature space, a co-clustering algorithm that simultaneously clusters objects and features is often more suitable as compared to one that is restricted to traditional “one-sided” clustering. We propose Robust Overlapping Co-clustering (ROCC), a scalable and very versatile framework that addresses the problem of efficiently mining dense, arbitrarily positioned, possibly overlapping co-clusters from large, noisy datasets. ROCC has several desirable properties that make it extremely well suited to a number of real life applications. Through extensive experimentation we show that our approach is significantly more accurate in identifying biologically meaningful co-clusters in microarray data as compared to several other prominent approaches that have been applied to this task. We also point out other interesting applications of the proposed framework in solving difficult clustering problems.

1. Motivation

When clustering certain real world datasets, it has been observed that only a part of the data forms cohesive clusters. For example, in the case of microarray data, typically only a small subset of the genes cluster well and the rest can be considered non-informative [16]. Problems addressed by eCommerce businesses, such as market basket analysis and fraud detection involve huge, noisy datasets with coherent patterns occurring only in small pockets of the data. Moreover, for such data, coherent clusters could be arbitrarily positioned in subspaces formed by different, possibly over-

lapping subsets of features, e.g., different subsets of genes may be correlated across different subsets of experiments in microarray data. Additionally, it is possible that some features may not be relevant to any cluster.

Traditional clustering algorithms like k -means do not address either issue since they assign every data point to a cluster based on a similarity measure computed across all the features. Feature selection or feature clustering [9, 11] improve clustering results on high dimensional and noisy datasets, but do not allow clusters existing in different subsets of the feature space to be detected easily. Co-clustering simultaneously clusters the data along multiple axes, e.g., in the case of microarray data it simultaneously clusters the genes as well as the experiments [6] and can hence detect clusters existing in different subspaces of the feature space. In this paper we focus on real life datasets, where co-clusters are arbitrarily positioned in the data matrix, could be overlapping and are obfuscated by the presence of a large number of irrelevant points. Our goal is to discover dense, arbitrarily positioned and overlapping co-clusters in the data, while simultaneously pruning away non-informative objects and features.

2. Related Work

Density based clustering algorithms have a motivation similar to our proposed approach and use the notion of local density to cluster only a relevant subset of the data into multiple dense clusters. DBSCAN [12] and its improved versions such as OPTICS rely on the notion of density to find arbitrarily shaped clusters in large spatial databases in the presence of noise. These approaches however, are not scalable to high dimensional datasets and are limited to Euclidean or related distance measures. The One Class Information Bottleneck algorithm [8] and the Batch Ball One Class Clustering algorithm [15] are efficient and scalable algorithms that can work with a large class of distance measures. However, both algorithms find only a single dense region. The Bregman Bubble Clustering (BBC)

technique [16] addresses the problem of discovering multiple, dense regions in a dataset while discarding the relatively non-coherent parts of the data. BBC provides a robust, scalable framework for clustering only a relevant fraction of the data. However, all of these approaches are developed for one-sided clustering only, where the data points are clustered based on their similarity across the entire set of features.

In contrast, both co-clustering (biclustering) and subspace clustering approaches locate clusters in subspaces of the feature space. The literature in both areas is recent but explosive, so we refer to the surveys and comparative studies in [20, 22, 24] as good starting points. As we shall see in Section 3, none of the existing methods provide the full set of capabilities that the proposed method provides.

Co-clustering was first applied to gene expression data by Cheng and Church [6], who used a greedy search heuristic to generate arbitrarily positioned, overlapping co-clusters, based on a homogeneity constraint. However, their iterative insertion and deletion based algorithm is expensive, since it identifies individual co-clusters sequentially rather than all at once. The algorithm also causes random perturbations to the data while masking discovered biclusters, which reduces the clustering quality. The plaid model approach [18] improves upon this by directly modeling overlapping clusters, but still cannot identify multiple co-clusters simultaneously. These algorithms are not very general as they assume additive Gaussian noise models. Neither can they effectively handle missing data. The xMotif algorithm [21] is an iterative search method that identifies submatrices with genes that have near constant expression levels across a subset of experiments. Therefore, most generated biclusters are not very interesting from the biological point of view.

In addition to the greedy, iterative algorithms discussed above, deterministic algorithms such as BiMax [24] and OPSM [3] have also been proposed. The BiMax approach, proposed by Prelic et al. is based on a simple, binary data model. Accordingly, the expression data has to be discretized before application of the BiMax algorithm, which could cause loss of useful information. The algorithm exactly finds all the inclusion maximal biclusters consisting of genes that jointly respond across a subset of experiments. The number of co-clusters identified by BiMax is exponential in the number of genes and experiments, making it infeasible to store and evaluate all the co-clusters in case of large datasets. The order preserving sub matrix algorithm (OPSM) looks for submatrices in which the expression levels of all the genes induce the same linear ordering of the experiments. This algorithm although very accurate, is designed to identify only a single co-cluster. A recent extension to OPSM [30] finds multiple, overlapping co-clusters in noisy datasets, but is very expensive in the number of

features.

Bregman Co-clustering (BCC), proposed by Banerjee et al. [1], is a highly efficient, generalized framework for partitional co-clustering [20] that works with any distance measure that is a Bregman divergence, or equivalently any noise distribution from the regular exponential family. BCC includes several previously developed co-clustering algorithms like information theoretic co-clustering [10] and minimum sum squared co-clustering [7] as special cases. The BCC framework is however restricted to grid-based, partitional co-clustering and assigns every point in the data matrix to exactly one co-cluster i.e., the co-clustering is exhaustive and exclusive.

Parsons et al. [22] present a survey of subspace clustering algorithms, which includes bottom-up grid based methods like CLIQUE and iterative top-down algorithms like PROCLUS. However, most of them are computationally intensive, need extensive tuning to get meaningful results and identify uniform clusters with very similar values rather than clusters with coherent trends or patterns. The pCluster model proposed by Wang et al. [26] generalizes subspace clustering and aims at discovering scaling or shifting patterns in clusters. Yoon et al. [29] improve upon this model and propose a biclustering algorithm that uses specialized data structures such as zero-suppressed binary decision diagrams to improve efficiency. Their algorithm is however very complex and is exponential in the number of features in the worst case. The more recent reg-cluster model [28] is designed to identify arbitrary scaling and shifting co-regulations patterns along with negative correlations. However, unlike our proposed approach, these pattern based, heuristic approaches do not use a principled cost function and do not scale well due to high complexity in the number of features.

3. Our Contributions

We propose Robust Overlapping Co-clustering (ROCC), a novel approach for discovering dense, arbitrarily positioned co-clusters in large, possibly high dimensional datasets. Our approach is robust in the presence of noisy and irrelevant objects as well as features, which our algorithm automatically detects and prunes during the clustering process. ROCC is based on a systematically developed objective function, which is minimized by an iterative procedure that provably converges to a locally optimal solution. ROCC is also robust to the noise model of the data and can be tailored to use the most suitable distance measure for the data, selected from a large class of distance measures known as Bregman divergences.

The final objective of ROCC is achieved in two steps. In the first step, the Bregman Co-clustering algorithm is adapted to automatically prune away non-informative data

points and perform feature selection by eliminating non-discriminative features and hence cluster only the relevant part of the dataset. This step finds co-clusters arranged in a grid structure, but only a predetermined number of rows and columns are assigned to the co-clusters. Note however that this result cannot be achieved by simply removing some rows/columns from the BCC result. An agglomeration step then appropriately merges similar co-clusters to discover dense, arbitrarily positioned and even overlapping co-clusters. Figure 1 contrasts the nature of the co-clusters identified by ROCC with those found by BCC and illustrates the way in which they are conceptually derived from the partitioned model of BCC.

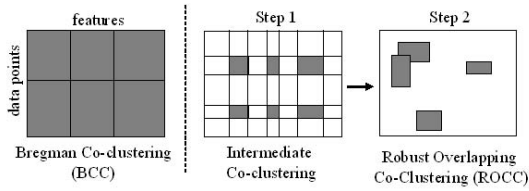


Figure 1. Nature of clusters identified by BCC and ROCC. Shaded areas represent clustered elements, rearranged according to cluster labels, while non-shaded areas denote discarded values.

The ROCC framework has the following key features that distinguish it from existing co-clustering algorithms.

1. The ability to mine the most coherent co-clusters from large and noisy datasets.
2. Detection of arbitrarily positioned and possibly overlapping co-clusters in a principled manner by iteratively minimizing a suitable cost function.
3. Generalization to all Bregman divergences, including squared Euclidean distance, commonly used for clustering microarray data and I-divergence, commonly used for text data clustering [10].
4. The ability to naturally deal with missing data values, without introducing random perturbations or bias in the data.
5. Efficient detection all the co-clusters simultaneously rather than sequentially, enabling scalability to large and high-dimensional datasets.

As far as we know, no existing co-clustering algorithm [30, 28, 1, 6] has all of the above features. Our contribution is significant, since as described in Section 1 there exist several applications where all these features are necessary for discovering meaningful patterns. Through extensive experimental evaluation on synthetic and real microarray data we illustrate that ROCC does significantly

better than a variety of prominent, previously proposed co-clustering techniques. The application of ROCC to the simultaneous feature selection and clustering problem discussed in Section 7.1 highlights the ability of our approach to improve upon human curated feature selection in a completely unsupervised manner.

4. Problem Definition

We begin with the formulation of the first step of the ROCC algorithm that prunes away irrelevant data points and features and clusters the rest into a grid of co-clusters (Refer to Figure 1). Let m be the total number of rows (data points) and n the total number of columns (features). The data can be represented as an $m \times n$ matrix Z of data points and features. Let s_r and s_c be the specified number of rows and columns respectively, to be retained after pruning. If the exact values are not known, it is sufficient to set s_r and s_c conservatively to large values since the algorithm (Section 5.3) does a second round of pruning as needed. Our aim is to simultaneously cluster s_r rows and s_c columns of Z , into a grid of k row clusters and l column clusters. The co-clusters will hence be comprised of $s_r \times s_c$ entries selected from the $m \times n$ entries of Z . Let K and L denote the sets consisting of the s_r clustered rows and the s_c clustered columns respectively. Let ρ be a mapping from the s_r rows $\in K$ to the k row clusters and γ be a mapping from the s_c columns $\in L$ to the l column clusters. Let squared Euclidean distance be the selected distance measure¹. We want to find a co-clustering defined by (ρ, γ) and sets K and L for the specified s_r and s_c that minimize the following objective function

$$\sum_{g=1}^k \sum_{h=1}^l \sum_{u \in K: \rho(u)=g} \sum_{v \in L: \gamma(v)=h} w_{uv} (z_{uv} - \hat{z}_{uv})^2, \quad (1)$$

where z_{uv} is the original value in row u , column v of the matrix, assigned to row cluster g and column cluster h and \hat{z}_{uv} is the value approximated within co-cluster $g-h$. w_{uv} is the non-negative weight associated with matrix entry z_{uv} , which allows the algorithm to deal with missing values and data uncertainties. For example, the weights for known values can be set to 1 and missing values can be effectively ignored by setting their weights to 0. The objective function is hence the element-wise squared error between the original and the approximated value, summed only over the clustered elements ($s_r \times s_c$) of the matrix Z . The value \hat{z}_{uv} can be approximated in several ways, depending on the type of summary statistics that each co-cluster preserves. This is analogous to the notion of approximation schemes used by

¹A more general description, which allows any Bregman divergence as the loss function, is given in Section 5.3.

Bregman co-clustering [1], where the data matrix is reconstructed to preserve a certain set of statistics within each co-cluster. Banerjee et al. [1] identify six possible sets of summary statistics, of increasing complexity, that one might be interested in preserving in the reconstructed matrix \hat{Z} , which lead to 6 different co-clustering bases. Two of these approximation schemes for \hat{z}_{uv} are described in Section 5.3.

In the next step of ROCC, the goal is to agglomerate similar co-clusters to recover the arbitrarily positioned co-clusters. In order to agglomerate co-clusters, we first define a distance measure between two candidate co-clusters ($cc1$ and $cc2$) as follows. Let cc denote the co-cluster formed by the union of the rows and columns in $cc1$ and $cc2$. The matrix entries \hat{z}_{uv} in cc are approximated using the selected approximation scheme. The average element-wise error e for cc is computed as $e = \frac{1}{N} \sum_{z_{uv} \in cc} (z_{uv} - \hat{z}_{uv})^2$, where N is the number of elements in cc . The error e is defined to be the distance between $cc1$ and $cc2$.

5. ROCC Algorithm

5.1. Solving Step 1 of the ROCC Problem

A co-clustering (ρ, γ) , that minimizes the objective function (1), can be obtained by an iterative algorithm. The objective function can be expressed as a sum of row or column errors, computed over the s_r rows and s_c columns assigned to co-clusters. If row u is assigned to row cluster g , i.e., $\rho(u) = g$, the row error is the error summed over the appropriate s_c elements in the row as

$$E_u(g) = \sum_{h=1}^l \sum_{v \in L: \gamma(v)=h} w_{uv} (z_{uv} - \hat{z}_{uv}(g))^2.$$

For a fixed γ , the best choice of the row cluster assignment for row u is the g that minimizes this error, i.e., $\rho^{new}(u) = \operatorname{argmin}_g E_u(g)$. After computing the best row cluster assignment for all the m rows, the rows are sorted in increasing order of their row errors and the top s_r rows with minimum error are selected to participate in the current row clusters. A similar approach is used to assign columns to column clusters. Note that the rows/columns that are not included in the current s_r/s_c rows/columns assigned to co-clusters are still retained since they could be included in the co-clusters in future iterations. The row cluster update step hence selects that fraction of rows that minimize the error among the entire set of rows. Also, the assignment of the selected rows to their corresponding row clusters is done in such a way that it directly minimizes the error. Hence, row cluster updates and similarly column cluster updates reduce the objective function. Updating column cluster assignments could cause the best row assignments to change

and visa versa. Optionally, the row and column cluster re-assignment steps can be repeated several times, in arbitrary order until row and column cluster memberships converge.

Given the current row and column cluster assignments (ρ, γ) , the values \hat{z}_{uv} within each co-cluster have to be updated by recomputing the required co-cluster statistics based on the approximation scheme. This problem is identical to the Minimum Bregman Information (MBI) problem presented in [1] for updating the matrix reconstruction \hat{Z} . Solving the MBI problem for the update step is guaranteed to decrease the objective function.

This iterative procedure is described in Figure 2. Step 1(i) decreases the objective function due to the property of the MBI solution, while Steps 1(ii) and 1(iii) directly decrease the objective function. The objective function hence decreases at every iteration. Since this function is bounded from below by zero, the algorithm is guaranteed to converge. Note that the co-clustering problem is NP-complete, so convergence to only a local minimum is possible.

5.2. Solving Step 2 of the ROCC Problem

We now provide a heuristic to hierarchically agglomerate similar co-clusters. The detailed steps are as follows.

1. **Pruning co-clusters.** Since the desired number of co-clusters is expected to be significantly smaller than the number of co-clusters at this stage of the algorithm, co-clusters with the largest error values can be filtered out in this step. Filtering also reduces the computation effort required by the following merging step. If one has no idea of the final number of co-clusters, a simple and efficient filtering heuristic is to select the error cut-off value as the one at which the sorted co-cluster errors show the largest increase between consecutive values. The co-clusters with errors greater than the cut-off are filtered out. Alternatively, if the final number of co-clusters to be found is pre-specified, it can be used to prune away an appropriate number of co-clusters with the largest errors.
2. **Merging similar co-clusters.** This step involves hierarchical, pairwise agglomeration of the co-clusters left at the end of the pruning step (Step 1) to recover the true co-clusters. Each agglomeration identifies the “closest” pair of co-clusters that can be well represented by a single co-cluster model and are thus probably part of the same original co-cluster, and merges them to form a new co-cluster². “Closest” here is in terms of the smallest value of distance as defined in Section 4. The rows and columns of the new co-cluster

²A variant of this algorithm can be derived by adopting Ward’s method [27] to agglomerate co-clusters. Empirically we found little difference between the two approaches.

consist of the union of the rows and columns of the two merged co-clusters. Merging co-clusters in this manner allows co-clusters to share rows and columns and hence allows partial overlap between co-clusters. If the number of co-clusters to be identified is pre-specified, one can stop merging when this number is reached. If not, merging is continued all the way until only a single co-cluster (or a reasonably small number of co-clusters) is left. The increase in the distance between successively merged co-clusters is then computed and the set of co-clusters just before the largest increase is selected as the final solution.

5.3. Overall ROCC Meta-Algorithm

In this Section we put together the procedures described in Sections 5.1 and 5.2 and present the complete ROCC algorithm. The key idea is to over-partition the data into small co-clusters arranged in a grid structure and then agglomerate similar, partitioned co-clusters to recover the desired co-clusters. The iterative procedure (Section 5.1) is run with large enough values for the number of row and column clusters (k and l). Similarly, the s_r and s_c input parameters are set to sufficiently large values. Since the pruning step (Step 2 in Section 5.2) takes care of discarding less coherent co-clusters, setting $s_r \geq s_r^{\text{true}}$ and $s_c \geq s_c^{\text{true}}$ is sufficient. The resulting $k * l$ clusters are then merged as in hierarchical agglomerative clustering until a suitable stopping criterion is reached. The pseudo-code for the complete algorithm is illustrated in Figure 2.

Approximation Schemes. The ROCC algorithm can use each of the 6 schemes (co-clustering bases) listed by Banerjee et al. [1] for approximating the matrix entries \hat{z}_{uv} . For concreteness, we illustrate two specific approximation schemes with squared Euclidean distance, which give rise to block co-clusters and pattern-based co-clusters respectively³. The meta-algorithm in Figure 2 uses C to refer to the selected co-clustering basis.

Block co-clusters. Let the co-cluster row and column indices be denoted by sets U and V respectively. In this case, a matrix entry is approximated as $\hat{z}_{uv} = z_{UV}$, where $z_{UV} = \frac{1}{|U||V|} \sum_{u \in U, v \in V} z_{uv}$ is the mean of all the entries in the co-cluster.

Pattern-based co-clusters. z_{uv} is approximated as $\hat{z}_{uv} = z_{uV} + z_{UV} - z_{uV}$, where $z_{uV} = \frac{1}{|V|} \sum_{v \in V} z_{uv}$ is the mean of the entries in row u whose column indices are in V and $z_{UV} = \frac{1}{|U|} \sum_{u \in U} z_{uv}$ is the mean of the entries in column v whose row indices are in U . This approximation can identify co-clusters that show a coherent trend or pattern in the data values, making it suitable for clustering gene expression data [7].

³These co-cluster definitions correspond to basis 2 and basis 6 defined by the BCC framework [1] respectively.

Distance Measures. In Section 4 we developed the objective function (1) assuming squared Euclidean distance as the distance measure. The objective function and the iterative procedure to minimize it can be generalized to all Bregman divergences [1]. The selected Bregman divergence is denoted by d_ϕ in Figure 2.

Algorithm: ROCC
Input: $Z_{m \times n}$, s_r , s_c , k , l , basis C , d_ϕ
Output: Set of co-clusters

Step 1
 Begin with a random co-clustering (ρ, γ)
Repeat
Step (i): Update co-cluster models. $\forall [g]_1^k, [h]_1^l$,
 Update statistics for co-cluster (g, h) based on basis C to compute new \hat{z} values

Step (ii): Update ρ
(iia). $\forall [u]_1^m$,
 $\rho(u) = \operatorname{argmin}_g \sum_{h=1}^l \sum_{v \in L: \gamma(v)=h} w_{uv} d_\phi(z_{uv}, \hat{z}_{uv}(g))$
(iib). K = the set of s_r rows with least error from among the m rows

Step (iii): Update γ
(iiia). $\forall [v]_1^n$,
 $\gamma(v) = \operatorname{argmin}_h \sum_{g=1}^k \sum_{u \in K: \rho(u)=g} w_{uv} d_\phi(z_{uv}, \hat{z}_{uv}(h))$
(iiib). L = the set of s_c columns with least error from among the n columns

until convergence
Step 2: Post-process (see text for details)
 (i) Prune co-clusters with large errors.
 (ii) Merge similar co-clusters until stopping criterion is reached.
return identified co-clusters.

Figure 2. Pseudo-code for ROCC Meta-Algorithm

Time Complexity. The iterative procedure (Step 1) for squared Euclidean distance and I-divergence is linear in the size of the input data with a time complexity of $O(mn + mkl + nkl)$ per iteration, for all six bases [1]. $O(mn)$ operations are due to the co-cluster model update step (Step (i)), the $O(mkl)$ operations are due to the computation of the distance of all the m rows from the k row clusters. Note that each such computation involves $O(l)$ operations rather than $O(n)$. Similarly, column reassignment involves $O(nkl)$ operations. The selection steps in the iterative procedure (Steps (iib) and (iiib)) can be performed in linear time by the efficient order statistics based algorithm proposed by Blum et al. [5]. The size of the data that Step 2 operates on is small as compared to the original dataset since only the most relevant matrix entries are retained at the end of Step 1. Although the merging procedure in Step 2 requires $O(M^2)$ operations for computation of all pairwise co-cluster distances, where M is the number of co-clusters to be merged, M is typically small since several co-clusters are discarded in the pruning step (Step 2(i)).

Special Cases. The ROCC framework is very general with a lot of known algorithms as its special cases. Step 1 of ROCC with $s_r = m$ and $s_c = n$ is the same as BCC for all six bases and a selected Bregman divergence. With $l = n$, $s_c = n$ and basis 2 it is identical to Bregman Bubble Clustering [16]. With $l = n$, $s_c = n$, $s_r = m$ and basis 2 it becomes Bregman Clustering and additionally with squared Euclidean distance it is the same as k -means.

5.4. ROCC with Pressurization

The iterative minimization procedure in Step 1 of the ROCC algorithm begins with random initialization for ρ and γ , which could lead to poor local minima. The problem is particularly severe for small s_r and s_c since row and column clusters may not move much and the final clustering may be very heavily influenced by the initialization. This problem can be addressed by using the pressurization technique applied by BBC [16]. Pressurization begins by clustering all the data and iteratively shaving off data points and features till s_r rows and s_c columns are left. Let $s_r^{press(j)}$ and $s_c^{press(j)}$ denote the number of data points and features to be clustered using the Step 1 procedure (Figure 2) in the j^{th} iteration of pressurization. $s_r^{press(1)}$ and $s_c^{press(1)}$ are initialized to m and n respectively, after which these parameters are decayed exponentially till $s_r^{press(j)} = s_r$ and $s_c^{press(j)} = s_c$. The rate of decay is controlled by parameters β_{row} and β_{col} , which lie between 0 and 1. At iteration j , $s_r^{press(j)} = s_r + \lfloor (m - s_r) * \beta_{row}^{j-1} \rfloor$ and $s_c^{press(j)} = s_c + \lfloor (n - s_c) * \beta_{col}^{j-1} \rfloor$. The intuition behind pressurization is that by beginning with all the data being clustered and then slowly reducing the fraction of data clustered, co-clusters can move around considerably from their initial positions to enable the discovery of small, coherent patterns. For speed, each pressurization iteration need not be run to convergence and can be run for a small, fixed number of iterations. This gives competitive results in practice, without a large increase in run time.

5.5. Generative Model for Soft ROCC

The formulation of Step 1 of the ROCC algorithm is based on a very intuitive generative model, which consists of a mixture of $k \times l$ exponential family distributions, corresponding to the $k \times l$ co-clusters, and a background uniform distribution, corresponding to the non-informative data matrix entries. Each element z_{ij} of the data matrix Z is assumed to be generated from this mixture model as follows

$$P(z_{ij}) = \sum_{I=1}^k \sum_{J=1}^l \alpha_{IJ} f_{\psi}(z_{ij} | \theta_{i,j,I,J}) + \alpha_0 p_0, [i]_1^m, [j]_1^n,$$

where α_{IJ} denotes the co-cluster priors, f_{ψ} is an exponential family distribution with cumulant $\psi(\cdot)$ and $\theta_{i,j,I,J}$ is

the natural parameter. The form of $\theta_{i,j,I,J}$ depends on the co-clustering basis. α_0 and p_0 denote the prior probability and the probability density of the uniform distribution. Note that the co-cluster assignments here are soft, i.e., each matrix element belongs to multiple co-clusters and to the background with different probabilities. By assuming that the matrix elements are generated i.i.d. with weights w_{ij} , the data log-likelihood is given by

$$L(\Theta|Z) = \sum_{i=1}^m \sum_{j=1}^n w_{ij} \log P(z_{ij}), \quad (2)$$

where Θ denotes all the model parameters. An EM based technique, with latent variables for cluster memberships can be used to fit this mixture model.

6. Experimental Results

6.1. Synthetic Datasets

We first tested the ROCC algorithm in a controlled setting using synthetically generated data, with true co-cluster labels available for the data matrix entries. The synthetic data is generated by creating a matrix Z of values randomly selected from a uniform distribution (range 0-10), which forms the background. Rectangular blocks of coherent values representing co-clusters are arbitrarily placed in Z and are made to overlap in some datasets. The embedded co-clusters are either block (basis 2) or pattern-based (basis 6). The rows and columns of Z are then randomly permuted to generate the data matrix. Table 1 describes the synthetic datasets used for experimentation.

Dataset	m, n	# co-clusters	basis
1	500, 500	3	2
2	500, 200	4	2
3	500, 500	3	6
4	500, 200	4	6

Table 1. Synthetic Datasets

The evaluation metric used is the relative non-intersection area (RNIA) metric [23] that compares co-clustering solutions when cluster labels for individual matrix entries are available. Given two co-clustering solutions S_1 and S_2 , $RNIA(S_1, S_2)$ is given by $\frac{|U| - |I|}{|U|}$, where $|U|$ and $|I|$ are the number of matrix elements in the union and intersection of S_1 and S_2 respectively. RNIA can be applied to overlapping co-clusters by defining $|U| = \sum_{i,j} \max(n_{ij}^{S_1}, n_{ij}^{S_2})$ and $|I| = \sum_{i,j} \min(n_{ij}^{S_1}, n_{ij}^{S_2})$, where $n_{ij}^{S_1}$ and $n_{ij}^{S_2}$ are the number of co-clusters in S_1 and S_2 respectively that contain the matrix element z_{ij} . The RNIA is hence 0 if S_1 and S_2 are identical and 1 if they are completely disjoint.

Figure 3 compares the RNIA of the co-clusters identified by ROCC (with pressurization (Section 5.4)) and Cheng and Church’s Biclustering algorithm [6] with the true co-clusters, across the 4 synthetic datasets. Note that the objective function optimized by the Biclustering algorithm is equivalent to the ROCC basis 6 objective with squared Euclidean distance, computed over a single bicluster, so datasets 3 and 4 actually match the underlying generative model of Biclustering. On these datasets, the ROCC algorithm is run with an additional local refinement step, which uses each identified co-cluster to seed the iterative procedure described in Section 5.1 with $k = 1$ and $l = 1$. This can help each co-cluster to move around locally, leading to a better solution. This step can be run in parallel for all the co-clusters, to refine all of them simultaneously. The true s_r and s_c values are the only inputs to the ROCC algorithm. The error cut-offs for the ROCC pruning and merging steps (Steps 2(i) and 2(ii)) are determined as described in Section 5.2 and the number of co-clusters is automatically discovered. On the other hand the number of co-clusters is input to Biclustering. For each dataset, the α parameter for Biclustering is set to the average H-score [6] of the true co-clusters. On datasets 1 and 2 which are relatively easy, ROCC does slightly better than Biclustering, while on datasets 3 and 4 ROCC performs significantly better.

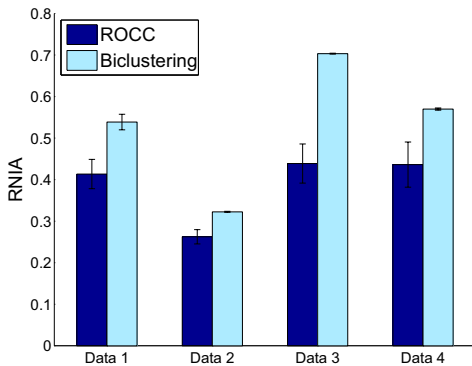


Figure 3. Comparison of ROCC and Biclustering with respect to RNIA on the 4 synthetic datasets.

Figure 4 displays the RNIA values for BCC and ROCC as s_r and s_c are varied from m and n to small values on synthetic datasets 3 and 4. The X-axis represents different fractions of the data being clustered, as rows and columns are pruned. Since BCC clusters all the data, pruning is carried out by a post-processing step. This step sorts the rows and columns by their distance to the corresponding cluster representatives and selects the s_r rows and s_c columns with smallest errors. Note that in case of ROCC, the actual fraction of the data clustered could be smaller than the corresponding X-axis value, due to the pruning of irrelevant co-clusters. The RNIA of the Biclustering solution, given the true number of clusters and a suitable α , is also

displayed as a line on the same plot for comparison. Each plotted value is averaged over 10 randomly initialized runs. On these datasets the ROCC approach consistently achieves a significantly lower RNIA than BCC and does better than Biclustering, in the region near the true s_r and s_c values. We observed that on all the synthetic datasets, ROCC is able to identify the co-clusters and reconstruct the original data matrix very accurately. In almost all cases, ROCC also correctly recovers the true number of co-clusters.

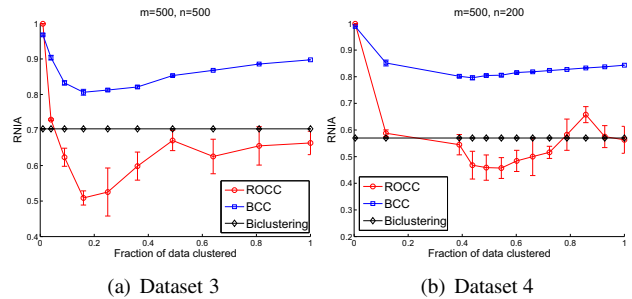


Figure 4. RNIA Plots comparing co-clustering approaches on synthetic data.

6.2. Real Microarray Datasets

We now evaluate the performance of ROCC on two yeast microarray datasets, the Lee dataset [19] and the Gasch dataset [13]. The Lee dataset consists of gene expression values of 5612 yeast genes across 591 experiments and can be obtained from the Stanford Microarray Database (<http://genome-www5.stanford.edu/>). The Gasch dataset consists of the expression values of 6151 yeast genes under 173 environmental stress conditions and is available at http://genome-www.stanford.edu/yeast_stress/. Since the ground truth for both datasets is available only in the form of pairwise linkages between the genes that are known to be functionally related, we compare the quality of the co-clusters identified by different co-clustering algorithms by computing the overlap lift [16] for the genes in each co-cluster. Overlap lift measures how many times more correct links are predicted as compared to random chance and is related to a normalized version of the proportion of disconnected genes measure used by [24]. On these datasets, the aim is to find the most coherent and biologically useful 150 to 200 co-clusters. We run ROCC (with pressurization) on the Lee dataset with the input parameters set to $s_r = 2000$, $s_c = 400$, $k = 50$ and $l = 10$ and on the Gasch dataset with $s_r = 500$, $s_c = 120$, $k = 80$, $l = 15$. Based on the final number of clusters to be identified, Step 2 of ROCC prunes all but the best 200 co-clusters and then continues merging until 150 co-clusters are left. The set of co-clusters just before the largest increase in merge distance is returned as the

solution.

Figure 5 compares the performance of ROCC with prominent co-clustering algorithms, i.e., Cheng and Church’s Biclustering algorithm, the Order Preserving Submatrix algorithm (OPSM) [3], the BiMax algorithm [24], and the BCC algorithm on the Lee and Gasch microarray datasets. Through extensive experimentation, Prelic et al. [24] show that the OPSM and the BiMax algorithms outperform other well known co-clustering algorithms like Samba [25], ISA [4] and xMotif [21] on real microarray data. The BiMax and OPSM results were generated using the BicAT software (<http://www.tik.ee.ethz.ch/sop/bicat/>) [2]. For application of the BiMax algorithm, we discretized the datasets using the discretization threshold suggested in [24]. Since it would be infeasible to evaluate the exponential number of co-clusters identified by BiMax, we selected the first 200 co-clusters for comparison. Though OPSM is designed to return only the best co-cluster, it is extended in BicAT to return upto 100 largest co-clusters among those that achieve the optimal score. The value of the l parameter for OPSM was set to 10. The Biclustering algorithm⁴ is run with the number of clusters equal to 200. The value of the parameter α for the Lee dataset is set to the average H-score [6] of the clusters in the ROCC co-clustering solution with the highest overlap lift from Figure 7(a), i.e. $\alpha = 0.032$. Similarly α is set to 0.017 for the Gasch dataset⁵. We also attempted comparison with the state-of-art subspace clustering algorithm proposed by Yoon et al. [29], but due to its worst case exponential complexity in the number of features this technique did not scale to our datasets. In the Lee and Gasch datasets respectively, around 15% and 3% of the matrix entries are missing. As described in Section 5, ROCC and BCC can ignore missing entries by appropriately setting the weight matrix. The missing entries in the data matrix input to the other algorithms are replaced by random values in the same range as the known expression values. Both ROCC and BCC use squared Euclidean distance and find pattern-based co-clusters. BCC uses the same s_r and s_c values as ROCC for the post processing step as described in Section 6.1. The ROCC, BCC and Biclustering results are averaged over 10 trials, while OPSM and BiMax are deterministic.

Figure 5 shows that on both datasets, ROCC does much better than the other co-clustering approaches in terms of the overlap lift of the gene clusters. The figure also displays above each bar, the percentage of the data matrix entries clustered by the corresponding algorithm. On the Lee dataset, it is interesting that although ROCC clusters a much larger fraction of the data matrix entries than Biclustering, OPSM and BiMax, the co-clusters are of superior quality.

⁴We used the implementation provided by Cheng and Church.

⁵We found the biclustering results to not be very sensitive to the choice of α (range of α values from 0.005 to 0.04 were tried).

The exact execution time values are not given since the algorithms were implemented in different languages and executed on different platforms, however it was evident that ROCC is faster than OPSM and BiMax, but slower than BCC and Biclustering.

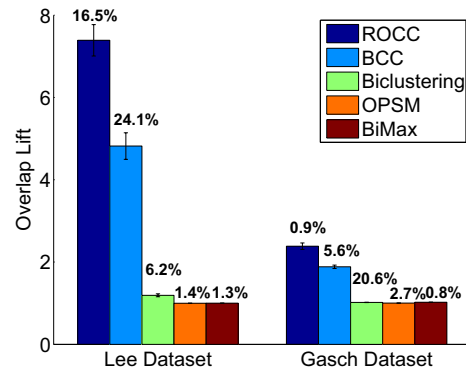


Figure 5. Comparison of ROCC with other co-clustering algorithms on the Lee and Gasch datasets. The number above each bar indicates the percentage of the data matrix entries clustered by each algorithm.

Most of the gene clusters identified by ROCC on the Lee dataset were biologically significant, with very low p-values⁶. Table 2 summarizes some of the identified high purity gene clusters. The coverage (x/y) indicates that x out of the y known members of a category were found. In contrast, the 10 best gene clusters identified by Biclustering had an *average p-value* of 5.50e-04.

# genes	Category(Coverage)	p-value
20	tRNA ligase (8/36)	6.63e-14
63	Endoplasmic reticulum membrane (14/84)	3.886e-14
108	Structural constituent of ribosome (104/206)	<1e-14
20	PF00270-DEAD (12/51)	<1e-14
12	Glycolysis (8/16)	<1e-14
37	Threonine endopeptidase (23/30)	<1e-14
24	PF00660-SRP1-TIP1 (22/30)	<1e-14

Table 2. Examples of biologically significant clusters found by ROCC on the Lee dataset.

7. Other Applications of the ROCC Algorithm

We now discuss two interesting applications of ROCC, which show that it can be very useful in diverse settings.

7.1. Simultaneous Feature Selection and Clustering

A particular instance of the ROCC algorithm can be used to perform feature selection along one axis, while simulta-

⁶The p-values were obtained using *Funspec* (<http://funspec.med.utoronto.ca/>)

neously clustering along the other. ROCC interleaves feature selection with clustering and iteratively improves both, which is intuitively better than independently performing feature selection *a priori* and then clustering using the identified features [17]. Additionally, ROCC also clusters related features, achieving simultaneous dimensionality reduction.

We now consider an exemplary application of ROCC in the above context to a lung cancer microarray dataset [14] represented as a matrix of 12533 genes and 181 human tissue samples. The samples belong to two lung cancer classes, malignant pleural mesothelioma (31 samples) and adenocarcinoma (150 samples). In this application, the aim is to cluster the samples, to recover the 2 existing sample groups in an unsupervised manner, using the expression values of the genes as features. Of the thousands of genes present, many of them are known to be non-informative, redundant and have noisy expression values, which makes feature selection an important issue. We use a version of the dataset that is pre-processed based on domain knowledge, where genes that do not show substantial variation in expression values across the samples are removed as described in [7], resulting in a set of 2401 genes. Even though the preprocessing step results in removing several non-discriminative genes, we apply ROCC to test if any more genes can be identified, that on pruning will improve sample cluster accuracy further. For this application, ROCC (with pressurization) is set up to cluster all the samples and prune along the “gene” axis. Note that the agglomeration procedure (Step 2) is not required for this application.

Sample clustering solutions are evaluated by computing the accuracy of the cluster labels with respect to the true class labels as defined in [7]. Figure 6 displays the sample cluster accuracy of ROCC at different fractions of genes clustered. For comparison, the sample cluster accuracy values of BCC, which uses all the genes to obtain a co-clustering of genes and samples, and *k*-means, which uses all the genes as features to cluster the samples are also plotted as straight lines in the same Figure. These experiments are performed on the column standardized dataset, where every column has zero mean and unit variance. BCC and ROCC use basis 6 with squared Euclidean distance and $k = 20$ and $l = 2$. The results are averaged over 20 runs. One can see that ROCC gives almost perfect clustering, even with only 10% of the genes selected, significantly better than BCC and *k*-means.

7.2. Simultaneous Pruning of Irrelevant Data Points and Features

Identifying arbitrarily positioned, overlapping clusters in noisy datasets is a very challenging problem. However, in several datasets, simply identifying and pruning irrelevant data points and features results in substantially better clus-

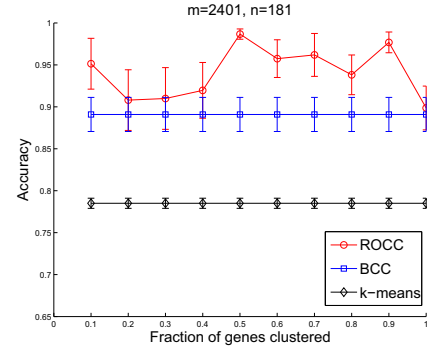


Figure 6. Lung Cancer data: sample clustering accuracy

ters than with clustering all the data. We illustrate the ability of ROCC to do this on the Lee and Gasch datasets. Figure 7 compares the overlap lift of the gene clusters identified by ROCC (with pressurization), BCC, BBC [16] and *k*-means at varying fractions of the data matrix clustered. The number of genes and experiments clustered are gradually increased from very small values to the entire set of genes and experiments and the fraction of the data matrix entries actually assigned to clusters is represented along the X-axis. Since BCC, BBC and *k*-means assign every row to some row cluster, for a fair comparison of these algorithms with ROCC, pruning is carried out by a post-processing step as described in Section 6.1. Note that in these experiments we only use Step 1 of ROCC, thus needing much less computation.

The missing gene expression values in both datasets are set to zero in the data matrix input to BBC and *k*-means. All the algorithms are run with squared Euclidean distance, ROCC and BCC are run with basis 6, and $k = 20$, $l = 10$ for the Lee dataset and $k = 50$, $l = 5$ for the Gasch dataset. ROCC shows a dramatic improvement in performance over all the other approaches. One can observe from Figure 7 that the overlap lift reduces as the fraction of genes and experiments clustered increases, highlighting the fact that the dense clusters existing in the data involve only a small fraction of the dataset. P-value evaluation of the clusters identified in the Lee dataset shows that ROCC has a higher percentage (65%) of clusters with p-values below 10^{-4} as compared to the other approaches.

8. Concluding Remarks

In this paper, we have presented Robust Overlapping Co-clustering as a comprehensive framework capable of dealing with several challenges in clustering real life datasets. ROCC is robust to the presence of irrelevant data points and features and discovers coherent co-clusters very accurately as illustrated in Section 6. Moreover, though ROCC requires several input parameters to be supplied, i.e., s_r , s_c , k and l , it is relatively very robust to the choice of these pa-

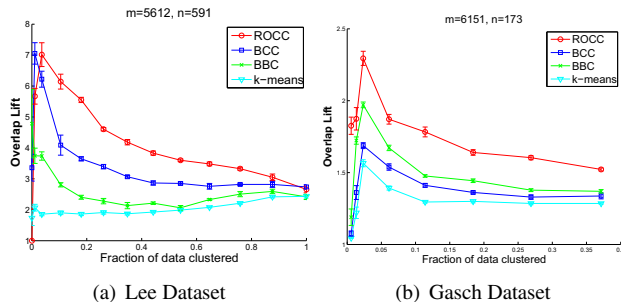


Figure 7. Overlap lift at different fractions of the data clustered on the Lee and Gasch datasets.

rameters because of the post-processing steps as detailed in Section 5.3. While in this paper we focused on clustering microarray data, it would be worthwhile to investigate the applicability of suitable instances of the ROCC framework to clustering problems in different domains like text mining and market basket analysis.

Acknowledgments. This research was supported by NSF grants IIS 0325116, IIS 0307792, IIS 0713142 and CCF 0431257.

References

- [1] A. Banerjee, I. Dhillon, J. Ghosh, S. Merugu, and D. Modha. A generalized maximum entropy approach to bregman co-clustering and matrix approximation. *JMLR*, 8:1919–1986, 2007.
- [2] S. Barkow, S. Bleuler, A. Prelic, P. Zimmermann, and E. Zitzler. Biccat: a biclustering analysis toolbox. *Bioinformatics*, 22(10):1282–1283, 2006.
- [3] A. Ben-Dor, B. Chor, R. Karp, and Z. Yakhini. Discovering local structure in gene expression data: the order-preserving submatrix problem. In *Proc. RECOMB '02*, pages 49–57, 2002.
- [4] S. Bergmann, J. Ihmels, and N. Barkai. Iterative signature algorithm for the analysis of large-scale gene expression data. *Phys. Rev. E. Stat. Nonlin. Soft Matter Phys.*, 67, 2003.
- [5] M. Blum, R. Floyd, V. Pratt, R. Rivest, and R. Tarjan. Time bounds for selection. *J. Comput. Syst. Sci.*, 7(4):448–461, 1973.
- [6] Y. Cheng and G. M. Church. Biclustering of expression data. In *Proc. ICMB '00*, pages 93–103, 2000.
- [7] H. Cho and I. Dhillon. Co-clustering of human cancer microarrays using minimum sum-squared residue co-clustering. To appear in *IEEE/ACM TCBB*, 2008.
- [8] K. Cramer and G. Chechik. A needle in a haystack: Local one-class optimization. In *Proc. ICML '04*, 2004.
- [9] I. Dhillon, S. Mallela, and R. Kumar. A divisive information-theoretic feature clustering algorithm for text classification. *JMLR*, 3:1265–1287, 2003.
- [10] I. Dhillon, S. Mallela, and D. Modha. Information-theoretic co-clustering. In *Proc. KDD '03*, pages 89–98, 2003.
- [11] J. Dy and C. Brodley. Feature selection for unsupervised learning. *JMLR*, 5:845–889, 2004.
- [12] M. Ester, H. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proc. KDD '96*, 1996.
- [13] A. Gasch, P. Spellman, C. Kao, O. Carmel-Harel, et al. Genomic expression program in the response of yeast cells to environmental changes. *Molecular Cell Biology*, 11:4241–4257, 2000.
- [14] G. J. Gordon, R. V. Jensen, L. Hsiao, S. R. Gullans, et al. Translation of microarray data into clinically relevant cancer diagnostic tests using gene expression ratios in lung cancer and mesothelioma. *Cancer Research*, 62:4963–4967, 2002.
- [15] G. Gupta and J. Ghosh. Robust one-class clustering using hybrid global and local search. In *Proc. ICML '05*, pages 273–280, 2005.
- [16] G. Gupta and J. Ghosh. Bregman bubble clustering: A robust, scalable framework for locating multiple, dense regions in data. In *Proc. ICDM '06*, pages 232–243, 2006.
- [17] M. Law, M. Figueiredo, and A.K.Jain. Simultaneous feature selection and clustering using a mixture model. *IEEE Trans. PAMI*, 26(9):1154–1166, 2004.
- [18] L. Lazzeroni and A. B. Owen. Plaid models for gene expression data. *Statistica Sinica*, 12(1):61–86, 2002.
- [19] I. Lee, S. Date, A. Adai, and E. Marcotte. A probabilistic functional network of yeast genes. *Science*, 306:1555–1558, 2004.
- [20] S. C. Madeira and A. L. Oliveira. Biclustering algorithms for biological data analysis: A survey. *IEEE/ACM TCBB*, 1(1):24–45, 2004.
- [21] T. Murali and S. Kasif. Extracting conserved gene expression motifs from gene expression data. *Pacific Symposium on Biocomputing*, 8:77–88, 2003.
- [22] L. Parsons, E. Haque, and H. Liu. Subspace clustering for high dimensional data: a review. *SIGKDD Explor. Newsl.*, 6(1):90–105, 2004.
- [23] A. Patrikainen and M. Meila. Comparing subspace clusterings. *IEEE TKDE*, 18(7):902–916, 2006.
- [24] A. Prelic, S. Bleuler, P. Zimmermann, A. Wille, and et. al. A systematic comparison and evaluation of biclustering methods for gene expression data. *Bioinformatics*, 22(9):1122–1129, 2006.
- [25] A. Tanay, R. Sharan, and R. Shamir. Discovering statistically significant biclusters in gene expression data. *Bioinformatics*, 18:136–144, 2002.
- [26] H. Wang, W. Wang, J. Yang, and P. Yu. Clustering by pattern similarity in large data sets. In *Proc. SIGMOD '02*, pages 394–405, 2002.
- [27] J. Ward. Hierarchical grouping to optimize an objective function. *Journal of American Statistical Association*, 58(301):236 – 244, 1963.
- [28] X. Xu, Y. Lu, A. Tung, and W. Wang. Mining shifting-and-scaling co-regulation patterns on gene expression profiles. In *Proc. ICDE '06*, page 89, 2006.
- [29] S. Yoon, C. Nardini, L. Benini, and G. D. Micheli. Discovering coherent biclusters from gene expression data using zero-suppressed binary decision diagrams. *IEEE/ACM TCBB*, 2(4):339–354, 2005.
- [30] M. Zhang, W. Wang, and J. Liu. Mining approximate order preserving clusters in the presence of noise. In *Proc. ICDE '08*, pages 160–168, 2008.