A Comparison of Spectral Clustering Algorithms

Deepak Verma
Department of CSE
University of Washington
Seattle, WA 98195-2350
deepak@cs.washington.edu

Marina Meilă Department of Statistics University of Washington Seattle, WA 98195-4322 mmp@stat.washington.edu

Abstract

Spectral Clustering has become quite popular over the last few years and several new algorithms have been published. In this paper, we compare several of the best-known algorithms from the point of view of clustering quality over artificial and real datasets. We implement many variations of the existing spectral algorithms and compare their performance to see which features are more important. We also demonstrate that spectral methods show competitive performance on real dataset with respect to existing methods.

1 Introduction

Clustering has always been a hard problem and an active topic of research. Recently, a new approach has started to get a lot of attention namely spectral methods. The spectral methods for clustering usually involve taking the top eigen vectors of some matrix based on the distance between points (or other properties) and then using them to cluster the various points.

Spectral clustering techniques have seen an explosive development and proliferation over the past few years. They promise to become strong competitors for other clustering methods. Several successes have already been registered (LSA,[9]). Spectral methods are attractive because they are easy to implement and are reasonably fast (for *sparse* data sets up to several thousands). Also they do not intrinsically suffer from the problem of local optima. (Though depending on the exact algorithm some local optima might be there.)

In spite of the large number of papers on spectral clustering, so far no systematic comparison between the existing algorithms has been published. This is what we set out to do here. We intend to take a look at four spectral algorithms and take them apart. We generate a list of algorithms which are made from different parts of different algorithms and compare their performance. We hope to be able to find out which of these sub-components are important and which not, and to see if some combination of these works the best.

2 Spectral Clustering Algorithms

The algorithms we selected are at this date among the most popular of the published ones. We have also aimed at representing a diverse set of algorithmic features. The algorithms are: (1) the image segmentation algorithm introduced by Shi and Malik (SM) [9], (2)A variant by Kannan, Vempala and Vetta (KVV) [2], (3) the algorithm of Ng, Jordan and Weiss (NJW) [8], (4) an algorithm suggested by Meilă and Shi (Multicut) [5]. We also present the results obtained with the Single and Ward linkage algorithms (denoted by MST, Ward) as a "strawman" in order to demonstrate that the clustering tasks that we chose are not trivial.

Before describing the algorithms in more detail, we introduce some notation. Then we would describe the four spectral and two grouping algorithms.

2.1 Notation

The set of data points to be clustered will be denoted by I, with |I| = n. For each pair of points $i, j \in I$ a similarity $S_{ij} = S_{ji} \ge 0$ is given. The similarities S_{ij} can be viewed as weights on the undirected edges ij of a graph G over I. The matrix $S = [S_{ij}]$ plays the role of a "real-valued" adjacency matrix for G. Let $D_i = \sum_{j \in I} S_{ij}$ be called the degree of node i, and the volume of a set $A \subset I$ be $Vol\ A = \sum_{i \in A} D_i$. The set of edges between two disjoint sets $A, B \subseteq I$ is called the edge cut or in short the cut between A, B.

A clustering $C = \{C_1, C_2, \dots C_K\}$ is a partitioning of I into the nonempty mutually disjoint subsets $C_1, \dots C_K$. In the graph theoretical paradigm a clustering represents a multiway cut in the graph G.

All the algorithms that we use here just need a similarity matrix between points (except anchor. But we never us it on the original (unmapped) points). So all we need is a data with similarity matrix and there may not be an actual set of distinct points in the initial domain or the points may even come from an infinite dimensional space. (something an output or an kernel function).

2.2 Overview

The algorithms presented here can be thought of as consisting of 3 stages:

- **Preprocessing:** This is a form of normalization of the similarity matrix S. We did some smoothing initially to make sure that matrix is not too ill conditioned. However to make results more relevant w.r.t. to other papers, we eventually dropped the smoothing.
- **Spectral Mapping:** Some eigenvectors of the preprocessed similarity matrix are computed. Each data point *i* is mapped to a tuple representing the values of component *i* in the aforementioned eigenvectors.
- **Postprocessing/Grouping:** A (usually simple) grouping algorithm clusters the data (in the original or spectral domain).

There are three kind of algorithms presented here.

- Recursive Spectral: These algorithms try to split the data into two partitions based on a single eigenvector and are then are recursively used to generate more partitions.
- Multiway Spectral: These use more information in multiple eigenvectors to do a direct multiway partition of data.
- Non spectral: A (usually simple) grouping algorithm that clusters the data quickly. These are used in conjunction with the Multiway spectral algorithms and also provide a baseline for performance.

For some of the algorithms (SM, KVV) heuristic methods for finding the number of clusters K have been suggested [9, 2]. In this work, to provide for a fair comparison, we have assumed for all the algorithms that the number of clusters K is given in advance.

2.3 The Shi and Malik (SM) algorithm

This algorithm was introduced by [9] as a heuristic algorithm aimed to minimize the *Normalized Cut* criterion proposed by the same authors. The normalized cut between two sets $A, B \subseteq I$ is defined as

$$NCut(A, B) = Cut(A, B) \left(\frac{1}{Vol A} + \frac{1}{Vol B} \right)$$

According to [9] the set I is partitioned into two clusters C, $C' = I \setminus C$ that minimize NCut(C, C') over all possible two way partitions of I. This problem is provably NP-hard, but the authors show that under certain special conditions a spectral algorithm exists that finds the optimum.

Algorithm SM

1. Compute

$$P = D^{-1}S \tag{1}$$

- 2. Let $1 = \lambda_1 \ge \lambda_2 \ge \ldots \ge \lambda_n$ be the eigenvalues of P and v^1, v^2, \ldots, v^n the corresponding eigenvectors¹ Compute v^2 .
- 3. Min-Ratio-Cut
 - (a) Sort the elements of v^2 in increasing order. Denote by v_i^2 the *i*-th element in the sorted list.
 - (b) For i = 1, ..., n 1Compute $NCut(C_i, C'_i)$ where $C = \{1, ..., i\}, C' = \{i + 1, i + 2 ..., n\}$
 - (c) Partition I into the two clusters C_{i_0}, C'_{i_0} where $i_0 = \min NCut(C_i, C'_i)$
- 4. Repeat steps 1–3 recursively on the cluster with the largest λ_2 until K clusters are obtained.

MIN-RATIO-CUT is inspired by [10] which suggests that a line search along the vector produces better cuts than the original heuristic. This algorithm is almost identical with spectral algorithm II in [2]. That is discussed in more detail in 2.4. Based on that another version for this algorithm in implemented using the conductance as a criterion in Min-Ratio-Cut.

2.4 The Kannan, Vempala and Vetta Algorithm (KVV)

The KVV is very similar to to SM algorithm with two differences. One difference is in step 3, where the optimal cut is found with respect to the Cheeger conductance $\phi(C_i, C'_i)$, another measure of cut quality. The conductance of a clustering $\{C, I \setminus C\}$ is defined as

$$\phi(C, I \setminus C) = \frac{Cut(C, I \setminus C)}{\min Vol C, Vol I \setminus C)}$$

Also in the recursive step the kvv variant decides the next cluster is the one with the minimum conductance. This variant of step 3 will be called Min-Conductance.

Another difference is in the "normalization" past the first iteration of step 1. The SM algorithm just takes the block of S corresponding to the current cluster. The variant in [2] always uses blocks from the P computed at the first iteration. To ensure that the row sums of the blocks equal 1, the diagonal elements P_{ii} of the current block are adjusted. Thus, this variant ignores the self-similarity of the data points in all but the first iteration.

To adjust the P_{ii} there are two possibilities. We can either scale up all the entries in row to sum up to one or add the extra weight to the diagonal element. We call the first variant kvv_mult and the second kvv_add.

2.5 The Ng, Jordan and Weiss (NJW) algorithm

Algorithm NJW

- 1. Set the diagonal elements S_{ii} to 0.
- 2. Compute the matrix

$$L = D^{-\frac{1}{2}}SD^{-\frac{1}{2}} \tag{2}$$

- 3. Let $1 = \mu_1 \ge \mu_2 \ge \ldots \ge \mu_K$ be the K largest eigenvalues of L and u^1, u^2, \ldots, u^K the corresponding eigenvectors All eigenvectors are normalized to have unit length. Form the matrix $U = [u^1 \ u^2 \ \ldots \ u^K]$ by stacking the eigenvectors in columns.
- 4. Form the matrix Y from U by renormalizing each of U's rows to have unit length (i.e $Y_{ij} = U_{ij} / \sqrt{\sum_{j} U_{ij}^2}$).
- 5. K-Means-Orthogonal Treating each row of Y as a point in K dimensions, cluster them by the K-means algorithm to obtain the final clustering. The K-means algorithm is initialized by the Orthogonal-Initialization method described in [8].

¹If the eigenvalues are not distinct, pick the eigenvectors such that $v^{i}^{T}Dv^{j}=0$ for $i\neq j$. This is always possible and the Matlab implementation that we used does it automatically.

² If the eigenvalues are not distinct, choose u^k 's that are orthogonal to each other. L is related to the *Laplacian* of S. See e.g [9] for details.

In [8] it is proved that if the clusters are well separated in the sense that the similarity matrix S is almost block diagonal, and if the sizes of the clusters and degrees of individual nodes don't vary too much, the rows of the Y matrix cluster near K orthonormal vectors in R^K . This fact suggested the orthogonal initialization method.

We implemented a slightly different version of this algorithm which we think has greater numerical stability. The details are in the appendix.

2.6 The Meila-Shi algorithm

This algorithm was suggested in [6]. Algorithm MULTICUT

- 1. Compute the stochastic matrix P as in (1).
- 2. Compute $v^1, \ldots v^K$ the eigenvectors of P corresponding to the K largest eigenvalues. Form the matrix V whose columns are $v^1, \ldots v^K$.
- 3. Cluster the rows of V as points in a K-dimensional space.

2.7 Anchor Algorithm

A "flat" version of **the Anchor algorithm** of [7], also very related to the minimum diameter clustering method of [1].

Algorithm ANCHOR

- 1. Choose a point at random. Set k = 1, k' = 0. Choose anchor x_1 to be the farthest point from the initial point.
- 2. Construct C_k the cluster associated with x_k as a list of points that are closer to x_k than to any other anchor. The list is sorted by decreasing distance from x_k .
- 3. Test if x_k has enough points. If $|C_k| < n_{min}$ then k' = k' + 1.
- 4. Set k = k + 1. Choose anchor x_k to be the farthest point from all existing anchors.
- 5. If k k' < K go to step 2.

Note that the algorithm may produce more than K clusters. It is also possible that it never produces K clusters having more than n_{min} points. In our experiments, all performed with $n_{min} = 3$, the latter was never observed.

2.8 Linkage algorithms

These are the hierarchical clustering algorithm which work on distances between the points. Since we have similarities, we need to choose a way of mapping similarities to distances. We chose to use the inverse of similarity as the measure (A very small value was added it to similarity to ensure that inverse of zero is not taken).

We used two methods: single linkage and ward linkage. Single linkage is same as performing the MST on the *dissimilarities* graph of the points and ward linkage is similar except that distance is the inner square distance. For more details on the ward algorithms see [11].

As it will be shown in the following sections, the different algorithms are much closer then they initially appear. Both the experiments and the theory suggest that the differences between algorithms will depend strongly on the quality of the postprocessing step. Therefore, we experimented with several different clustering methods that will be described here.

3 Theoretical Results

In this section we would present some theoretical results concerning the various spectral algorithms The main result is the (near) equivalence of the two (non recursive) spectral methods (NJW, Multicut) First we present a modification to the two algorithms to make them numerically more stable and then we present he conditions and proof of similarity.

3.1 A modification to the NJW and Multicut method

As discussed in section 2.5, NJW use the top k eigenvectors of $L = D^{-\frac{1}{2}}SD^{-\frac{1}{2}}$ to map data. We propose using the the top K eigen vectors of the generalized eigen system $Sx = \lambda Dx$. We believe this is a numerically more stable method to compute the eigen vectors as there is no division by D involved, which could contain very small values. (An outlier for example would have a very small degree as it is not close to another point and self similarity is defined to be zero in [8]).

To prove why this gives the same vector Y consider the following:

$$D^{-\frac{1}{2}}SD^{-\frac{1}{2}}u = \lambda u, \text{ Premultiplying by } D^{-\frac{1}{2}} \text{ we get}$$
 (3)

$$D^{-1}S(D^{-\frac{1}{2}}x) = \lambda(D^{-\frac{1}{2}}x)$$
 Putting $v = D^{-\frac{1}{2}}x$ (hence $x = D^{\frac{1}{2}}v$) we get (4)

$$D^{-1}Sv = \lambda v \text{ and hence } Sv = \lambda Dv.$$
 (5)

Now consider the top k eigen vectors $U_{n\times k}$ and $V_{n\times k}$. $U=D^{\frac{1}{2}}V$. (U,V) have the same meaning as in the previous section. Recall also that the rows of Y are unit vectors in K dimensions.) Since D is a diagonal matrix this means that the i^{th} row of U is same as the i^{th} row of V scaled by $D_{ii}^{\frac{1}{2}}$. So after the row is normalized to length 1, the Y obtained from V is identical to the Y obtained from U.

For the multicut method observe that $P = D^{-1}S$ so the generalized eigen value system $Sv = \lambda Dv$ is mathematically equivalent and numerically more stable than the computing P and its eigenvectors.

3.2 Preliminaries: The perfect S

When each of the clusters in the postprocessing step of a spectral algorithm is reduced to a distinct point we say that S is *perfect* for the respective algorithm. For example, a block diagonal similarity matrix is perfect for all algorithms. A perfect S represents the ideal situation for a clustering algorithm. Here we essentially show that when S is such that the resulting P is block-stochastic, a term that will defined below, then S is perfect for NJW and Multicut and should give good performance on recursive algorithms.

Definition 1 Let P be a stochastic matrix. We say P is block stochastic w.r.t a clustering $\Delta = \{C_1, \ldots C_K\}$ iff $\sum_{i \in C_k} P_{ij}$ has the same value for all points $i \in C_{k'}$ for all $k, k' = 1, \ldots K$.

Definition 2 A vector $v = [v_1, v_2, \dots, v_n]^T$ is piecewise constant (PC) w.r.t a clustering iff $v_i = v_j$ whenever i, j are in the same cluster.

Proposition 3 A matrix $A_{n \times k}$ with orthonormal columns and atmost k unique rows would have exactly k unique orthogonal rows.

Proof: First of since A has rank k it has to have k independent rows, which means that it has exactly k unique rows.

Now, Rearrange the rows so that the identical rows are next to each other. (This does not affect the result so we assume that A has this property). So now A can written as $A = C_{n \times k} B_{k \times k}$ where B is the matrix with the k unique rows and $C_{ij} = 1$ if the i^{th} row of A is same as the j^{th} row of B. We just need to prove that B has orthogonal rows.

Since the columns of A are orthonormal, $A^TA = I_{k \times k}$. This implies $B^TC^TCB = I$, Now it is easy to see that C^TC is a diagonal matrix (say D). Define $Z = D^{\frac{1}{2}}B$. This gives us $Z^TZ = I$ which means that Z is orthonormal with orthogonal rows. Which proves that B has orthogonal rows. (Premultiplying a matrix with a diagonal matrix just scales its rows). QED.

Consider the algorithms NJW and Multicut. They use the top K eigen vectors which are orthonormal. So if these eigenvectors are piecewise constant w.r.t to a clustering Δ then they have atmost K unique rows and hence would need have exactly K unique rows which would make the S perfect. So all we need is for the eigenvectors to be PC w.r.t Δ . Also note that these rows would be orthogonal.

Here is a key result from [6].

Proposition 4 Let P be a stochastic matrix with rows and columns indexed by I that has independent eigenvectors. Let $\Delta = \{C_1, C_2, \dots C_K\}$ be a partition of I. Then, P has K PCE w.r.t. Δ corresponding to non-zero eigenvalues if and only if the sums $P_{ik} = \sum_{j \in C_k} P_{ij}$ are constant for all $i \in C_{k'}$ and all $k, k' = 1, \dots K$ (i.e P is block stochastic) and the matrix $R = [P_{kk'}]_{k,k'=1,\dots K}$ (with $R_{kk'} = \sum_{j \in C_k'} P_{ij}$, $i \in C_k$) is non-singular

It is easy to see that when S produces a block stochastic P then S is perfect for the Multicut algorithm. The case of the NJW algorithm is characterized in the next section.

3.3 Equivalence Results

In this section we show a kind of equivalence between NJW and Multicut which characterizes the ideal case for NJW as well.

. We assume that P and L are obtained from the same symmetric S with non-negative elements by (1,2) and that U,V have the same meaning as in the section 3.1. Whenever we mention clustering it would assumed that it stand for a clustering $\Delta = \{C_1, C_2, \dots, C_K\}$ Recall also that the rows of Y are unit vectors in K dimensions.

Proposition 5 Let v^1, v^2, \ldots, v^K be the top K eigen vectors of P such that they are PC w.r.t a clustering. Then the rows of Y will be grouped in K groups of identical vectors $\tilde{y}_1, \tilde{y}_2, \ldots, \tilde{y}_K$ in R^K such that the groups are distinct and orthonormal.

Proof: For this proposition and next one define y_i to be the i^{th} row of Y and x_i to be the i^{th} row of V. By the results in the section 3.1, $y_i = x_i/|x_i|$ i.e. y_i is just x_i scaled to be on the unit sphere.

Let V be PCE w.r.t. Δ . Now consider arbitrary $i, j \in C_s$ where s is arbitrarily chosen from $1, 2, \ldots k$. Since V is PCE w.r.t Δ , $x_i = x_j = \tilde{x}_s(\text{Let})$. This implies $y_i = y_j = \tilde{y}_s$. Since \tilde{x}_i 's are orthogonal so are \tilde{y}_i (and distinct). Also since they lie on the unit sphere they are orthonormal.

Proposition 6 If the rows of Y are equal for all the points in a cluster, then P has K PCE w.r.t. Δ .

Proof: This is the reverse direction of the above proposition and somewhat more tricky to prove. Assume that rows of Y are equal w.r.t to Δ . Now consider $arbitrary\ i,j\in C_s$ where s is arbitrarily chosen from $1,2,\ldots k$. By the assumption, $y_i=y_j=\tilde{y}_s(\text{Let})$. This implies $x_i=|x_i|y_i=|x_i|\tilde{y}_s$ and $x_j=|x_j|y_j=|x_j|\tilde{y}_s$. In particular the first column of the x_i and x_j (which are also in first column of V) are proportional to $|x_i|$ and $|x_j|$. i.e. $V_{i1}=|x_i|\tilde{y}_s^1$ and $V_{j1}=|x_j|\tilde{y}_s^1$.

But the first column of V is also the first eigenvector of $P=D^{-1}S$ which is stochastic matrix. So $V_{l1}=V_{m1}$ for all l,m in $1,2,\ldots n$. In particular $V_{i1}=V_{j1}$. This implies $|x_i|=|x_j|=|\tilde{x}_s|(\text{Let})$. So $x_i=x_j=|\tilde{x}_s|\tilde{y}_s$. Since x_i and x_j are rows of $V=[v^1,\ldots v^K]$, this implies that $[v^1,\ldots v^K]$ are PC and P has PCE w.r.t Δ .

This gives us the following Theorem

Theorem 7 A similarity matrix with a block stochastic P is perfect for NJW and Multicut algorithms. All the points in the same cluster in the original domain would be mapped to to single point in the spectral domain with points corresponding to different clusters orthogonal to each other.

In other words, S is perfect for the NJW algorithm iff it is perfect for the Multicut algorithm as well.

Therefore, our experiments with the artificial data will focus less on the ideal case that can be studied theoretically and more on robustness in noise and real data. Note that in spite of the this similarity there is a difference in the spectral domain (the scale of the vectors) which could be important. So we still implement two separate versions for NJW and Multicut.

3.4 Theoretical Comparison

In the previous section we proved that a block stochastic S is perfect for NJW and Multicut. What about the recursive spectral algorithms?

Take the case when we are splitting the (potentially subset of original) points into two clusters. (At top level or in any of the recursive step). Let the P that we have (for these points) is block stochastic w.r.t Δ . In that case the second eigen vector would be piece wise constant w.r.t. to the Δ so when reorder the points based on this eigen vector we would have the points in the right order. That is, if we chose the right point to partition the two partitions, none of the clusters in Δ would be split. This would also make the P for these two split cluster would be block stochastic (w.r.t to the two parts of Δ) setting the optimal stage for the recursive sub steps.

The criterion that the two algorithms choose do not ensure that this optimal point would be optimally chosen but still should do a good job. In fact if P is block diagonal, then both the algorithms would choose the right position as Ncut and conductance is minimized (zero) at only these positions.

One possibility is to choose the point which maximizes the gap in the v^2 vector. We intend to explore this as future work.

4 Datasets

As discussed above there is a whole lot of spectral methods each with its little variation that needs to compared to each other. Since it is not possible to visually compare them we needed datasets which are "pre-classified" So that it is possible to compute the clustering error and the VI w.r.t. true clustering. (see section 5) We used both artificial and real datasets as described below. The artificial datasets were primarily used to demonstrate the robustness of algorithms to noise.

4.1 S100: A Block Stochastic Matrix

This is the "ideal" case for the multicut and ang-based spectral algorithms. Here we constructed an 100×100 matrix which consists of 5 clusters. There are five clusters present of size 10,20,30,20 and 20 respectively. The similarity matrix is also slightly block diagonal. The purpose of using this dataset was to demonstrate the stability of the algorithms w.r.t to noise. This data file is called block-diagonal-hard.

4.2 Handwritten digits

This is the set of optical handwritten digit recognition that is available in the NIST site. There are lots of version available with different preprocessing. In particular we used the data set and preprocessing as mentioned in [3]

Here is the description of preprocessing done by . They used preprocessing programs made available by NIST to extract normalized bitmaps of handwritten digits from a preprinted form. From a total of 43 people, 30 contributed to the training set and different 13 to the test set. 32x32 bitmaps are divided into non overlapping blocks of 4x4 and the number of on pixels are counted in each block. This generates an input matrix of 8x8 where each element is an integer in the range 0..16. This reduces dimensionality and gives invariance to small distortions.

We further down sampled the dataset to 100 elements per digit giving a total of thousand 64 dimensional points and 10 clusters. We call this dataset digit1000. Some of digits were more easy to distinguish from another and in particular digits 0,2,4,6,7 were a lot more easier to distinguish than the others. So we created another dataset containing the hundred instances of just these five digits. We called this data set digitFive1000. (The 1000 just to remind that this is the same the digit1000 database.

4.3 Gene Expression Data

DNA microarrays provide a way to the biologists to study the variation of many genes together. Using these has been a plethora of gene expression data generated in the community. This has led to the great need for the data to be analyzed. ([12]). We used one such dataset, the yeast cell cycle data, which is publically available at [13]. It shows the fluctuation of the gene expression levels of over 6000 genes over the two cellcycles (which has 17 time points). The dataset is restricted to the 384 genes who's expression level peak

at different points corresponding to the five phases of the cell cycle. The objective given these expression levels is to be able to cluster them into clusters corresponding to the five phases.

There are two kinds of pre processing that are suggested in [12]. First is to the take the logarithm of the expression level and second to "standardize" the mean to be zero and variance 1. These data transformations were done so as to make the data fit better to the gaussian model (They were using mixture models to cluster the data. See [12] for more details). We call the first dataset cellcycle and the second cellcycle-std.

5 Evaluating Clustering Performance

Measuring a clustering performance in general is a very hard problem. The notion of good clustering is intrinsically tied with the definition of what a cluster is which in itself is a big research topic. In our case measuring clustering performance is easier as in all the datasets we have the "true" clustering available. Given that the clustering performance is just a measure of how "different" is the clustering produced w.r.t the true clustering. There are three kind of measures that we used: Clustering Error and Variation of Information . In this section \mathcal{C}^{true} would represent the true (given) clustering and \mathcal{C} the clustering produced by the clustering algorithm.

5.1 Clustering Error

Clustering error is defined as the number of "misclassification" This the error induced in the clustering w.r.t. to the true clustering.

Let Confusion be the confusion matrix of two clusterings. $(Confusion(k_{true}, k) = |\mathcal{C}_{k_{true}}^{true} \cap \mathcal{C}_{k}|$ i.e. number of points x that are cluster k_{true} in true clustering and cluster k in the clustering produced. Then

$$CE(\mathcal{C}, \mathcal{C}^{true}) = \left(\sum_{k_{true}} \sum_{k \neq k_{true}} Confusion(k_{true}, k)\right) / n \tag{6}$$

where n is the total number of points.

There is a subtle problem with this naive definition. This does not take into account the renumbering that might happen while clustering. Cluster 1 in the true clustering might be assigned cluster 3 in the clustering produced and so on. To counter that the CE is computed for all possible renumbering of the clustering produced and the minimum of all those is taken. (This is computed efficiently by modeling the problem as a maximum weighted bipartite matching problem and then computing the solution using linear programming).

5.2 Variation of Information

Variation of Information (VI) is a *metric* introduced in [4] to compare two clusterings. It measures the amount of information that is lost/gained from going from one clustering to another and vice versa. To define it let introduce some more notation as used in [4]. (with some changes).

Let n_k be the number of point in k^{th} cluster in C. Let $P(k) = \frac{n_k}{n}$. Then the entropy associated with clustering C is defined as

$$H(\mathcal{C}) = -\sum_{k=1}^{K} P(k) \log P(k)$$

Define $n_{k_{true}k} = |\mathcal{C}_{k_{true}}^{true} \cap \mathcal{C}_{k}|$ and $P(k_{true}, k) = \frac{n_{k_{true}k}}{n}$. Then mutual information between the clustering $I(\mathcal{C}^{true}, \mathcal{C})$ is defined as

$$I(\mathcal{C}^{true}, \mathcal{C}) = \sum_{k_{true}=1}^{K} \sum_{k=1}^{K} P(k_{true}, k) \log \frac{P(k_{true}, k)}{P^{true}(k_{true})P(k)}$$

Given these the variation of information is

$$VI(\mathcal{C}^{true},\mathcal{C}) = H(\mathcal{C}) + H(\mathcal{C}^{true}) - 2I(\mathcal{C}^{true},\mathcal{C})$$

See [4] for various properties of this measure.

6 EXPERIMENTAL SETUP

In this section we provide the specific details on how the experiments were run. Throughout this section we use AffinityMatrix of a group vectors $x_1, x_2 \dots x_n$ to be the similarity matrix S such that $S_{ij} = \exp(-||x_i - x_j||^2/2\sigma^2)$ where σ would be the parameter used (We would specify the value used). Also K would refer to the input the clustering algorithms to specify the number of cluster to generate. Each algorithm was run multiple number of times and the average taken.

6.1 Exact Algorithms Used

In section 2 we described the algorithms as presented in paper the so called classic version. However to exactly distinguish between the effects of the various components of the spectral algorithms we implemented a whole range of algorithms containing most of the variations of the algorithms mentioned above.

The list of all the algorithms implemented is shown below. In each of the algorithm listed here, the various components represent the application of a particular algorithm in that "stage". So ang and mcut refer to the spectral mapping using the NJW and Multicut methods into a domain. After mapping the similarity matrix (if required) is obtained by computing the AffinityMatrix with $\sigma = 0.2$ This choice was straight forward in case of NJW algorithm as points lie on a unit sphere. We used the same value for Multicut as well.

The anchor, ward, kmeans refer to the respective algorithms applied after the spectral mappings. In kmeans we performed kmeans with 5 runs of initializing with orthogonal centers and 20 runs initialized with random centers.

We also had the intuition that the spectral methods might be more effective in clustering points after mapping them in the spectral domain. To explore that possibility we implemented the double spectral methods like ang_mcut_ward in which the first we map the points in the ang spectral domain and then those points in the mcut spectral domain, finally grouping them using the ward method.

Linkage Algorithms:	cluster_single_linkage	cluster_ward_linkage	
Multiway Spectral Algorithms:	ang_ward	ang_kmeans	ang_anchor
	mcut_ward	mcut_kmeans	mcut_anchor
	ang_ang_ward	ang_ang_kmeans	ang_mcut_ward
	ang_mcut_kmeans	mcut_mcut_ward	mcut_mcut_kmeans
Recursive Spectral Algorithms:	shi_r_ncut	kvv_mult_ncut	kvv_add_ncut
	shi_r_cond	kvv_mult_cond	kvv_add_cond

Table 1: List of Algorithms

6.2 S100

We used this database primarily to compare the robustness of algorithms w.r.t. to noise. We took the original block stochastic matrix and added uniform noise of increasing magnitude from $\eta = 10^{-0.1}$ to $\eta = 10^{0.7}$ in steps of 0.1 (in the exponent). The noise added was to preserve the *signal-to-noise* ratio in the sense that noise added in S_{ij} was made proportional to degrees of points i, j. More precisely, the new S_{ij} was calculated as follows:

$$S_{ij} = S_{ij} + \left(U(0,1) \times \eta \times \sqrt{D_i \times D_j}\right) / n$$

Where U(0,1) is number between 0 and 1 chosen at random.

For each noise levels 10 such matrices were generated and the average performance of each algorithm taken.

6.3 Handwritten digits

This dataset consisted of vector in the 64 dimensional space ranging from 0 to 16. The similarity matrix was computed as Affinity matrix with $\sigma = 10$. (We experimented with various sigma and the value of 10 seemed to give reasonable results).

For the dataset digit1000, we ran each algorithm for 5 iterations for K ranging from 8 to 12. Where as for digitFive1000, 10 iterations for K = 3 to 7 were executed.

6.4 Gene Expression Data

For both the datasets, cellcycle and cellcycle-std, the similarity was computed as the correlation coefficients between the gene expression levels of the different genes. (plus 1 to make the similarity matrix positive. So the similarities ranged from 0 to 2.) Five runs were executed for K varying from 3 to 7.

6.5 Implementation

The algorithms are very simple to implemented and we were able to implement each of them using only a few lines of code of matlab. The majority of the time taken was for the eigen decomposition. A full eigen decomposition (using eig function of matlab) would take $O(n^3)$ time. However since we just needed the top K eigen vectors, we used eigs function to reduce the time taken.

7 Performance Graphs

In this section we present the graphs for the various algorithms on the five datasets. Since there are so many algorithms we do not show them all on the same graph. For all the datasets we present six graphs. Three each for the two metrics: Clustering error (CE) and Variation of Information (VI) shown one above the another.

In first column we have the various versions of the multiway spectral algorithms. In the second column the recursive spectral algorithms and the third columns the best five. The best five are chosen as follows: First we pick the best algorithm amongst the linkage, recursive, and multiway spectral classes of algorithms. The other two are the best two of the remaining. (The "best" method was picked by looking at individual graphs) In many of the cases when there were a lot of methods with very similar performance we just chose two which looked better (or arbitrarily if that was hard to decide).

This way we can see how the various classes of spectral methods compare within themselves and w.r.t to each other. Note that **y-axis** of the graphs are *not* the same. And hence different graphs should be compared by just looking at their heights or levels. (This to done to show better contrast in between a particular class, esp. when performance within the class is near identical)

(The next Five pages contain the graphs in Landscape mode).

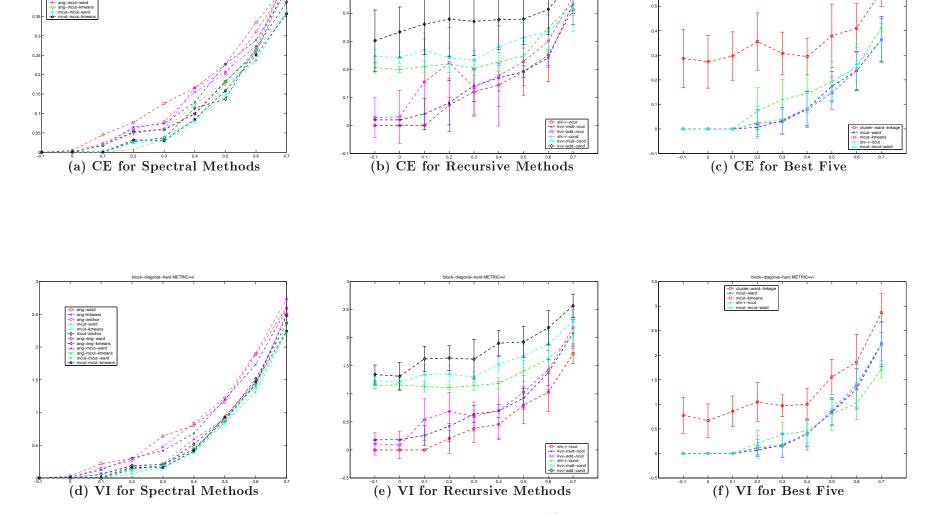
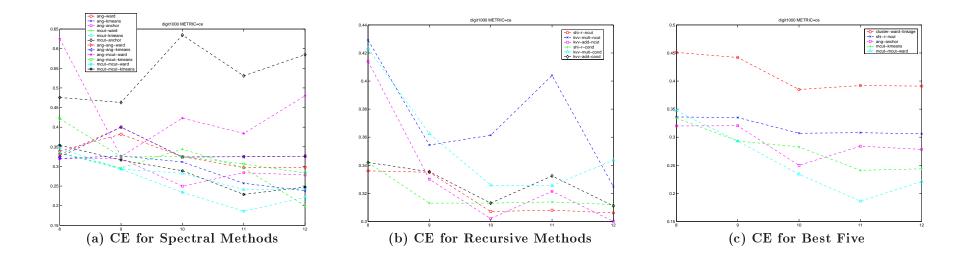


Figure 1: Block Diagonal Hard Dataset. The x-axis is the $log_{10}(\eta)$ where η is the noise added.



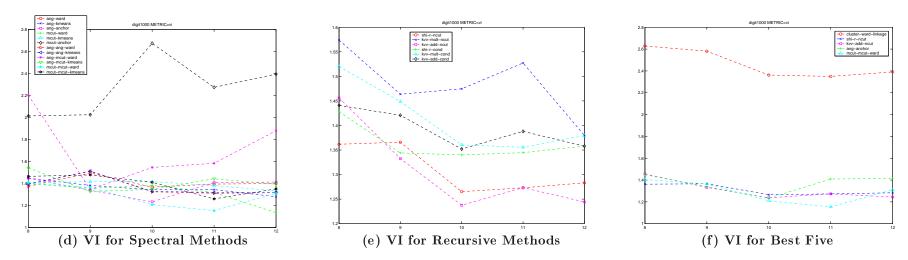
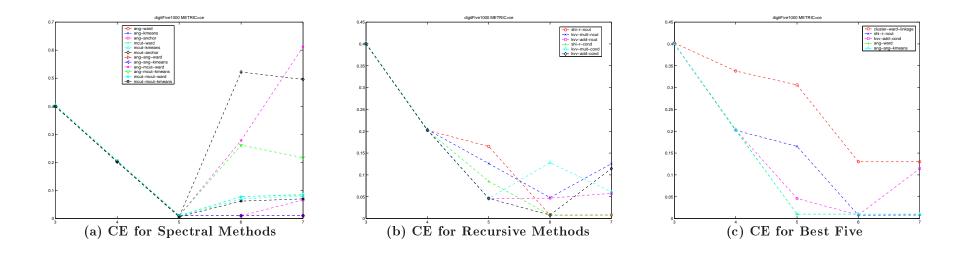


Figure 2: Handwritten digits (digit1000)



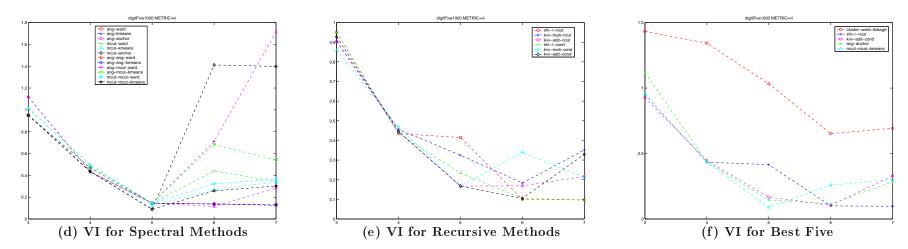
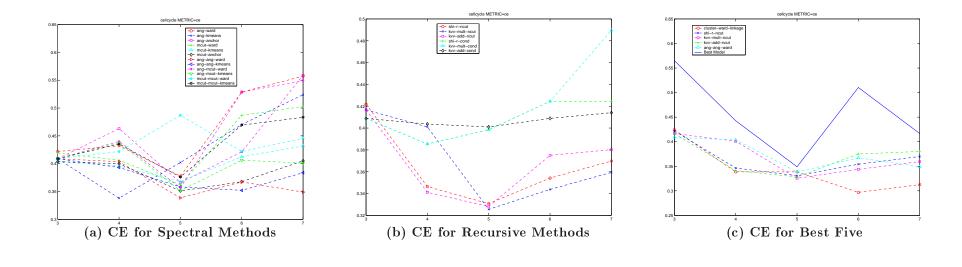


Figure 3: Handwritten digits: Five Digits (0,2,4,6,7) (digitFive1000



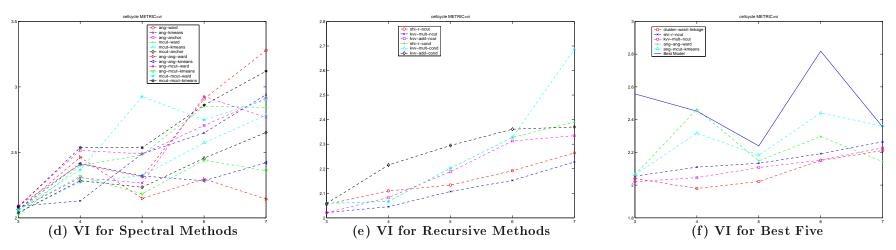
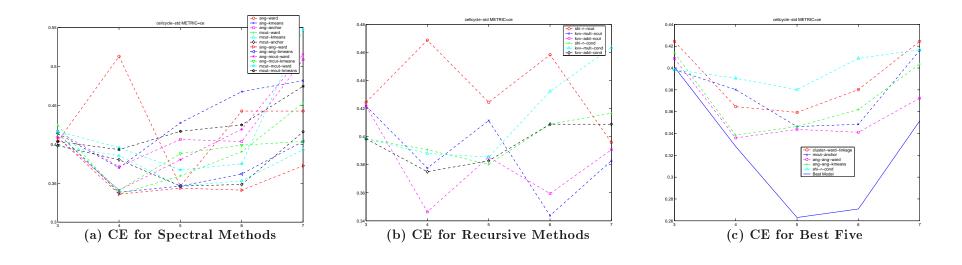


Figure 4: Log normalized yeast cell cycle data (cellcycle)



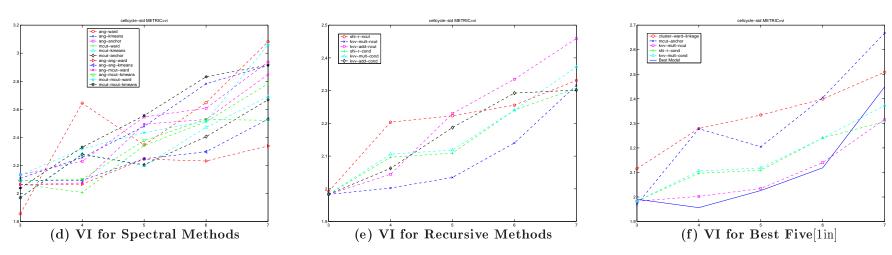


Figure 5: Standardized yeast cell cycle data cellcycle-std

8 Results and Discussion

8.1 S100

The purpose of this dataset was to demonstrate the robustness to noise. So this is the only dataset on which the error bars are shown (except for the first column in which all methods performed nearly the same with similar error bars. We omitted them to make the graphs more clear).

As we can see in figure 1 (c) and (f) linkage algorithm are too sensitive to noise and infact could not find out the correct clustering even when (almost) no noise was added to the block stochastic matrix. The multiway spectral methods as expected perform the best as this is their perfect S. Within this class All the algorithms seems to perform nearly same with mcut methods performing slightly better than the ang methods. This suggests that for the block stochastic similarity matrices it might be slightly preferable to use the Multicut base algorithms. The reason for this might be the that NJW maps the clusters to the unit sphere and this might blow up distances between points that are in the same cluster. However the experimental proof is not conclusive.

In the recursive algorithms only the shi-r-ncut gets the perfect clusters in case of low noise though other variants based on also perform well. It is interesting to note that conductance based performs significantly worse. This is again expected as the the conductance only takes the smaller cluster size into account while *Ncut* is based on both the cluster sizes.

8.2 Handwritten Digits

This is the first real dataset the we tested the algorithm on. On the complete dataset digit1000 (figure 2) the multiway spectral methods perform slightly better than the recursive spectral. However the performance difference is not that significant (and in case of the VI measure almost zero) to make any conclusive statement. The linkage algorithms performed a lot worse. Within the multiway spectral algorithms, mcut_mcut_ward seems to be the clear winner.

The results on the digitFive1000 dataset are much more interesting. The performance is near perfect (at K=5) and hence the comparison could be done in light of a dataset with well established structure. If we take a look at figure 3 (a) then is is easy to see that all the multiway spectral algorithm give nearly identical results from K=3 to 5. This is a strong empirical justification of the similarity of the NJW and Multicut which was theoretically proved above. Also in this particular the clusters are obviously well formed as the result in this sections are independent of the grouping algorithm that is used in the third stage.

This is also one dataset in which the multiway spectral methods seem to dominate over the recursive methods. The linkage methods are as expected far behind. One surprising thing observed is that the *Ncut* methods are lagging behind in the performance as compared to those based on *conductance*.

8.3 Gene Expression Data

This dataset was more interesting of the two real datasets we used. There are a variety of reasons. First of all, since we had results from the model based algorithms for this dataset (from [12]) there was something to compare the spectral algorithms with rather than just amongst themselves. Secondly this contained the same dataset with different data transformation applied to them (See section 4.3), we could see how much the clustering algorithms are dependent on preprocessing.

For the cellcycle dataset the best of the spectral algorithms perform slightly better than the model based algorithms. This is encouraging as this shows that spectral methods are competitive even on real dataset and not just the perfect case. The recursive algorithm show similar performance as the multiway algorithms except that Ncut based algorithms are a little better and the conductance based a little worse. While the ordering within the recursive algorithms is expected it is not clear why some of them are better than the multiway algorithms. It is possible that in presence of noise depending on the later eigenvectors is not always the best thing to do and it is better to do the process recursively which ensures that atleast the first few partitions are correct.

However what is even more surprising is the performance of cluster-ward-linkage. This simple linkage algorithm gives nearly the best performance on both measures!. We think that in case of such high error

rates as we are observing here it is really anybody's game unless there is a dominant structure known to be in the data which corresponds to the clustering algorithm.

In comparison the situation for cellcycle_std is completely reversed! The model based algorithm performs the best. The reason for this is that this data transformation is known to fit better to gaussian model and hence the better performance. The performance of best spectral methods remains the same, though the multiway methods perform better now than recursive ones. (with conductance based methods now just slightly worse and even better at K = 5.).

8.4 Future work

In this paper we did not address the problem of how to go about choosing the number of clusters. We intend to explore methods which could find the number of clusters based on the data.

There also are two other algorithms that we did not implement for the lack of time. The first one is another variant of the SM algorithm which theoretically should perform very well on block stochastic matrix. We did not use it because we think it might be too sensitive to noise. The second is a non-spectral method based on single linkage and runt analysis which we expect to be a lot more robust to noise. We wish to explore how using this algorithm as the grouping algorithm after spectral mapping affects the performance of various methods.

9 Conclusion

The goal of the present paper was to analyze comparatively the features of a number of published spectral clustering algorithms. Rather than establishing which of the published algorithms is better, we aimed at evaluating what features make a spectral clustering more valuable.

Because for clustering a data set the "goodness" is in the eye of the beholder, one should look at clustering algorithms not only as competing with each other but also as *complementing* each other's strengths and weaknesses. Hence, a second goal of our research, was to see how different the various algorithms are in their approach.

The answer to the second question is largely negative. The theory predicts that the perfect S for all three algorithms is the same, a result that is strongly supported by the experiments. All algorithms work very well in the cases when S is almost perfect and there is not clear winner in case it is not. We did find the multiway spectral clustering algorithm to be slightly better performing especially when there is structure to be easily found in the data. For the recursive methods we recommend using the Ncut measure over others though other than that there is no clear winner. As compared to other method we showed that spectral methods give competitive performance to the existing methods and are definitely worth further exploration.

Acknowledgements

. I am thankful to Marina Meila, as my collaborator, in the project without whose immense contribution this work would not have been possible, and as an advisor who has been very supportive and thoughtful. Special thanks to Thomas Richardson for providing support for Deepak Verma from NSF grant DMS-9972008. We would also like to thank Ka Yee Yeung for providing the yeast cell cycle dataset and the classifications. And finally to Jayant Madhavan who provided feedback on earlier version of this paper.

References

- [1] Sanjoy Dasgupta. Performance guarantees for hierarchical clustering. In *Proceedings of the 15th Annual Conference on Computational Learning Theory (COLT 2002)*, pages 351–363. Springer, 2002.
- [2] Ravi Kannan, Santosh Vempala, and Adrian Vetta. On clusterings good, bad and spectral. In *FOCS*, pages 367–377, 2000.

- [3] C. Kaynak. Methods of combining multiple classifiers and their applications to handwritten digit recognition. Master's thesis, Institute of Graduate Studies in Science and Engineering, Bogazici University., 1995.
- [4] Marina Meila. Comparing clusterings. Technical Report 418, UW Statistics Department, 2002.
- [5] Marina Meila and Jianbo Shi. Learning segmentation by random walks. In NIPS, pages 873–879, 2000.
- [6] Marina Meila and Jianbo Shi. A random walks view of spectral segmentation, 2001.
- [7] Andrew Moore. The anchors hierarchy: Using the triangle inequality to survive high-dimensional data. In *Proceedings of the Twelfth Conference on Uncertainty in Artificial Intelligence*, pages 397–405. AAAI Press, 2000.
- [8] A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, Advances in Neural Information Processing Systems 14, pages 849–856, Cambridge, MA, 2002. MIT Press.
- [9] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- [10] Daniel A. Spielman and Shang-Hua Teng. Spectral partitioning works: Planar graphs and finite element meshes. In *IEEE Symposium on Foundations of Computer Science*, pages 96–105, 1996.
- [11] J.H. Ward. Hierarchical grouping to optimize an objective function. J. Amer. Statist. Assoc., pages 236 244, 1963.
- [12] K. Yeung, C. Fraley, A. Murua, A. Raftery, and W. Ruzzo. Model-based clustering and data transformations for gene expression data. Technical Report UW-CSE-01-04-02, Dept. of Computer Science and Engineering, University of Washington, 2001.
- [13] Ka Yee Yeung. Model-based clustering and data transformations for gene expression data. http://staff.washington.edu/kayee/model/.