

UNIFIED APPROACH TO DEPENDENT AND DISPARATE
CLUSTERING OF NONHOMOGENOUS DATA

David R. Easterling^{1 §}, Layne T. Watson^{2,3,4},
Naren Ramakrishnan², Richard F. Helm⁵, Satish Tadepalli⁷,
M. Shahriar Hossain⁶

¹ University of Dayton Research Institute
University of Dayton

300 College Park, Dayton, OH, 45469, USA

Departments of ² Computer Science, ³ Mathematics, ⁴ Aerospace and Ocean
Engineering, and ⁵ Biochemistry

Virginia Polytechnic & State University
Blacksburg, VA 24061, USA

⁶ Department of Computer Science
University of Texas at El Paso
El Paso, TX, 79968, USA

⁷ Bloomberg LP
New York City, NY, 10022, USA

Abstract: There are many data mining settings that involve a combination of attribute-valued descriptors over entities as well as specified relationships between these entities. We present an approach to cluster such nonhomogeneous datasets by using the relationships to impose either dependent clustering or disparate clustering constraints. Unlike prior work that views constraints as Boolean criteria, we present a formulation that allows constraints to be satisfied or violated in a smooth manner. This enables us to achieve dependent clustering and disparate clustering using the same optimization framework by merely maximizing versus minimizing the objective function. We present results on both synthetic data as well as several real-world datasets.

AMS Subject Classification: 62H30, 65H20

Key Words: clustering, disparate clustering, nonhomogeneous clustering, contingency tables, multivariate information bottleneck, homotopy tracking

Received: January 8, 2019

© 2019 Academic Publications

§Correspondence author

1. Introduction

This paper focuses on algorithms for mining nonhomogeneous data involving attribute-valued descriptors over objects from different domains and connected through a relationship. Consider, for instance, the schematic in Figure 1 (top) that reveals many-many relationships between companies and countries. Each company is characterized by a vector indicating stock values, profit margins, earnings ratios, and other financial indicators. Similarly, countries are characterized by vectors in a different space, denoting budget deficits, inflation ratio, unemployment rate, etc. Each company is also related to the countries that it conducts business in.

Since Figure 1 (top) has two different vector spaces and one relation, there can be diverse objectives for clustering such a nonhomogeneous dataset. We study two broad objectives here, which correspond to what we term dependent and disparate clustering. In Figure 1 (bottom left), we seek to cluster companies (by their financial indicators) and cluster countries (by their economic indicators) such that the relationships between individual entities are preserved at the cluster level. In other words, companies within a cluster tend to do business exclusively with countries in a corresponding cluster. In Figure 1 (bottom right), we identify clusters of companies and clusters of countries where the original relationships between companies and countries are actually violated at the cluster level. In other words, clusters in the company space tend to do business with (almost) all clusters in the country space. These two conflicting goals of clustering are meant to reflect two competing hypotheses about companies and their economic performances:

1. **Dependent clustering:** Fortunes/troubles of individual companies are intertwined with the fortunes and woes of the countries they do business in. This school of thought would support the contention that General Motors' (GM) financial troubles began with the collapse of the mortgage industry in the United States.
2. **Disparate clustering:** Diversification helps prepare companies for bad economic times, and hence performance of companies may not necessarily be tied to (and is, hence, independent of) country indicators. An oft cited example here is that Google is well positioned to weather economic storms because its advertisers are broad based.

Observe that in either case, the clusters are still local in their respective attribute spaces, i.e., points within a cluster are similar whereas points across

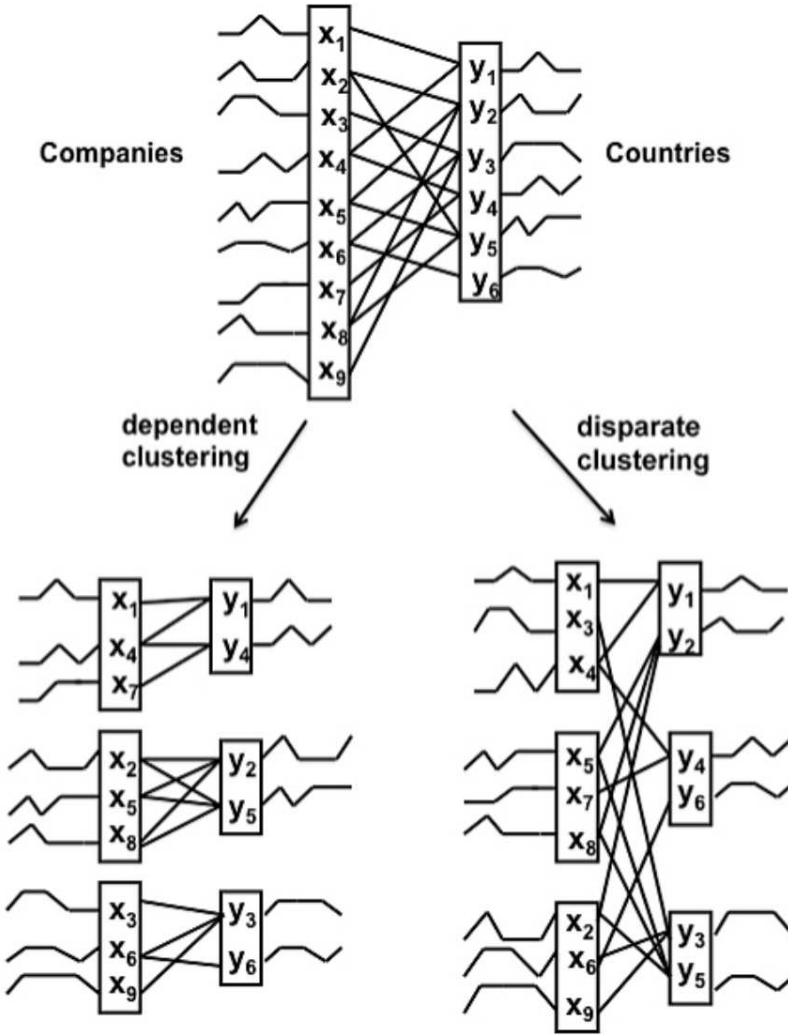


Figure 1: Clustering nonhomogeneous data with two different criteria. Here both domains are clustered into three clusters each based on their attribute vectors. (left) Dependent clustering. (right) Disparate clustering.

clusters are dissimilar. Without advocating any point of view, we posit that it is important to design clustering algorithms that can support both the above analysis objectives. The need for clustering nonhomogeneous data with such conflicting criteria arises in many more contexts, including bioinformatics (stud-

ied here), hypertext modeling, social networks [1], recommender systems, paleontology, and epidemiology.

The chief contributions of this paper are:

- An integrated framework for clustering that unifies dependent clustering and disparate clustering for nonhomogeneous datasets. Unlike prior work that views constraints as Boolean criteria, we present a formulation that allows constraints to be satisfied or violated in a smooth manner.
- While the idea of dependent clustering through a relation has been studied previously [2], the idea of disparate clustering where the objects are of different spaces has not been studied before. We propose this problem here and, moreover, show how we can view dependent clustering and disparate clustering as two sides of the same coin. We propose an integrated objective function whose minimization or maximization leads us to disparate or dependent clustering (respectively).
- The idea of disparate clustering through a relation is closely connected to the current topic of mining multiple, alternative, clusterings [3, 4]. Alternative clusterings are to be expected in high dimensional datasets where different explanations of the data may involve using distinct subspaces of the data. For instance, Figure 1 (right) can be viewed as finding alternative clusterings for different types of objects (companies and countries). The clusterings of (i) the companies and (ii) the countries are alternative in the sense that we cannot use the relational information to recover one from the other and hence they are alternatives with respect to the relational information. To our knowledge, the literature on alternative clustering has not explored this scenario of alternative clustering of objects of different types.
- In addition, a novel approach for studying the intermediate clusterings present along the transition between clusterings, clustering via tracking a continuous homotopy between alternates, is presented and examples are given to show the value of this approach to mining multiple clusterings from a single dataset.

This paper extends [5] by including the additional homotopy research (in Section 5.2) and related experimental findings. This work is the first application of probability-one homotopy methods to clustering, and provides a mathematically rigorous and computationally efficient way to generate and explore multiple alternate clusterings. Besides being more computationally efficient than

a standard parameter sweep, the homotopy method provides insight into the nature of the transition between alternate clusterings.

2. Contingency table clustering

As stated in the introduction, we require our clusters to have two key properties. First, the individual clusters must be local in the respective attribute spaces. Second, when compared across relationships, the clusters must be either highly dependent on each other, or highly independent of each other. We present a uniform framework based on contingency tables that works for both dependent and disparate clusterings.

4	0	0
0	6	0
0	0	4

2	1	1
2	1	2
1	1	3

Table 1: Contingency tables for (left) dependent clustering and (right) disparate clustering for the scenarios from Figure 1.

Table 1 presents contingency tables for the two clusterings from Figure 1. The tables are 3×3 , where the rows denote the clusters from the left domain (here, company clusters) and the columns denote the clusters from the right domain (here, country clusters). The cells indicate the number of entries from the corresponding clusters that are related in the original dataset. For instance, cell (1,1) of Table 1 (left) indicates that there are four relationships between entities in Cluster 1 of the “companies” dataset and entities in Cluster 1 of the “countries” dataset. Observe that the actual sizes of the clusters are not reflected in this matrix, just the number of relationships. Contrast this cell with the corresponding entry of the disparate case, which shows the smaller number of relationships (viz. two) obtained from a different clustering.

Thus the ideal dependent case is best modeled by a diagonal or permutation contingency matrix. In practice, we can aim to achieve a diagonally dominant matrix. Similarly, the disparate case is modeled by a uniform (or near uniform) distribution over all the contingency table entries. It is important to note, however, that we do not have direct control over the contingency table entries. These entries are computed from the clusters, which are in turn defined by the prototype vectors. So the only free variables are the prototype vectors p_1, p_2 but the optimization criteria must be stated in terms of the resulting contingency tables. Mathematically,

$$\mathcal{F}(c(q(d_1, p_1), q(d_2, p_2), r)) \quad (1)$$

is minimized over (p_1, p_2) for disparate clustering or maximized for dependent clustering, where \mathcal{F} is the objective function that evaluates the contingency matrix for either a dependent or a disparate clustering (more on this later), q is the clustering (assignment) function that finds (separate) clusterings of the two datasets using prototypes, c brings the clusterings together and prepares the contingency table with respect to the underlying relation r , d_1 and d_2 are datasets, and p_1 and p_2 are prototypes.

3. Formalisms

Let \mathcal{X} and \mathcal{Y} be two datasets, where $\mathcal{X} = \{x_s \mid s = 1, \dots, n_x\}$ is a set of (real-valued) vectors, each of dimension l_x , i.e., $x_s \in \mathbb{R}^{l_x}$ (likewise, $\mathcal{Y} = \{y_t, t = 1, \dots, n_y\}$, $y_t \in \mathbb{R}^{l_y}$). The many-to-many relationships between \mathcal{X} and \mathcal{Y} are represented by an $n_x \times n_y$ binary matrix B , where $B(s, t) = 1$ if x_s is related to y_t , else $B(s, t) = 0$. Let $C_{(x)}$ and $C_{(y)}$ be the cluster indices, i.e., indicator random variables, corresponding to the datasets \mathcal{X} and \mathcal{Y} and let k_x and k_y be the corresponding number of clusters. Thus, $C_{(x)}$ takes values in $\{1, \dots, k_x\}$ and $C_{(y)}$ takes values in $\{1, \dots, k_y\}$.

3.1. Assigning data vectors to clusters

Let $m_{i,\mathcal{X}}$ be the prototype vector for cluster i in dataset \mathcal{X} (similarly $m_{j,\mathcal{Y}}$). (These are precisely the quantities we wish to estimate/optimize for, but assume they are given in this section). Let $v_i^{(x_s)}$ (likewise $v_j^{(y_t)}$) be the cluster membership indicator variables, i.e., the probability that data sample x_s (y_t) is assigned to cluster i (j) in dataset \mathcal{X} (\mathcal{Y}). Thus, $\sum_{i=1}^{k_x} v_i^{(x_s)} = \sum_{j=1}^{k_y} v_j^{(y_t)} = 1$. The traditional K-means *hard* assignment is given by

$$v_i^{(x_s)} = \begin{cases} 1, & \|x_s - m_{i,\mathcal{X}}\| \leq \|x_s - m_{i',\mathcal{X}}\|, i' = 1, \dots, k_x, \\ 0, & \text{otherwise.} \end{cases}$$

(Likewise for $v_j^{(y_t)}$.) Ideally, we would like a continuous function that tracks these hard assignments to a high degree of accuracy. A standard approach is to use a Gaussian kernel to smooth out the cluster assignment probabilities. Here, we present a novel smoothing formulation that provides tunable guarantees on its quality of approximation and for which the Gaussian kernel is a special case. First we define

$$\gamma_{(i,i')}(x_s) = \frac{\|x_s - m_{i',\mathcal{X}}\|^2 - \|x_s - m_{i,\mathcal{X}}\|^2}{D}, 1 \leq i, i' \leq k_x,$$

where

$$D = \max_{1 \leq s, s' \leq n_x} \|x_s - x_{s'}\|^2, 1 \leq s, s' \leq n_x$$

is the pointset diameter. We now use $\operatorname{argmin}_{i'} \gamma_{(i,i')}(x_s)$ for cluster assignments so the goal is to track $\min_{i'} \gamma_{(i,i')}(x_s)$ with high accuracy. The approach we take is to use the Kreisselmeier-Steinhauser envelope function [6] given by

$$\mathcal{E}_i(x_s) = \frac{-1}{\kappa} \ln \left[\sum_{i'=1}^{k_x} \exp(-\kappa \gamma_{(i,i')}(x_s)) \right],$$

where $\kappa \gg 0$. The \mathcal{E} function is a smooth function that is infinitely differentiable (second, third, ... derivatives exist). Using this the cluster membership indicators are redefined as

$$\begin{aligned} v_i^{(x_s)} &= \frac{\exp[\kappa \mathcal{E}_i(x_s)]}{\sum_{i'=1}^{k_x} \exp[\kappa \mathcal{E}_{i'}(x_s)]} \\ &= \frac{\exp(-\frac{\kappa}{D} \|x_s - m_{i,\mathcal{X}}\|^2)}{\sum_{i'=1}^{k_x} \exp(-\frac{\kappa}{D} \|x_s - m_{i',\mathcal{X}}\|^2)}. \end{aligned} \tag{2}$$

An analogous equation holds for $v_j^{(y_t)}$. The astute reader would notice that this is really the Gaussian kernel approximation with κ/D being the width of the kernel. However, this novel derivation helps tease out how the width must be set in order to achieve a certain quality of approximation. Notice that D is completely determined by the data but κ is a user-settable parameter, and precisely what we can tune.

3.2. Preparing contingency tables

Preparing the $k_x \times k_y$ contingency table (to capture the relationships between entries in clusters across \mathcal{X} and \mathcal{Y}) is now straightforward. We simply iterate over every combination of data entities from \mathcal{X} and \mathcal{Y} , determine whether they have a relationship, and suitably increment the appropriate entry in the contingency table

$$w_{ij} = \sum_{s=1}^{n_x} \sum_{t=1}^{n_y} B(s, t) v_i^{(x_s)} v_j^{(y_t)}. \tag{3}$$

We also define

$$w_{i\cdot} = \sum_{j=1}^{k_y} w_{ij}, \quad w_{\cdot j} = \sum_{i=1}^{k_x} w_{ij},$$

where $w_{i\cdot}$ and $w_{\cdot j}$ are the row-wise and column-wise counts of the cells of the contingency table respectively.

We will find it useful to define the row-wise random variables α_i , $i = 1, \dots, k_x$ and column-wise random variables β_j , $j = 1, \dots, k_y$ with probability distributions

$$p(\alpha_i = j) = p(C_{(y)} = j \mid C_{(x)} = i) = \frac{w_{ij}}{w_{i\cdot}}, \quad (4)$$

$$p(\beta_j = i) = p(C_{(x)} = i \mid C_{(y)} = j) = \frac{w_{ij}}{w_{\cdot j}}. \quad (5)$$

The row-wise distributions represent the conditional distributions of the clusters in dataset \mathcal{X} given the clusters in \mathcal{Y} ; the column wise distributions are also interpreted analogously.

3.3. Evaluating contingency tables

Now that we have a contingency table, we must evaluate it to see if it reflects a dependent or disparate set of clusterings (as the requirement may be). Ideally, we would like one criterion that when minimized leads to a disparate clustering and when maximized leads to a dependent clustering.

For this purpose, we compare the row-wise and column-wise distributions from the contingency table entries to the uniform distribution $U(\cdot)$. (In the example from Table 1, there are three row-wise distributions and three column-wise distributions.) For dependent clusters, the row-wise and column-wise distributions must be far from uniform, whereas for disparate clusters, they must be close to uniform. Identify the probability density functions for the random variables α_i and β_j with α_i and β_j , respectively, in the KL-divergences below. We use KL-divergences to define our unified objective function

$$\begin{aligned} \hat{\mathcal{F}} = & \frac{1}{k_x} \sum_{i=1}^{k_x} D_{KL} \left(\alpha_i \parallel U \left(\frac{1}{k_y} \right) \right) \\ & + \frac{1}{k_y} \sum_{j=1}^{k_y} D_{KL} \left(\beta_j \parallel U \left(\frac{1}{k_x} \right) \right), \end{aligned} \quad (6)$$

where the KL-divergence between distributions $p_1(x)$ and $p_2(x)$ over the sample space X is given by

$$D_{KL}[p_1||p_2] = \sum_{x \in X} p_1(x) \log \frac{p_1(x)}{p_2(x)}.$$

$D_{KL}[p_1||p_2]$ measures the inefficiency of assuming that the distribution is p_2 when the true distribution is actually p_1 .

Note that the row-wise distributions take values over the columns $1, \dots, k_y$ and the column-wise distributions take values over the rows $1, \dots, k_x$. Hence the reference distribution for row-wise variables is over the columns, and vice versa. Also, observe that the row-wise and column-wise KL-divergences are averaged to form $\hat{\mathcal{F}}$. This is to mitigate the effect of lopsided contingency tables ($k_x \gg k_y$ or $k_y \gg k_x$) wherein it is possible to optimize $\hat{\mathcal{F}}$ by focusing on the “longer” dimension without really ensuring that the other dimension’s projections are close to uniform.

Finally observe that the KL-divergence of any distribution with respect to the uniform distribution is proportional to the negative entropy ($-H$) of the distribution. Thus we are essentially aiming to minimize or maximize (for dependent or independent clusters) the entropy of the cluster conditional distributions between pairs of two datasets.

4. Algorithms

Now we are ready to formally present our data mining algorithms as optimization over the space of prototypes.

4.1. Disparate clustering

Here the goal is to minimize $\hat{\mathcal{F}}$, a nonlinear function of $m_{i,x}$ and $m_{j,y}$. For this purpose, we adopt an augmented Lagrangian formulation with a quasi-Newton trust region algorithm. We require a flexible formulation with equality constraints (i.e., that mean prototypes lie on the unit hypersphere) and bound constraints (i.e., that the prototypes are bounded by the max and min (componentwise) of the data, otherwise the optimization problem has no solution). The LANCELOT software package [7] provides just such an implementation.

For ease of description, we “package” all the mean prototype vectors for clusters from both datasets (there are $\eta = k_x + k_y$ of them) into a single vector v of length t . The problem to solve is then

$$\begin{aligned} \min \hat{\mathcal{F}}(v) \quad \text{subject to} \quad & h_i(v) = 0, \quad i = 1, \dots, \eta, \\ & L_j \leq v_j \leq U_j, \quad j = 1, \dots, t, \end{aligned}$$

where ν is a t -dimensional vector and $\hat{\mathcal{F}}$, h_i are real-valued functions continuously differentiable in a neighborhood of the box $[L, U]$. Here the h_i ensure that the mean prototypes lie on the unit hypersphere (i.e., they are of the form $h_1 = \|m_{1,\mathcal{X}}\| - 1$, $h_2 = \|m_{2,\mathcal{X}}\| - 1$, \dots , $h_\eta = \|m_{(k_x+k_y),\mathcal{Y}}\| - 1$). The bound constraints are uniformly set to $[-1, 1]$. The augmented Lagrangian L is defined by

$$L(\nu, \mu, \varphi) = \hat{\mathcal{F}}(\nu) + \sum_{i=1}^{\eta} (\mu_i h_i(\nu) + \varphi h_i(\nu)^2), \quad (7)$$

where the μ_i are Lagrange multipliers and $\varphi > 0$ is a penalty parameter. The augmented Lagrangian method (implemented in LANCELOT) to solve the constrained optimization problem above is given in *OptPrototypes*.

In Step 1 of *OptPrototypes*, we initialize the prototypes using a K-means algorithm (i.e., one which separately finds clusters in each dataset without coordination), package them into the vector ν , and use this vector as starting points for optimization. For each iteration of the augmented Lagrangian method, we require access to $\hat{\mathcal{F}}$ and $\nabla \hat{\mathcal{F}}$ which we obtain by invoking Algorithm *Problem-Setup*.

This routine goes step-by-step through the framework developed in earlier sections to link the prototypes to the objective function. There are no parameters in these stages except for κ which controls the accuracy of the KS approximations. It is chosen so that the KS approximation error is commensurate with the optimization convergence tolerance. Gradients (needed by the trust region algorithm) are mathematically straightforward but tedious, so are not explicitly given here (see [8]).

Algorithm 1 *OptPrototypes*

1. Choose initial values $\nu_{(0)}$ (e.g., via a K-means algorithm), $\mu_{(0)}$, set $k := 0$, and fix $\varphi > 0$.
 2. For fixed $\mu_{(k)}$, update $\nu_{(k)}$ to $\nu_{(k+1)}$ by using one step of a quasi-Newton trust region algorithm for minimizing $L(\nu, \mu_{(k)}, \varphi)$ subject to the constraints on ν . Call *ProblemSetup* with ν as needed to obtain F and ∇F .
 3. Update μ by $\mu_{(k+1)_i} = \mu_{(k)_i} + 2\varphi h_i(\nu_{(k)})$ for $i = 1, \dots, \eta$.
 4. If $(\nu_{(k)}, \mu_{(k)})$ has converged, stop; else, set $k := k + 1$ and go to (2).
 5. Return ν .
-

Algorithm 2 *ProblemSetup*

1. Unpackage ν into values for mean prototype vectors.
2. Use (2) (and its analog) to compute $v_i^{x_s}$ and $v_j^{y_t}$.
3. Use (3) to obtain contingency table counts w_{ij} .
4. Use (4) and (5) to define r.v.s α_i and β_j .
5. Use (6) to compute $\hat{\mathcal{F}}$ and $\nabla\hat{\mathcal{F}}$ (see [8]).
5. Return $\hat{\mathcal{F}}, \nabla\hat{\mathcal{F}}$.

Modulo the time complexity of K-means (which is used for initializing the prototypes), the per-iteration complexity of the various stages of our algorithm can be given as

Step	Time Complexity
Assigning vectors to clusters	$\mathcal{O}(n_x l_x k_x + n_y l_y k_y)$
Preparing contingency tables	$\mathcal{O}(k_x k_y n_x n_y)$ (naïve) $\mathcal{O}(k_x k_y \beta)$ (replicated)
Evaluating contingency tables	$\mathcal{O}(k_y k_x)$
Optimization	$\mathcal{O}((\eta + 1)t^2)$

First, observe that this is a continuous, rather than discrete, optimization algorithm, and hence the overall time complexity depends on the number of iterations, which is an unknown function of the requested numerical accuracy. The step of assigning vectors to clusters takes place independently in the two datasets, so the time complexity has two components. For each vector, we compare it to each mean prototype, and an inner loop over the dimensionality of the vectors gives $\mathcal{O}(n_x l_x k_x + n_y l_y k_y)$. The straightforward way to prepare contingency tables as suggested by (3) gives rise to a costly computation, since for each cell of the contingency table (there are $k_x k_y$ of them), we will expend $\mathcal{O}(n_x n_y)$ computations. In [8] we show how we can reduce this by an order of magnitude using a method of ‘replicating’ vectors which helps us treat the relationship matrix β as if it were one-to-one. In this case, the per-cell complexity will simply be a linear function of the nonzero entries in β , i.e., $|\beta|$. Evaluating the contingency tables requires us to calculate KL-divergences which are dependent on the sample space over which the distributions are compared and the number of such comparisons. There are two terms, one for row-wise distributions, and one for column-wise distributions. Finally, the time complexity of the optimization is $\mathcal{O}((\eta + 1)t^2)$ per iteration, and the space complexity is also $\mathcal{O}((\eta + 1)t^2)$, mostly for storage of Hessian matrix approximations of $\hat{\mathcal{F}}$ and h_i . Note that $t = k_x l_x + k_y l_y$. In practice, to avoid sensitivity to local minima, we perform several random restarts of our approach, with different initializations

of the prototypes.

4.2. Dependent clustering

This proceeds exactly as above except the goal now is to $\min -\hat{\mathcal{F}}$ (i.e., to maximize $\hat{\mathcal{F}}$). Simply replacing $\hat{\mathcal{F}}$ with $-\hat{\mathcal{F}}$ in the above algorithm conducts dependent clustering. For ease of description later, we henceforth refer to $\hat{\mathcal{F}}$ as $\hat{\mathcal{F}}_i$ (for independent) and to $-\hat{\mathcal{F}}$ as $\hat{\mathcal{F}}_d$ (for dependent).

1	0	0
0	1	0
0	0	14

1.5	1.5	1.5
1.5	1.5	1.5
1.5	1.5	1.5

Table 2: Degenerate contingency tables for (left) dependent clusters and (right) disparate clusters. These are bad solutions to be avoided because the clusters in (a) are highly imbalanced and (b) is obtained by trivially assigning all points to all clusters.

4.3. Regularization

Degenerate situations can arise as shown in Table 2. In the dependent case, we might obtain a diagonal contingency table but with imbalanced cluster sizes. In the independent case, the data points are assigned with equal probability to every cluster, resulting in a trivial solution for ensuring that the contingency table resembles a uniform distribution. See [8] for how to add additional terms in the objective function to alleviate both these issues.

4.4. How Many Clusters?

Selecting the number of clusters here has a direct mapping to the sufficient statistics of contingency tables necessary to capture differences between distributions. We have used the minimum discrimination information (MDI) principle (discussed later) for model selection. Since choosing the number of clusters is a research issue all unto itself, this is not pursued further here.

5. Experiments

We evaluate our approach using both synthetic and real datasets. The questions answered through the experiments in this section are:

1. Can we realize classical constrained clustering and alternative clustering scenarios (i.e., over a single dataset) using our framework? (Section 5.1)
2. Can we realize the same classical constrained clustering scenarios using the alternative homotopy framework? (Section 5.2)
3. How much does our emphasis on clustering relations compromise locality of clusters in the respective attribute spaces? (Section 5.4)
4. How does our approach (of defining an integrated objective function and locally minimizing it) scale? (Section 5.4)
5. As the number of clusters increases, does it become easier or more difficult to achieve dependent and disparate clusterings? (Section 5.4)
6. Can we pose integrated dependent and disparate clustering formulations over nonhomogeneous data involving multiple datasets and relations? (Section 5.5)
7. In mining nonhomogeneous datasets with multiple criteria, what is the effect of varying the emphasis of different criteria on the clustering results? (Sections 5.2 and 5.5)

5.1. Constrained clustering

In constrained clustering, we are given a single dataset D with instance level constraints such as must-link and must-not-link constraints [9, 10]. We can model such problems in our relational context as shown in Figure 2 (a), (b), and (c). We create two copies of D into D_1 and D_2 . In the case with only must-link (ML) constraints (Figure 2 (a)), such as between x_1 and x_3 , we create a relation between the entries: x_1 of D_1 and x_3 of D_2 , and between entries: x_3 of D_1 and x_1 of D_2 . In addition we include relations between the same instances in D_1 and D_2 . Applying the dependent clustering criterion $\hat{\mathcal{F}}_d$ on this dataset will realize the constrained clustering scenario. Conversely, as shown in Figure 2 (b), for must-not-link (MNL) constraints we would create relations between the entries that should not be brought together, and use $\hat{\mathcal{F}}_i$ as the optimization criterion. In Figure 2 (a), the relations would force the clusterings to be dependent and as a result, either clustering would respect the ML constraints. In Figure 2 (b), the F_i objective will force the clusterings to violate the relations (which are really MNL constraints).

Going further, we can combine the above modeling approaches in Figure 2

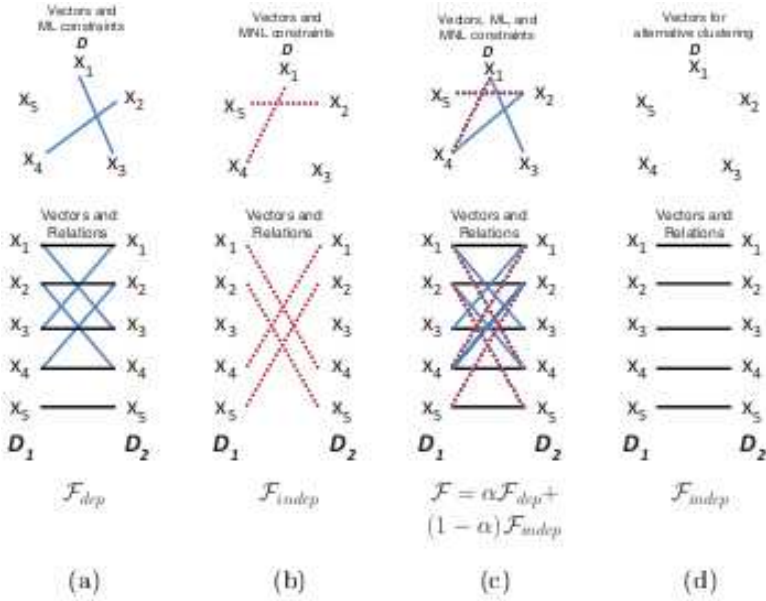
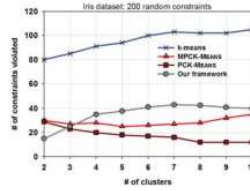
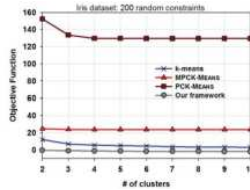


Figure 2: Realizing classical single-dataset clustering scenarios using our framework. (a) Clustering with must-link constraints. (b) Clustering with must-not-link constraints. (c) Clustering with both must-link and must-not-link constraints. (d) Finding alternative clusterings.

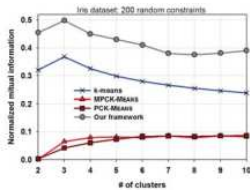
(c) which has both ML and MNL constraints. For this scenario, the optimization criterion is essentially a convex combination of both $\hat{\mathcal{F}}_d$ and $\hat{\mathcal{F}}_i$. As we vary α smoothly from 0 to 1, we increase our emphasis from satisfying the ML constraints to satisfying the MNL constraints. Here we set α to 0.5 (and explore other settings in future sections). We compare our constrained clustering framework with simple unconstrained K-means and two constrained K-means algorithms (MPCK-Means and PCK-Means) from [13]. Overall, the number of constraint violations from our approach (Figure 3 (a)) is worse than that of either MPCK-Means and PCK-Means, except for a small number of clusters. This is to be expected since our method does not take a strict (Boolean) view of constraint satisfactions. Conversely, the objective function in our approach is the best possible value (Figure 3 (b)) when compared with the solutions obtained by the other three algorithms. Finally, as shown in Figure 3 (c), the normalized mutual information score (between the cluster assignments and the class labels) is best for our approach compared to the other three algorithms.



(a)



(b)



(c)

Figure 3: Comparison of our approach with unconstrained K-means and two other constrained clustering formulations. We cluster the Iris dataset with randomly generated 100 ML and 100 MNL constraints. Results are averaged over 20 runs each. (a) number of constraints violated. (b) Objective function. (c) Normalized mutual information.

This shows that taking a soft view of constraints does not compromise the locality of the mined clusters.

5.2. Homotopy tracking

In this section, we consider the homotopy tracking method of [11]. Homotopy methods are systematic approaches to characterize solution sets by smoothly tracking solutions from one formulation to another (in this case, from an unconstrained formulation to a constrained formulation). This can allow the effect of changing λ on the quality and nature of the solutions to be mathematically

characterized. Smoothly tracking solutions as λ varies provides a holistic understanding of the interplay between the algorithm and a dataset. The resulting tradeoff curve can yield information about the nature of the problem and the probability of improvement offered by constraints.

For the purposes of this section, let superscripts denote vector indices and subscripts denote components of vectors and scalar indices unless otherwise indicated. Let all norms be 2-norms unless otherwise indicated and let all distances be Euclidean distances. Let \mathbb{R}^n denote n -dimensional real Euclidean space and let $\mathbb{R}^{n \times m}$ be the set of real $n \times m$ matrices. Let the i th row of a matrix $A \in \mathbb{R}^{n \times m}$ be denoted by A_i and the j th column by A_j . Finally, for a vector $x \in \mathbb{R}^n$, $x > 0$ means all $x_i > 0$, $x \geq 0$ means all $x_i \geq 0$, and $x \geq 0$ means $x \geq 0$ but $x \neq 0$.

Given a set $\hat{X} = \{x^i \mid x^i \in \mathbb{R}^d, i = 1, 2, \dots, k\}$ of k points (cluster representatives) in d dimensions, let $X = \text{vec}(x^1, x^2, \dots, x^k) \in \mathbb{R}^{kd}$. Given a set $\hat{Y} = \{y^i \mid y^i \in \mathbb{R}^d, i = 1, 2, \dots, n\}$ of n data points in d dimensions, let $Y = \text{vec}(y^1, y^2, \dots, y^n) \in \mathbb{R}^{nd}$. Represent a constraint by the vector $c = (a, b, z, w) \in \mathbb{R}^{2d+2}$ of two data points $a, b \in \hat{Y}$, an identifier $z = \pm 1$, and a degree-of-belief weight $\mathbb{R} \ni w > 0$, where an identifier of $z = 1$ means that a and b are bound by a must-link constraint (i.e., must be in the same cluster) and an identifier of $z = -1$ means that a and b are bound by a cannot-link constraint (can not be in the same cluster). Given a set $\hat{C} = \{c^i \mid c^i \in \mathbb{R}^{2d+2}, i = 1, 2, \dots, q\}$ of q constraints, let $C = \text{vec}(c^1, c^2, \dots, c^q) \in \mathbb{R}^{q(2d+2)}$.

For a data point $y \in \hat{Y}$ and two cluster prototypes x^i, x^j define the comparator function $D_H : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ by

$$D_H(x^i, x^j, y) = (\max \{0, \|x^i - y\|^2 - \|x^j - y\|^2\})^4.$$

Note that D_H is three times continuously differentiable, $D_H \geq 0$, and $D_H(x^i, x^j, y) > 0$ if and only if the distance between y and x^i is larger than the distance between y and x^j .

Given $a, b \in \hat{Y}$, let the must-link function

$$F_m : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^{kd} \rightarrow \mathbb{R}$$

be defined by

$$F_m(a, b, X) = \prod_{i=1}^k \left(\sum_{j=1, j \neq i}^k D_H(x^i, x^j, a) + D_H(x^i, x^j, b) \right)$$

and let the cannot-link function $F_c : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^{kd} \rightarrow \mathbb{R}$ be defined by

$$F_c(a, b, X) = \sum_{i=1}^k \left(\prod_{j=1, j \neq i}^k D_H(x^j, x^i, a) D_H(x^j, x^i, b) \right).$$

Then the following observations are easily verified.

Observation 1. F_m and F_c are nonnegative and three times continuously differentiable.

Observation 2. For any must-link constraint $c = (a, b, 1, w) \in \hat{C}$, the must-link function $F_m(a, b, X) = 0$ if and only if constraint c is satisfied.

Observation 3. For any cannot-link constraint $c = (a, b, -1, w) \in \hat{C}$, the cannot-link function $F_c(a, b, X) = 0$ if and only if constraint c is satisfied.

Observation 4. The penalty function

$$F(C, X) = \sum_{\{i:z_i=1\}} w_i F_m(a^i, b^i, X) + \sum_{\{i:z_i=-1\}} w_i F_c(a^i, b^i, X)$$

is zero if and only if all the constraints in \hat{C} are satisfied.

This penalty function is not infallible; degenerate solutions are still problematic, as they can cause a zero in the penalty function, and the homotopy map also needs to be bounded on the solution curve.

Consider again the bounding constraint. Instead of forcing each prototype to exist on the unit hypersphere as was done previously, a straightforward concave function $\Psi : \mathbb{R}^{kd} \rightarrow \mathbb{R}$ to achieve bounding is $\Psi(X) = B - \sum_{i=1}^n \|x^i\|^2 \geq 0$, where $B \in \mathbb{R}$ is a given large constant. Second, to prevent the degenerate condition noted above, a set of constraints $g_i : \mathbb{R}^{kd} \rightarrow \mathbb{R}$ can be constructed as $g_i(X) = \epsilon_g - \|x^{i1} - x^{i2}\|^2 \leq 0$, where $1 \leq i \leq \binom{k}{2}$, $x^{i1}, x^{i2} \in \hat{X}$ are different cluster representatives and $\epsilon_g > 0$ is a small constant. Note that these constraints are differentiable everywhere, and satisfy the reverse convex constraint qualification at \mathcal{X} if $\Psi(\mathcal{X}) > 0$ is inactive. If the active constraints at \mathcal{X} satisfy a constraint qualification (e.g., Arrow-Hurwicz-Uzawa), then the resulting optimization problem

$$\begin{aligned} & \min_X F(C, X) \\ & \text{subject to } -\Psi(X) \leq 0, \\ & g_i(X) \leq 0, \quad 1 \leq i \leq \binom{k}{2} \end{aligned} \tag{8}$$

satisfies the Karush-Kuhn-Tucker (KKT) necessary conditions and may be considered as a potential minimum at \mathcal{X} .

Let \mathcal{E} be redefined to fit this formulation as

$$\mathcal{E}(X) = \frac{1}{\kappa} \ln \left[\sum_{i=1}^{\binom{k}{2}} \exp(-\kappa g_i(X)) \right].$$

Let $g_{\max}(X) = \max_{1 \leq i \leq \binom{k}{2}} g_i(X)$ and note that

$$g_{\max}(X) \leq \mathcal{E}(X) \leq g_{\max}(X) + \frac{\ln(k(k-1)/2)}{\kappa},$$

which means that if $\mathcal{E}(X) \leq 0$ then $g_i(X) \leq 0$ for all i . Thus a reasonable approximation for (8) utilizing \mathcal{E} to reduce the number of inequality constraints can be defined by

$$\begin{aligned} & \min_X F(C, X) \\ & \text{subject to } -\Psi(X) \leq 0, \\ & \mathcal{E}(X) \leq 0. \end{aligned} \tag{9}$$

Let $\hat{\Phi} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be the continuous positively oriented nonlinear complementarity function

$$\hat{\Phi}(a, b) = -|a - b|^3 + a^3 + b^3.$$

Let $\Phi : \mathbb{R} \times \mathbb{R} \times [0, 1) \times (0, \infty) \rightarrow \mathbb{R}$ be the λ -dependent approximation of $\hat{\Phi}$

$$\Phi(a, b, \lambda, h) = -|a - b|^3 + a^3 + b^3 - (1 - \lambda)h.$$

Note that $\Phi \rightarrow \hat{\Phi}$ as $\lambda \rightarrow 1$. Let

$$\hat{K}(X) = \sum_{i=1}^{|\hat{Y}|} \frac{k}{\sum_{j=1}^k \frac{1}{\|y^i - x^j\|^2}}$$

approximate the K-means criterion (to allow for continuous derivatives) and note that the discontinuities in this formulation are removable. The Lagrangian function associated with (9) is

$$\tilde{L}(X, \tilde{\mu}, \tilde{\nu}) = F(C, X) - \tilde{\mu}\Psi(X) + \tilde{\nu}\mathcal{E}(X),$$

where $\tilde{\mu}$ and $\tilde{\nu}$ are the Lagrange multipliers, and a KKT point $(\bar{X}, \bar{\mu}, \bar{\nu})$ for (9) satisfies

$$\begin{aligned} \nabla_X \tilde{L}(\bar{X}, \bar{\mu}, \bar{\nu}) &= 0, \\ 0 &\leq \bar{\mu} \perp \Psi(\bar{X}) \geq 0, \\ 0 &\leq \bar{\nu} \perp -\mathcal{E}(\bar{X}) \geq 0. \end{aligned}$$

While this is not a KKT point (and thus a potential minimum) for (8), the distinction between the formulations is minimal. For a relatively small number

of clusters, κ can be tweaked to strengthen the approximation without encountering numerical difficulties.

Thus, we can define $\rho : \mathbb{R}^{kd} \times (0, \infty) \times (0, \infty) \times [0, 1) \times \mathbb{R}^{kd} \times \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^{kd+2}$ by

$$\tilde{\rho}(k^0, h_0, h_1, \lambda, X, \tilde{\mu}, \tilde{\nu}) = \begin{pmatrix} (1 - \tanh(60\lambda))(X - k^0) + \tanh(60\lambda)\varphi(\lambda, X, \tilde{\mu}, \tilde{\nu}) \\ \Phi(\lambda, \tilde{\mu}, \Psi(X), h_0) \\ \Phi(\lambda, \tilde{\nu}, -\mathcal{E}(X), h_1) \end{pmatrix},$$

where

$$\varphi(\lambda, X, \mu, \nu) = ((1 - \lambda)\nabla_X \mathcal{E}(X) + \lambda \nabla_X \tilde{L}(X, \mu, \nu))^T,$$

h_0 and h_1 are selected to make the initial Φ functions 0 at $\lambda = 0$, and k^0 is an initial solution to the K-means approximation criterion. ρ is thus a strong probability-one homotopy map for exploring clustering with constraints. Note that B is selected so that $\Psi(k^0) > 0$. Utilizing a homotopy tracker such as HOMPACT90 [12], we can track this map from $\lambda = 0$ to $\lambda = 1$ to find clusterings that fulfill the clustering constraints while maintaining the clustering hypothesis.

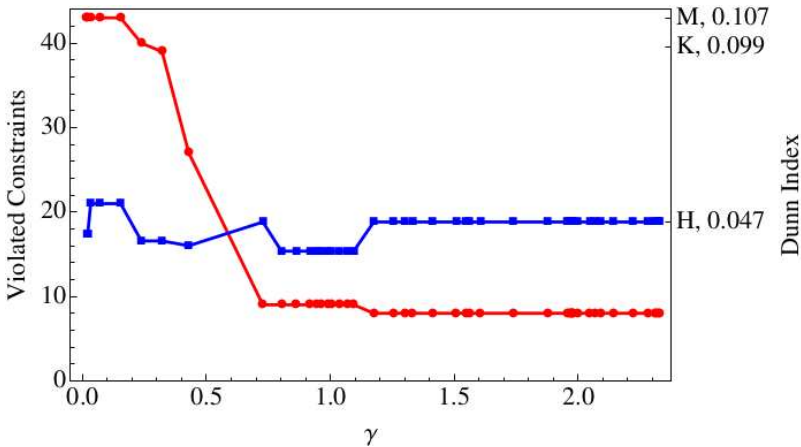


Figure 4: The iris dataset with “easy” constraints (no more than one MNL constraint per datapoint). The Dunn index is tracked against the arc length of γ in blue, while the satisfied constraints are tracked in red. The Dunn Indices for the final homotopy ($\tilde{\rho}$) clustering (“H”), MK-Means clustering (“M”) [13], and K-Means clustering (“K”) are also shown.

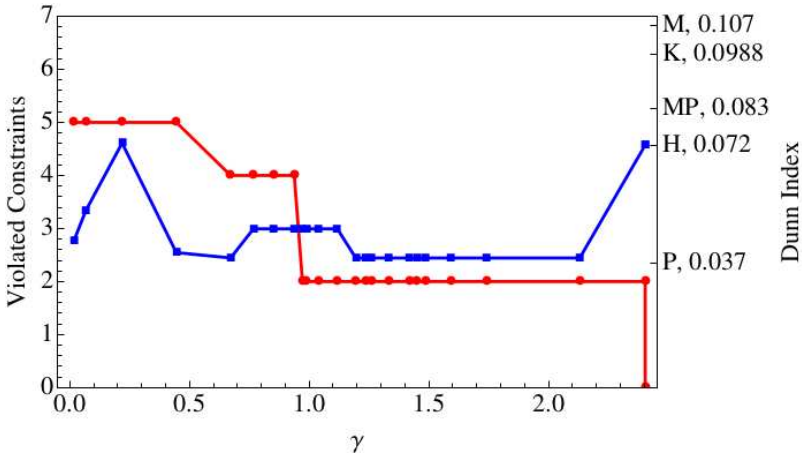


Figure 5: The iris dataset with “hard” constraints (consisting of multiple MNL constraints per point). The Dunn Indices for the final homotopy ($\tilde{\rho}$) clustering (“H”), MK-Means clustering (“M”), PK-Means clustering (“P”), MPK-Means clustering (“MP”), and K-Means clustering (“K”) are also shown.

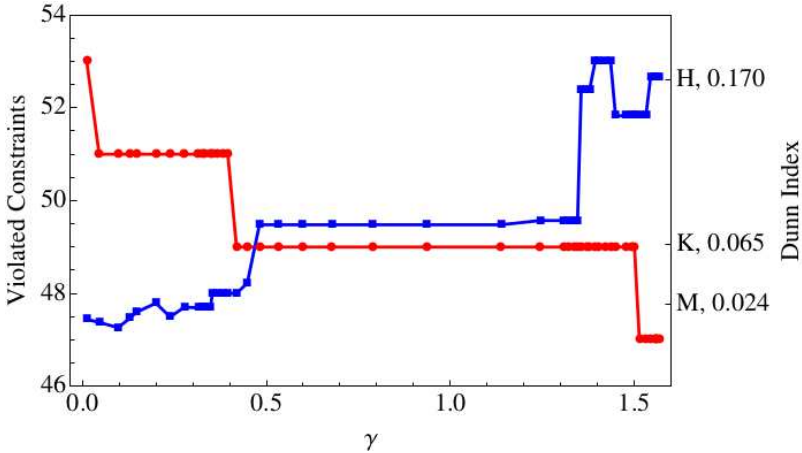


Figure 6: The liver dataset with “easy” constraints.

Figures 4–7 utilize datasets taken from the UCI Machine learning repository. One hundred random (valid) constraints are generated based on the true clustering for each of these datasets, and the homotopy is tracked from a random K-means solution to a local minimum based on the constraints. The Dunn

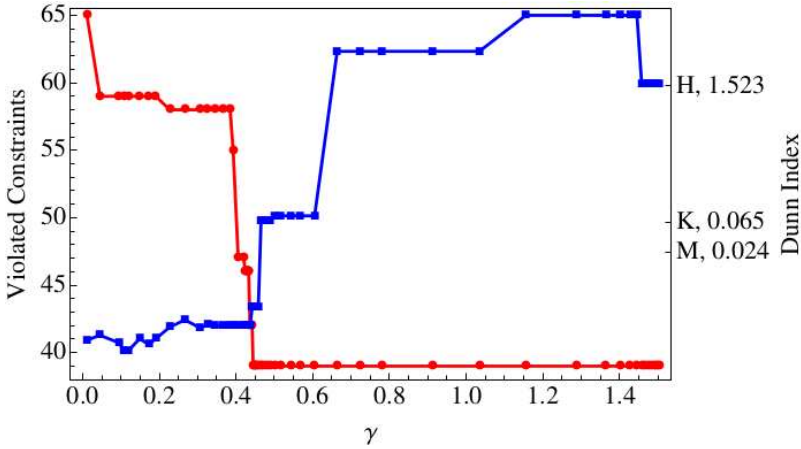


Figure 7: The liver dataset with “hard” constraints.

Index (blue) is tracked against the number of violated constraints (red). The Dunn Index is a fairly simple index for satisfaction of the clustering hypothesis, and it can't be used on nonconvex clusterings (such as those sometimes generated by the MPCK-means and PCK-means algorithms). However, where applicable, the result determined by those algorithms is shown as well as the MK-means algorithm and the best K-means result from 50 independent experiments (note that this is separate from the K-means solution used to prime the homotopy method for these examples).

Of particular interest is the fact that the change in the Dunn Index is not monotonic. As alternative clusters are encountered, their satisfaction of the Dunn Index may greatly increase or greatly decrease, depending on the strength of the constraints. Note that the final Dunn Index (as the clustering becomes closer to the “true” clustering) is always higher than the starting point, but local Dunn Indices may be even higher than the final. The strength of a constraint is not only reflected by whether or not it matches the true clustering (assuming one exists for a given dataset), but by how much satisfying the constraint increases the strength of the cluster.

One particular application of constrained clustering is the use of ϵ - and δ - constraints, that is, constraints used to reinforce the clustering hypothesis, rather than constraints based on some true clustering. Let C^1 and C^2 be constraints (must-link, cannot-link, or combinatorial) and let F^1 and F^2 be the corresponding penalty functions. Then $C^3 = C^1 \vee C^2$ has the corresponding penalty function $F^3 = F^1 F^2$. Similarly, $C^4 = C^1 \wedge C^2$ has the corresponding

penalty function $F^4 = F^1 + F^2$. Observe that $F^3 = 0$ if and only if C^3 is satisfied, and $F^4 = 0$ if and only if C^4 is satisfied. Finally, observe that any number of must-link and cannot-link constraints can thus be combined in conjunctive normal form by summing products of these penalty functions. As such, these penalty functions can easily be adapted to represent penalty functions for ϵ - and δ -constraints, which are always expressible as a combination of must-link and must-not-link constraints in conjunctive normal form.

	K	MK	PCK	MPCK	$\tilde{\rho}$
liver	1.7349	2.3067	1.7679	1.2516	0.8706
	1.7349	1.8801	1.4568	1.2516	0.8706*
	1.7349	1.6682	1.3542	1.2516	0.8706
pima	1.9995	0.9883	0.8762	0.8681	0.8094
	1.5653	1.9403	1.0585	1.4436	0.8601*
	1.5387	1.9316	1.0585	1.4436	0.8601
faults	0.9392	0.9883	0.8762	0.8681	0.8094
	0.9392	0.9652	0.8762	0.8681	0.8094*
	0.9392	0.9652	0.8637	0.8681	0.8094
wine	1.5126	1.6650	0.8185	1.5393	0.6604
	1.5126	1.5507	0.6542	1.4515	0.6097*
	1.5126	1.4506	0.6101	1.3447	0.4948
iris	0.7373	1.5023	1.4662	0.9612	0.9379
	0.7373	0.9455	0.8877	0.7175	0.6453*
	0.7373	0.7445	0.7041	0.6585	0.5776
iono	2.0706	2.0512	1.6898	1.6898	1.6188
	2.0706	1.8936	1.8936	1.6898	1.6188*
	2.0706	1.8936	1.8919	1.6898	1.6188
glass	3.4599	1.8348	1.0414	1.8284	2.2789
	2.2910	1.4204	1.0414*	1.2820	1.2204
	1.7415	1.0621	1.0414	1.0038	0.2403

Table 3: Table of Davies-Bouldin indices for lowest, median, and highest result over 50 experiments for each of the datasets listed. Note that lower is better. Best median results are marked with an asterisk.

Table 3 also uses datasets taken from the UCI Machine Learning Repository. It tracks the Dunn Index of the best result discovered by each of the listed algorithms and shows the worst, median, and best result over 50 experiments for each one. Here, the ability of the homotopy method to explore all clusterings

between two states, rather than simply striving to satisfy constraints, is clearly demonstrated.

5.3. Finding alternative clusterings

We investigate alternative clustering using the Portrait dataset as studied in [14]. This dataset comprises 324 images of three people each in three different poses and 36 illuminations. Pre-processing involves dimensionality reduction to a grid of 64×49 pixels. The goal of finding two alternative clusterings is to assess whether the natural clustering of the images (by person and by pose) can be recovered. We utilize the same 300 features as used in [14] and setup our framework as shown in Figure 2 (d). Two copies of the dataset are created with one-to-one relationships and we aim to cluster the dataset in a disparate manner.

Table 5 depicts the achieved accuracies on the Portrait dataset using simple K-means, convolutional-EM [14], decorrelated K-means [14] and our framework for disparate clustering. Our algorithm performs better than all other tested algorithms according to both person and pose clusterings.

(a)				(b)			
	C_1	C_2	C_3		C_1	C_2	C_3
C_1	0	0	72	C_1	36	36	36
C_2	63	64	0	C_2	36	36	36
C_3	3	8	114	C_3	36	36	36

Table 4: Contingency tables in analysis of the Portrait dataset. (a) After K-means with random initializations. (b) After disparate clustering.

Method	Person	Pose
k -means	0.65	0.55
Conv-EM [14]	0.69	0.72
Dec- k -means [14]	0.84	0.78
Our framework (disparate)	0.93	0.79

Table 5: Accuracy on the Portrait dataset.

5.4. Scalability and locality preservation

In this section, we consider two synthetic datasets with one (possibly many) relationship between them. The parameters we study are: l_x, l_y , the dimensions of the vectors (varied from 4 to 20), n_x, n_y , the number of vectors (fixed at 100, because as our time complexity analysis shows, they only effect the step of assigning vectors to clusters); $k_x; k_y$, the number of clusters sought (also varied from 4 to 20), and $|B|$, the number of relationships between the datasets (varied from a one-to-one case to about a density of 50%). The vectors were themselves sampled from (two different) mixture-of-Gaussians models.

Figure 8 (a) answers the question of whether our approach yields local clusters as the number of relationships increase (and hence each dataset is more influenced by the other). In this figure, we used settings of 4 and 20 clusters and used our framework to find dependent as well as disparate clusters, and also compared them with K-means (which doesn't use the relationship). Figure 8 (a) shows that even though the K-means algorithm is mining two separate datasets independently, our algorithms achieve very closely comparable results in spite of the co-ordination (dependence or disparate) requirements. Thus, locality of clusters in their respective attribute spaces is not compromised and unvarying with the sparsity of the relationship matrix. At the same time, as Table 6 shows (for the case of four clusters), we achieve the specified contingency table criteria.

Figure 8 ((b),(c)) shows the runtime for our algorithm as a function of attribute vector dimensions (i.e., l_x, l_y) and number of clusters (i.e., k_x, k_y). We vary one parameter, keeping the other fixed (k_x, k_y fixed at 8 versus l_x, l_y fixed at 12). In overall these plots track the complexity analysis presented earlier except for the higher dimension/cluster settings which show steeper increases in time. This can be attributed to the greater number of iterations necessary for convergence in these cases.

Finally, we explore how our results are influenced by the number of clusters, for both dependent as well as disparate clustering formulations (see Figure 8 (d)). As the number of clusters increases, both objective criteria ($\hat{\mathcal{F}}_d$ and $\hat{\mathcal{F}}_i$) become difficult to attain, but for different reasons (recall that the intent of both criteria is to be minimized). In the case of dependent clusters, although the problem gets easier as clusters increase (every point can become its own cluster), the objective function scores get lower due to our regularization as explained in Section 4. In the case of disparate clusters, as the number of clusters increases, the size of the contingency table increases quadratically with the number of samples staying constant. As a result, it becomes difficult to

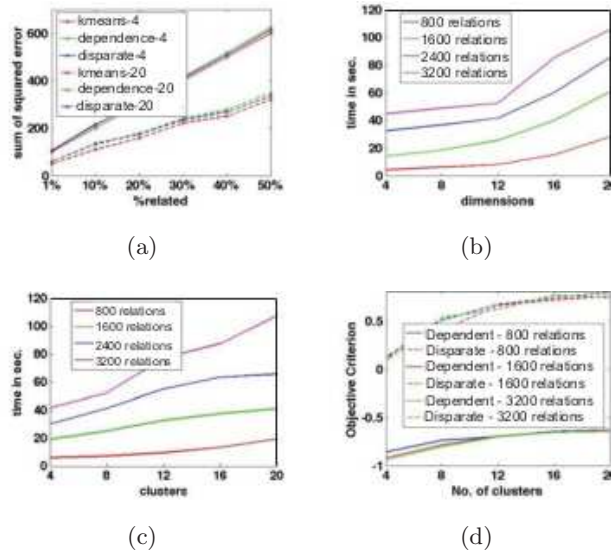


Figure 8: Synthetic data results. (a) Comparisons of SSE measures with K-means. (b, c) Time as number of attribute dimensions (b) or clusters (c) is increased. (d) Objective criterion as a function of the number of clusters for both dependent and disparate schemas of clustering.

distribute the samples across the contingency table entries without introducing some level of dependence (i.e., some entries must be zero implying dependence).

5.5. Comparing gene expression programs across yeast, worm, and human

In this study, we focus on time series gene expression profiles collected over heat shock experiments done on organisms of varying complexity: *H*: human cells (4 time points), *Y*: yeast (8 time points), and *W*: *C. elegans* (worm; 7 time points). We also gathered many-many (top-k) ortholog information between the three species. A typical goal in multispecies modeling is to identify both conserved gene expression programs as well as differentiated gene expression programs. The former is useful for studying core metabolism and stress response processes, whereas the latter is useful for studying species-specific functions (e.g., the yeast is more tolerant to desiccation stress, but the worm is the more complex eukaryote).

7	6	4	11
1	3	10	2
9	10	3	8
4	9	2	11

17	3	2	2
4	18	2	3
4	3	20	3
2	0	1	16

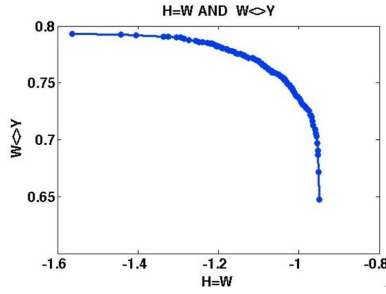
6	5	7	5
8	7	5	5
6	6	7	4
7	8	7	7

Table 6: Our approach helps drive a K-means cluster assignment (top) toward either dependent (bottom-left) or disparate (bottom-right) sets of clusters.

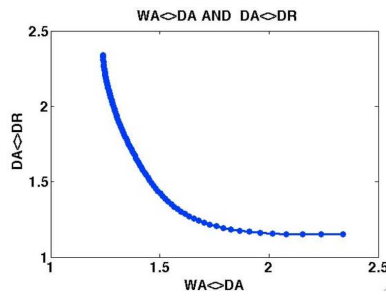
First we study a 3-way clustering setup with only two constraints, namely that clusters in H and W must be dependent, denoted by $H = W$, and that clusters in W and Y must be disparate, denoted by $W \langle \rangle Y$. See Figure 9 (top). As the balance between these criteria is varied from one extreme to another (via the convex combination formulation), this curve traces out the objective function values. The top left corner is the point where complete emphasis is placed on achieving the $H = W$ criterion (conversely for the bottom right corner). As we seek to improve the other criteria, note that we might (and will) sacrifice the already achieved criterion. The point of maximum curvature on this plot gives a ‘sweet spot’ so that any movement away from the sweet spot would cause a dramatic change in the objective function values. A qualitatively different type of plot is shown in Figure 9 (middle) for the case study described in the next section, but here again the point of maximum curvature reveals a balancing threshold of the two criteria. A 3-way clustering setup with three constraints is described in Figure 10 and its corresponding tradeoff plot is in Figure 9 (bottom). Here there are likely multiple points of interest depending on which criteria are sacrificed in favor of others.

5.6. Multiorganismal and multistress modeling case study

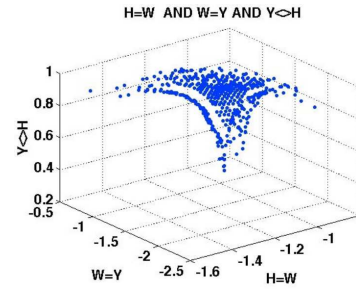
Finally, we present a case study that has a diversity of both organisms and stresses. To capture process-level similarities and differences, the data vectors we cluster here correspond to Gene Ontology categories rather than individual gene expression profiles. We used three time series datasets: WA—*C. elegans* aging (7 time points), DA—*D. melanogaster* aging (7 time points) and DR—



(a)



(b)



(c)

Figure 9: Balancing objectives in multicriteria clustering optimization. Points of maximum curvature on these plots reveal a balancing point between the conflicting criteria.

D. melanogaster caloric restriction (9 time points). Observe that the first two datasets share a similarity of process whereas the latter two share a similarity of organism. In a sense, the *D. melanogaster* aging dataset is squarely in the “middle” of the other two datasets. When subjected to clustering together, the inherent tradeoff is what we seek to capture.

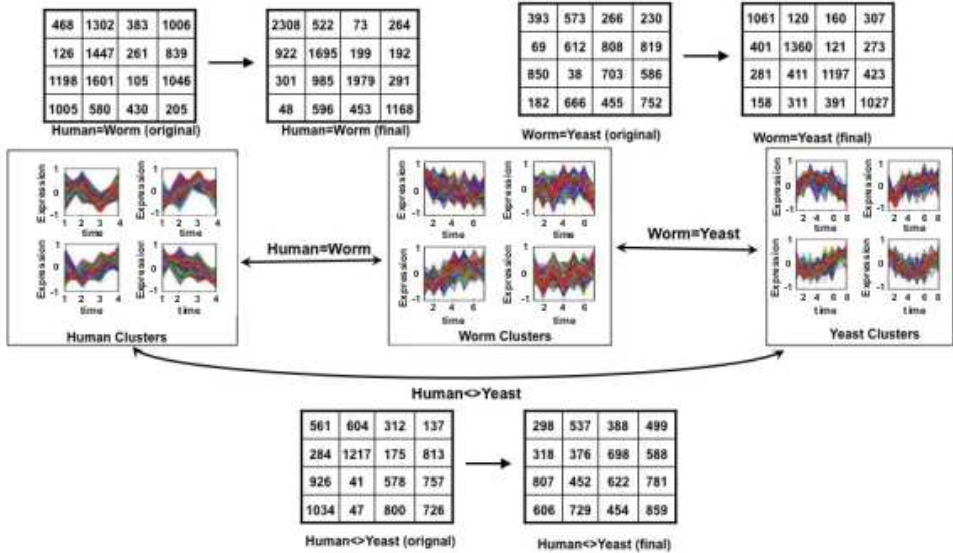


Figure 10: Clustering three datasets with three constraints between them. Two sets of clusters (between human/worm and between worm/yeast) are constrained to be similar whereas the third set (between human/yeast) is constrained to be dissimilar. Observe how the top two contingency tables are driven toward diagonal dominance whereas the bottom contingency table is driven toward a uniform distribution.

For this evaluation, we studied the enrichment of clusters obtained from our framework vis-a-vis K-means clustering. We set the number of clusters at 7 and evaluated GO terms for an FDR-corrected q -value of 0.05. First, we study the clustering setup so that $DA = DR$ and $WA = DA$, for a setting of $\alpha = 0$ (so that more emphasis is placed on achieving the dependent clustering $DA = DR$). Here, we observed 75 GO terms enriched (versus 37 for K-means). Similar improvements were seen for $\alpha = 0.5$ (55 versus 20) and for $\alpha = 1$ (89 versus 35). Observe the greater numbers of terms enriched in general for the extremes (which is to be expected). In terms of process-level similarities, the GO terms common across the aging datasets but which do not appear when we emphasize organism-level similarities are neuron recognition, embryonic pattern specification, aromatic compound catabolic process, somatic sex determination, and sulfur compound biosynthetic process.

The organism-level similarities are captured in chemo-sensory behavior, cell redox homeostasis, peptide metabolic process, regulation of cell proliferation, anatomical structure formation, and negative regulation of growth.

These results show that process-level similarities involve higher order functions whereas organism-level similarities involve growth and metabolism processes. The careful interplay between aging and caloric restriction, both at the organismal and at the interorganismal level, is an interesting conclusion from this study.

6. Related work

MDI: The objective functions defined here have connections to the principle of minimum discrimination information (MDI), introduced by Kullback for the analysis of contingency tables [15] (the minimum Bregman information (MBI) in [16] can be seen as a generalization of this principle). The MDI principle states that if q is the assumed or true distribution, the estimated distribution p must be chosen such that $D_{KL}(p||q)$ is minimized. In our objective functions the estimated distribution p is obtained from the contingency table counts. The true distribution q is always assumed to be the uniform distribution. We maximize or minimize the KL-divergence from this true distribution as required. Space restrictions prevent us from describing the connection to MDI in further detail.

Co-clustering binary matrices, cross-associations, and associative clustering: Identifying clusterings over a relation (i.e., a binary matrix) is the topic of many efforts [17, 18]. The former makes use of information-theoretic criteria to best approximate a joint distribution of two binary variables and the latter uses the MDL (minimum description length) principle to obtain a parameter-less algorithm by automatically determining the number of clusters. Our work is focused on not just binary relations but also attribute-valued vectors. The idea of comparing clustering results using contingency tables was first done in [19] although our work is the first to unify dependent and disparate clusterings in the same framework.

Finding disparate clusterings: The idea of finding disparate clusterings has been studied in [14]. Here only one dataset is considered and two dissimilar clusterings are sought simultaneously where the definition of dissimilarity is in terms of orthogonality of the two sets of basis vectors. This is an indirect way to capture dissimilarity whereas in our paper we use contingency tables to more directly capture the dissimilarity. Furthermore, our work enables the

combination of similar clusterings and disparate clusterings in a more expressive way. For instance, given just two datasets \mathcal{X} and \mathcal{Y} with two relationships R_1 and R_2 between them, our work can identify clusters in \mathcal{X} and \mathcal{Y} that are similar from the perspective of R_1 but dissimilar from the perspective of R_2 : it is difficult to specify such criteria in terms of the basis vectors since they will be the same irrespective of the relationship.

Clustering over relation graphs: Clustering over relation graphs is a powerful framework by Banerjee et al. [2] that uses the notion of Bregman divergences to unify a variety of loss functions and applies the Bregman information principle (from [16]) to preserve various summary statistics defined over parts of the relational schema. The key difference between this work and ours is that this framework is primarily targeted toward dependent clustering (compression) whereas our work targets both dependent and disparate clustering, over different parts of the relational schema as appropriate.

Multivariate information bottleneck: Our work is reminiscent of the multivariate information bottleneck (MIB) [20], a framework for specifying clusterings in terms of two conflicting criteria: compression (of vectors into clusters) and preservation of mutual information (of clusters with auxiliary variables that are related to the original vectors). We share with MIB the formulation of a multicriteria objective function derived from a clustering schema but differ in the specifics of both the intent of the objective function and how the clustering is driven based on the objective function. Furthermore, the MIB framework was originally defined for discrete settings whereas we support a mixed modality of datasets.

7. Conclusion

We have presented a very general and expressive framework for clustering non-homogeneous datasets. We have also shown how it subsumes many previously defined formulations and that it sheds useful insights into tradeoffs underlying complex relationships between datasets, especially when those tradeoffs are systematically explored with a homotopy algorithm. The mathematical theory behind a novel application of probability-one homotopy algorithms to constrained clustering was also given.

Our directions for future work are three fold. Thus far, we have used distinct relations to enforce disparate and dependent clusterings. One of the first directions for future work is to allow both types of clusterings to be captured in the same relation. This would help capture more expressive relationships

between datasets, such as a banded diagonal structure in the contingency table. Secondly, just as the theory of functional and multivalued dependencies (FDs and MDs) helps model relations in and between individual tuples, we aim to develop a theory of ‘clustering dependencies’ that can help model relations in the aggregate, e.g., between clusters. Thirdly, how to generalize homotopy algorithms to apply to multicriteria machine learning problems (more than two datasets and more than two relations) should be investigated.

References

- [1] B. Long, X. Wu, Z. Zhang, P.S. Yu, Unsupervised learning on k -partite graphs, In: *Proc. KDD '06* (2006), 317-326.
- [2] A. Banerjee, S. Basu, S. Merugu, Multi-way clustering on relation graphs, In: *Proc. SDM '07* (2007), 225-334.
- [3] E. Bae, J. Bailey, COALA: A novel approach for the extraction of alternate clustering of high quality and high dissimilarity, In: *Proc. ICDM '06* (2006), 53-62.
- [4] Z. Qi, I. Davidson, A principled and flexible framework for finding alternative clusterings, In: *Proc. KDD '09* (2009), 717-726.
- [5] M.S. Hossain, S. Tadepalli, L.T. Watson, I. Davidson, R.F. Helm, N. Ramakrishnan, Unifying dependent clustering and disparate clustering for nonhomogeneous data, In: *Proc. 16th ACM SIGKDD Conf. on Knowledge Discovery and Data Mining* (2010), 593-602.
- [6] G. Kreisselmeier, R. Steinhauser, Systematic control design by optimizing a vector performance index, In: *Proc. IFAC Symp. on Comp. Aided Design of Control Systems* (1979), 113-117.
- [7] A.R. Conn, N.I.M. Gould, P.L. Toint, *LANCELOT: A Fortran Package for Large-scale Nonlinear Optimization (Release A)*, Volume 17, Springer Verlag (1992).
- [8] S.S. Tadepalli, *Schemas of Clustering*, Virginia Tech. (2009).
- [9] I. Davidson, S.S. Ravi, Clustering with constraints: feasibility issues and the k-means algorithm, In: *Proc. SDM '05* (2005), 201-211.

- [10] K. Wagstaff, C. Cardie, S. Rogers, S. Schrödl, Constrained K -means clustering and background knowledge, In: *Proc. ICML '01* (2001), 577-584.
- [11] D.R. Easterling, *Solution of Constrained Clustering Problems through Homotopy Tracking*, Virginia Tech. (2014).
- [12] L.T. Watson, S.C. Billups, A.P. Morgan, Algorithm 652: HOMPACT: A suite of codes for globally convergent homotopy algorithms, *ACM Trans. Math. Software*, **13** (1987), 281-310.
- [13] M. Bilenko, S. Basu, R.J. Mooney, Integrating constraints and metric learning in semi-supervised clustering, In: *Proc. ICML '04* (2004), 11-18.
- [14] P. Jain, R. Meka, I.S. Dhillon, Simultaneous unsupervised learning of disparate clusterings, In: *Proc. SDM '08* (2008), 858-869.
- [15] S. Kullback, D.V. Gokhale, *The Information in Contingency Tables*, Marcel Dekker, Inc. (1978).
- [16] A. Banerjee, S. Merugu, I.S. Dhillon, J. Ghosh, Clustering with Bregman Divergences, *J. of Machine Learning Research*, **6** (2005), 1705-1749.
- [17] D. Chakrabarti, S. Papadimitriou, D.S. Modha, C. Faloutsos, Fully automatic cross-associations, In: *Proc. KDD '04* (2004), 79-88.
- [18] I.S. Dhillon, S. Mallela, D.S. Modha, Information theoretic co-clustering, In: *Proc. KDD '03* (2003), 89-98.
- [19] S. Kaski, J. Nikkilä, J. Sinkkonen, L. Lahti, J.E.A. Knuuttilla, C. Roos, Associative clustering for exploring dependencies between function genomics data sets, *IEEE/ACM TCBB*, **2(3)** (2005), 203-216.
- [20] N. Friedman, O. Mosenzon, N. Slonim, N. Tishby, Multivariate information bottleneck, In: *Proc. UAI '01* (2001), 152-161.