

# Learning a Bi-Stochastic Data Similarity Matrix

Fei Wang

*Department of Statistical Science  
Cornell University  
Ithaca, NY 14853, USA  
fw83@cornell.edu*

Ping Li

*Department of Statistical Science  
Cornell University  
Ithaca, NY 14853, USA  
pingli@cornell.edu*

Arnd Christian König

*Microsoft Research  
Microsoft Cooperation  
Redmond, WA 98052, USA  
chrisko@microsoft.com*

**Abstract**—An idealized clustering algorithm seeks to learn a cluster-adjacency matrix such that, if two data points belong to the same cluster, the corresponding entry would be 1; otherwise the entry would be 0. This integer (1/0) constraint makes it difficult to find the optimal solution. We propose a relaxation on the cluster-adjacency matrix, by deriving a bi-stochastic matrix from a data similarity (e.g., kernel) matrix according to the Bregman divergence. Our general method is named the *Bregmanian Bi-Stochastication* (BBS) algorithm.

We focus on two popular choices of the Bregman divergence: the Euclidian distance and the KL divergence. Interestingly, the BBS algorithm using the KL divergence is equivalent to the Sinkhorn-Knopp (SK) algorithm for deriving bi-stochastic matrices. We show that the BBS algorithm using the Euclidian distance is closely related to the relaxed  $k$ -means clustering and can often produce noticeably superior clustering results than the SK algorithm (and other algorithms such as Normalized Cut), through extensive experiments on public data sets.

## I. INTRODUCTION

Clustering [13], [6], which aims to organize data in an unsupervised fashion, is one of the fundamental problems in data mining and machine learning. The basic goal is to group the data points into clusters such that the data in the same cluster are “similar” to each other while the data in different clusters are “different” from each other.

In this paper, we view clustering from the perspective of matrix approximation. Suppose we are given a data set  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^n$ , which comes from  $k$  clusters. We can denote the cluster memberships by an  $n \times k$  matrix  $\mathbf{F}$ , such that

$$F_{ij} = \begin{cases} 1, & \text{if } \mathbf{x}_i \in \pi_j \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

where  $\pi_j$  denotes the  $j$ -th cluster. It is often more convenient to proceed with the scaled version  $\tilde{\mathbf{F}}$  [17][24], such that

$$\tilde{F}_{ij} = \begin{cases} 1/\sqrt{n_j}, & \text{if } \mathbf{x}_i \in \pi_j \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

where  $n_j = |\pi_j|$  is the cardinality of cluster  $\pi_j$ . Note that  $\tilde{\mathbf{F}}$  has (at least) the following properties (constraints):

$$\tilde{\mathbf{F}} \geq 0 \text{ (i.e., } \tilde{F}_{ij} \geq 0 \forall i, j\text{), } \tilde{\mathbf{F}}^\top \tilde{\mathbf{F}} = \mathbf{I}, \quad (\tilde{\mathbf{F}} \tilde{\mathbf{F}}^\top) \mathbf{1} = \mathbf{1},$$

where  $\mathbf{1} \in \mathbb{R}^{n \times 1}$  is an all-one vector, and  $\mathbf{I} \in \mathbb{R}^{n \times n}$  is an identity matrix.

If we define  $\mathbf{G} = \tilde{\mathbf{F}} \tilde{\mathbf{F}}^\top$ , we can hope to discover the cluster structure of  $\mathcal{X}$  from  $\mathbf{G}$ . The constraints on  $\tilde{\mathbf{F}}$  can be transferred to the constraints on  $\mathbf{G}$  as

$$\mathbf{G} \geq 0, \quad \mathbf{G} = \mathbf{G}^\top, \quad \mathbf{G}\mathbf{1} = \mathbf{1} \quad (3)$$

In other words,  $\mathbf{G}$  is a symmetric, nonnegative, and *bi-stochastic* (also called *doubly stochastic*) matrix [12].

### A. Deriving a Bi-Stochastic Matrix from a Similarity Matrix

The bi-stochastic matrix  $\mathbf{G}$ , constructed from the cluster-membership matrix ( $\mathbf{F}$  or  $\tilde{\mathbf{F}}$ ) can be viewed as a special type of *similarity matrix*. Naturally, one might conjecture: *If we relax the integer (0/1) constraint on  $\mathbf{F}$ , can we still derive a (useful) bi-stochastic matrix from a similarity matrix?*

For example, a popular family of data similarity matrix is the *Gaussian kernel* matrix,  $\mathbf{K} \in \mathbb{R}^{n \times n}$ , where each entry

$$K_{ij} = \exp\left(-\frac{1}{\gamma} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right), \quad \gamma > 0 \quad (4)$$

Here,  $\gamma$  is a tuning parameter. Obviously, an arbitrary similarity matrix can not be guaranteed to be bi-stochastic. For a given similarity matrix, there are multiple ways to derive a bi-stochastic matrix. We first review a straightforward solution known as the *Sinkhorn-Knopp (SK)* algorithm.

### B. The Sinkhorn-Knopp (SK) Algorithm

The following *Sinkhorn-Knopp Theorem* [18] says that, under mild regularity conditions, one can construct a bi-stochastic matrix from a similarity matrix.

**Theorem (Sinkhorn-Knopp)** *Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a nonnegative square matrix. A necessary and sufficient condition that there exists a bi-stochastic matrix  $\mathbf{P}$  of the form:  $\mathbf{P} = \mathbf{U}\mathbf{A}\mathbf{V}$ , where  $\mathbf{U}$  and  $\mathbf{V}$  are diagonal matrices with positive main diagonals, is that  $\mathbf{A}$  has total support. If  $\mathbf{P}$  exists, then it is unique.  $\mathbf{U}$  and  $\mathbf{V}$  are also unique up to a scalar multiple if and only if  $\mathbf{A}$  is fully indecomposable.*

Based on this theorem, [18] proposed a method called the *Sinkhorn-Knopp (SK)* algorithm to obtain a bi-stochastic matrix from a nonnegative matrix  $\mathbf{A}$ , by generating a sequence of matrices whose columns and rows are normalized alternatively. The limiting matrix is bi-stochastic. In particular,

if  $\mathbf{A}$  is symmetric, then the resulting matrix  $\mathbf{P} = \mathbf{UAV}$  is also symmetric with  $\mathbf{U}$  and  $\mathbf{V}$  being equal (up to a constant multiplier). The following example illustrates the procedure:

$$\begin{aligned}\mathbf{A} &= \begin{bmatrix} 1 & 0.8 & 0.6 \\ 0.8 & 1 & 0.4 \\ 0.6 & 0.4 & 1 \end{bmatrix} \\ \longrightarrow &\begin{bmatrix} 0.4167 & 0.3636 & 0.3000 \\ 0.3333 & 0.4545 & 0.2000 \\ 0.2500 & 0.1818 & 0.5000 \end{bmatrix} \longrightarrow \begin{bmatrix} 0.3857 & 0.3366 & 0.2777 \\ 0.3374 & 0.4601 & 0.2025 \\ 0.2683 & 0.1951 & 0.5366 \end{bmatrix} \\ \longrightarrow \dots & \\ \longrightarrow &\begin{bmatrix} 0.3886 & 0.3392 & 0.2722 \\ 0.3392 & 0.4627 & 0.1980 \\ 0.2722 & 0.1980 & 0.5297 \end{bmatrix} = \mathbf{P}\end{aligned}$$

The SK algorithm is not the unique construction. In statistics, this procedure is also known as the *iterative proportional scaling* algorithm [5], [20].

### C. Connection to the Normalized Cut (Ncut) Algorithm

Interestingly, the well-known *Normalized Cut (Ncut)* algorithm [17] can be viewed as a one-step construction towards producing bi-stochastic matrices. The Ncut algorithm normalizes a similarity matrix  $\mathbf{K} \in \mathbb{R}^{n \times n}$  with  $\mathbf{D} = \text{diag}(\mathbf{K}\mathbf{1})$ , where  $\mathbf{1} \in \mathbb{R}^{n \times 1}$  is an all-one vector, to be

$$\tilde{\mathbf{K}} = \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2} \quad (5)$$

[23] showed that if one keeps normalizing  $\mathbf{K}$  with

$$\mathbf{K}^{(t+1)} = (\mathbf{D}^{(t)})^{-1/2} \mathbf{K}^{(t)} (\mathbf{D}^{(t)})^{-1/2}, \quad (6)$$

then  $\mathbf{K}^{(\infty)}$  will be bi-stochastic.

### D. Our Proposed General Framework: BBS

In this paper, we propose to obtain a bi-stochastic matrix  $\mathbf{G} \in \mathbb{R}^{n \times n}$  from some initial similarity matrix  $\mathbf{K} \in \mathbb{R}^{n \times n}$ , by solving the following optimization problem

$$\begin{aligned}\min_{\mathbf{G}} \quad & D_\phi(\mathbf{G}, \mathbf{K}) = \sum_{ij} D_\phi(G_{ij}, K_{ij}) \\ \text{s.t.} \quad & \mathbf{G} \geq 0, \quad \mathbf{G} = \mathbf{G}^\top, \quad \mathbf{G}\mathbf{1} = \mathbf{1}\end{aligned} \quad (7)$$

where

$$D_\phi(x, y) \triangleq \phi(x) - \phi(y) - \nabla\phi(y)(x - y), \quad (8)$$

is the *Bregman divergence* between  $x$  and  $y$  with  $\phi$  being a strictly convex function. The problem (7) is a standard convex optimization program. We name the solution  $\mathbf{G}$  the *Bregmanian Bi-Stochasticization (BBS)* of  $\mathbf{K}$ .<sup>1</sup>

Two choices of the Bregman divergence  $D_\phi$  are popular:

1)  $\phi(x) = x^2/2$ : (squared) Euclidian distance,

2)  $\phi(x) = x \log x - x$ : Kullback-Leibler (KL) divergence.

It can be shown that the SK algorithm is equivalent to BBS using KL divergence. We will demonstrate that BBS with  $\phi(x) = x^2/2$  often produces superior clustering results over the SK algorithm (and other algorithms such as Ncut).

<sup>1</sup>Also see the work on matrix nearness in Bregman divergence without the bi-stochastic constraint [7].

## II. BREGMANIAN BI-STOCHASTICATION (BBS)

The BBS algorithm seeks a bi-stochastic matrix  $\mathbf{G}$  which optimally approximates  $\mathbf{K}$  in the Bregman divergence sense, by solving the optimization problem (7). For the two popular choices of the Bregman divergence  $D_\phi$  in Eq. (8), we study specially designed optimization strategies, for better insights.

### A. $\phi(x) = x^2/2$

For this choice of  $\phi(x)$ , we have

$$\begin{aligned}D_\phi(\mathbf{G}, \mathbf{K}) &= \sum_{ij} D_\phi(G_{ij}, K_{ij}) \\ &= \sum_{ij} \frac{1}{2} G_{ij}^2 - \frac{1}{2} K_{ij}^2 - K_{ij}(G_{ij} - K_{ij}) \\ &= \frac{1}{2} \|\mathbf{G} - \mathbf{K}\|_F^2 = \frac{1}{2} \text{tr} \left( (\mathbf{G} - \mathbf{K})^\top (\mathbf{G} - \mathbf{K}) \right) \\ &= \frac{1}{2} \text{tr} \left( \mathbf{K}^\top \mathbf{K} + \mathbf{G}^\top \mathbf{G} - 2\mathbf{K}^\top \mathbf{G} \right)\end{aligned} \quad (9)$$

Thus, the BBS problem with  $\phi(x) = x^2/2$  is equivalent to

$$\begin{aligned}\min_{\mathbf{G}} \quad & \text{tr} \left( \mathbf{G}^\top \mathbf{G} - 2\mathbf{K}^\top \mathbf{G} \right) \\ \text{s.t.} \quad & \mathbf{G} \geq 0, \quad \mathbf{G} = \mathbf{G}^\top, \quad \mathbf{G}\mathbf{1} = \mathbf{1}\end{aligned} \quad (10)$$

Problem (10) is a *Quadratic Programming* program [2], [16] and can be solved by standard methods such as the *interior point* algorithm. Here, we adopt a simple cyclic constraint projection approach, known as the *Dykstra algorithm* [10]. First we split the constraints into two sets  $\mathcal{C}_1$  and  $\mathcal{C}_2$ :

$$\mathcal{C}_1 : \quad \{\mathbf{G} | \mathbf{G} = \mathbf{G}^\top, \quad \mathbf{G}\mathbf{1} = \mathbf{1}\} \quad (11)$$

$$\mathcal{C}_2 : \quad \{\mathbf{G} | \mathbf{G} \geq 0\} \quad (12)$$

where  $\mathcal{C}_1$  defines an affine set and  $\mathcal{C}_2$  defines a convex set.

For the constraint set  $\mathcal{C}_1$ , we need to solve<sup>2</sup>

$$\begin{aligned}\min_{\mathbf{G}} \quad & \text{tr} \left( \mathbf{G}^\top \mathbf{G} - 2\mathbf{K}^\top \mathbf{G} \right) \\ \text{s.t.} \quad & \mathbf{G} = \mathbf{G}^\top, \quad \mathbf{G}\mathbf{1} = \mathbf{1},\end{aligned} \quad (13)$$

for which we first introduce a Lagrangian function

$$\begin{aligned}\mathcal{L}(\mathbf{G}) &= \text{tr} \left( \mathbf{G}^\top \mathbf{G} - 2\mathbf{K}^\top \mathbf{G} \right) - \boldsymbol{\mu}_1^\top (\mathbf{G}\mathbf{1} - \mathbf{1}) \\ &\quad - \boldsymbol{\mu}_2^\top (\mathbf{G}^\top \mathbf{1} - \mathbf{1})\end{aligned} \quad (14)$$

where  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^{n \times 1}$  are Lagrangian multipliers. By the constraint  $\mathbf{G} = \mathbf{G}^\top$ , we know  $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \boldsymbol{\mu}$ . Thus

$$\nabla_{\mathbf{G}} \mathcal{L}(\mathbf{G}) = 2(\mathbf{G} - \mathbf{K}) - \boldsymbol{\mu} \mathbf{1}^\top - \mathbf{1} \boldsymbol{\mu}^\top \quad (15)$$

Setting  $\nabla_{\mathbf{G}} \mathcal{L}(\mathbf{G}) = 0$  yields

$$\mathbf{G} = \mathbf{K} + \frac{1}{2} \boldsymbol{\mu} \mathbf{1}^\top + \frac{1}{2} \mathbf{1} \boldsymbol{\mu}^\top \quad (16)$$

Since  $\mathbf{G}$  must satisfy the constraint  $\mathbf{G}\mathbf{1} = \mathbf{1}$ , we can right-multiply  $\mathbf{1}$  on the both sides of Eq. (16) as

$$\mathbf{1} = \mathbf{G}\mathbf{1} = \mathbf{K}\mathbf{1} + \frac{n}{2} \boldsymbol{\mu} + \frac{1}{2} \mathbf{1} \mathbf{1}^\top \boldsymbol{\mu} \quad (17)$$

<sup>2</sup>It may be also formulated as an instance of the *Least Norm* problem [2].

Then we obtain

$$\boldsymbol{\mu} = 2(n\mathbf{I} + \mathbf{1}\mathbf{1}^\top)^{-1}(\mathbf{I} - \mathbf{K})\mathbf{1} \quad (18)$$

By making use of the *Woodbury formula* [11], we obtain

$$(n\mathbf{I} + \mathbf{1}\mathbf{1}^\top)^{-1} = \frac{1}{n} \left( \mathbf{I} - \frac{1}{2n}\mathbf{1}\mathbf{1}^\top \right) \quad (19)$$

We can then write the solution in a closed form:

$$\mathbf{G} = \mathbf{K} + \left( \frac{1}{n}\mathbf{I} - \frac{1}{n}\mathbf{K} + \frac{\mathbf{1}\mathbf{1}^\top\mathbf{K}}{n^2} \right) \mathbf{1}\mathbf{1}^\top - \frac{1}{n}\mathbf{1}\mathbf{1}^\top\mathbf{K} \quad (20)$$

For the constraint set  $\mathcal{C}_2$ , we need to solve another optimization problem:

$$\begin{aligned} \min_{\mathbf{G}} \quad & \frac{1}{2}\|\mathbf{G} - \mathbf{K}\|_F^2 \\ \text{s.t.} \quad & \mathbf{G} \geq 0 \end{aligned} \quad (21)$$

whose solution is simply

$$\mathbf{G} = \mathbf{K}^+ \quad (22)$$

where  $\mathbf{K}^+$  denotes the positive part of  $\mathbf{K}$ .

The overall algorithm of BBS with  $\phi(x) = x^2/2$  is summarized in Alg. 1. The total computational complexity of Alg. 1 is  $O(Tn^2)$  with  $T$  being the number of iterations needed for the algorithm to converge.

---

#### Algorithm 1 BBS WITH $\phi(x) = x^2/2$

---

**Require:** An initial similarity matrix  $\mathbf{K}$

- 1:  $t = 0$ ,  $\mathbf{G}^{(t)} = \mathbf{K}$ .
  - 2: **repeat**
  - 3:    $t \leftarrow t + 1$
  - 4:    $\mathbf{G}^{(t+1)} \leftarrow [\mathbf{G}^{(t)} + \frac{1}{n}(\mathbf{I} - \mathbf{G}^{(t)} + \frac{\mathbf{1}\mathbf{1}^\top\mathbf{G}^{(t)}}{n})] \mathbf{1}\mathbf{1}^\top - \frac{1}{n}\mathbf{1}\mathbf{1}^\top\mathbf{G}^{(t)}]^+$
  - 5: **until** Some convergence condition is satisfied
  - 6: Output  $\mathbf{G}^{(t)}$
- 

#### B. $\phi(x) = x \log x - x$

$$\begin{aligned} D_\phi(\mathbf{G}, \mathbf{K}) &= \sum_{ij} D_\phi(G_{ij}, K_{ij}) \quad (23) \\ &= \sum_{ij} G_{ij} \log \frac{G_{ij}}{K_{ij}} + K_{ij} - G_{ij} = KL(\mathbf{G} \parallel \mathbf{K}) \end{aligned}$$

The BBS problem becomes

$$\begin{aligned} \min_{\mathbf{G}} \quad & KL(\mathbf{G} \parallel \mathbf{K}) \\ \text{s.t.} \quad & \mathbf{G} \geq 0, \mathbf{G} = \mathbf{G}^\top, \mathbf{G}\mathbf{1} = \mathbf{1} \end{aligned} \quad (24)$$

We construct the following Lagrangian

$$\mathcal{L}(\mathbf{G}) = KL(\mathbf{G} \parallel \mathbf{K}) - \boldsymbol{\mu}_1^\top(\mathbf{G}^\top\mathbf{1} - \mathbf{1}) - \boldsymbol{\mu}_2^\top(\mathbf{G}\mathbf{1} - \mathbf{1}), \quad (25)$$

where we drop the constraint  $\mathbf{G} \geq 0$  for the time being, and we will later show it is automatically satisfied. Therefore

$$\nabla_{\mathbf{G}}\mathcal{L}(\mathbf{G}) = \log \mathbf{G} - \log \mathbf{K} - \boldsymbol{\mu}_1\mathbf{1}^\top - \mathbf{1}\boldsymbol{\mu}_2^\top \quad (26)$$

where  $\log$  represents the elementwise logarithm. Setting  $\nabla_{\mathbf{G}}\mathcal{L}(\mathbf{G}) = 0$  yields

$$\log G_{ij} - \log K_{ij} - \mu_{1i} - \mu_{2j} = 0 \quad (27)$$

Thus the solution satisfies

$$G_{ij} = e^{\mu_{1i}} K_{ij} e^{\mu_{2j}} \quad (28)$$

Next, we define the following two vectors

$$\boldsymbol{\pi}_1 = [e^{\mu_{11}}, e^{\mu_{12}}, \dots, e^{\mu_{1n}}]^\top \in \mathbb{R}^{n \times 1} \quad (29)$$

$$\boldsymbol{\pi}_2 = [e^{\mu_{21}}, e^{\mu_{22}}, \dots, e^{\mu_{2n}}]^\top \in \mathbb{R}^{n \times 1} \quad (30)$$

and two diagonal matrices  $\text{diag}(\boldsymbol{\pi}_1) \in \mathbb{R}^{n \times n}$ ,  $\text{diag}(\boldsymbol{\pi}_2) \in \mathbb{R}^{n \times n}$ . This way, we can express the solution to be

$$\mathbf{G} = \text{diag}(\boldsymbol{\pi}_1) \times \mathbf{K} \times \text{diag}(\boldsymbol{\pi}_2) \quad (31)$$

As  $\mathbf{G}$  is symmetric, we know  $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \boldsymbol{\mu}$  and  $\boldsymbol{\pi}_1 = \boldsymbol{\pi}_2 = \boldsymbol{\pi}$ . By comparing with the *Sinkhorn-Knopp Theorem*, we can immediately see that the BBS algorithm with  $\phi(x) = x \log x - x$  actually recovers the symmetric SK algorithm, and  $\text{diag}(\boldsymbol{\pi})$  is used for scaling  $\mathbf{K}$  to be bi-stochastic.

We should mention that, it appears that the fact that the symmetric SK algorithm minimizes the KL divergence was essentially discovered in statistics [4], [19].

### III. EXPERIMENTS

#### A. Data Sets

Table I summarizes the data sets used in our experiments:

- **MNIST**<sup>3</sup>: We randomly sampled 6000 data points from the original training set. We also created a smaller data set, **MNIST (0-4)**, using digits 0, 1, 2, 3, 4.
- **ISOLET**<sup>4</sup>: We took the original UCI training set and divided it into three smaller data sets so that the number of classes (clusters) for each set is not too large.
- **LETTER**<sup>5</sup>: We divided the original data into five sets.
- **NEWS20**<sup>6</sup>: The test set from the LibSVM site.
- **OPTDIGIT**<sup>7</sup>: We combined the original (UCI) training and test sets, as this data set is not too large.
- **PENDIGIT**<sup>8</sup>: The original (UCI) training set.
- **SATIMAGE**<sup>9</sup>: The original (UCI) training set.
- **SHUTTLE**<sup>10</sup>: The test set from the LibSVM site.
- **VEHICLE**<sup>11</sup>: The version from the LibSVM site.
- **ZIPCODE**<sup>12</sup>: We used the training set and also constructed a smaller data set using digits 0, 1, 2, 3, 4.

<sup>3</sup><http://yann.lecun.com/exdb/mnist/>

<sup>4</sup><http://archive.ics.uci.edu/ml/machine-learning-databases/isolet/isolet1+2+3+4.data.Z>

<sup>5</sup><http://archive.ics.uci.edu/ml/machine-learning-databases/letter-recognition/letter-recognition.data>

<sup>6</sup><http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multiclass/news20.t.scale.bz2>

<sup>7</sup><http://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits>

<sup>8</sup><http://archive.ics.uci.edu/ml/machine-learning-databases/pendigits/pendigits.tra>

<sup>9</sup><http://archive.ics.uci.edu/ml/machine-learning-databases/statlog/satimage/sat.trn>

<sup>10</sup><http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multiclass/shuttle.scale.t>

<sup>11</sup><http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multiclass/vehicle.scale>

<sup>12</sup><http://www-stat.stanford.edu/~tibs/ElemStatLearn/datasets/zip.train.gz>

Table I  
DATA SETS

Data Set	# Samples ( $n$ )	# Dimensions ( $d$ )	# Classes ( $k$ )
MNIST	6000	784	10
MNIST (0-4)	3031	784	5
ISOLET (A-I)	2158	617	9
ISOLET (J-R)	2160	617	9
ISOLET (S-Z)	1920	617	8
LETTER (A-E)	3864	16	5
LETTER (F-J)	3784	16	5
LETTER (K-O)	3828	16	5
LETTER (P-T)	3888	16	5
LETTER (U-Z)	4636	16	6
NEWS20	3993	62060	20
OPTDIGIT	5620	64	10
PENDIGIT	7494	16	10
SATIMAGE	4465	36	6
SHUTTLE	14500	9	7
VEHICLE	846	18	4
ZIPCODE	7291	256	10
ZIPCODE (0-4)	4240	256	5

## B. Experiment Procedure

For all data sets, we always normalized each data point (vector) to have a unit  $l_2$  norm, and we always used the Gaussian kernel Eq. (4) to form the initial similarity matrix  $\mathbf{K}$ . For the tuning parameter  $\gamma$  in Eq. (4), we experimented with  $\gamma \in \{1024, 256, 64, 32, 16, 8, 4, 2, 1, 0.5, 0.25\}$ .

We ran the BBS algorithm with  $\phi(x) = x^2/2$  for 1000 iterations at each  $\gamma$ . We also ran the SK algorithm (i.e., BBS with  $\phi(x) = x \log x - x$ ) for 1000 iterations at each  $\gamma$ .

We eventually used spectral clustering [17], [3], [8], [15] to evaluate the quality of the produced bi-stochastic matrices. In particular, we used the procedure described in [15]. That is, we computed the top- $k$  eigenvectors of the bi-stochastic matrix to form a new  $n \times k$  matrix and normalized each row to have a unit  $l_2$  norm. Denote the resulting new “data matrix” by  $\mathbf{Z}$ . We then used Matlab *kmeans* function:

```
kmeans(Z, k, 'MaxIter', 1000, 'EmptyAction', 'singleton')
```

We ran *kmeans* 100 times and reported both the *average* and *maximum* clustering results.

However, we would like to first introduce two measures that may allow us to directly assess the quality of the bi-stochastic matrices independent of the clustering algorithms.

### C. Quality Measurements of the Bi-Stochastic Matrices

After we have generated a (hopefully) bi-stochastic matrix  $\mathbf{P} \in \mathbb{R}^{n \times n}$  from a similarity matrix  $\mathbf{K}$ , we can compute:

$$M_B = \frac{1}{n} \sum_{i=1}^n \left| \sum_{j=1}^n P_{ij} - 1 \right|, \quad (32)$$

$$M_C = \sum_{c=1}^k \frac{1}{n} \sum_{i=1}^n \left| \sum_{j=1, x_j \in \pi_c}^n P_{ij} - 1 \right| \quad (33)$$

Basically,  $M_B$  measures how far  $\mathbf{P}$  is from being a bi-stochastic matrix, and  $M_C$  roughly measures the potential of producing good clustering results. Lower values of  $M_B$  and  $M_C$  are more desirable. We use the  $M_C$  measure because it is independent of the specific clustering algorithms.

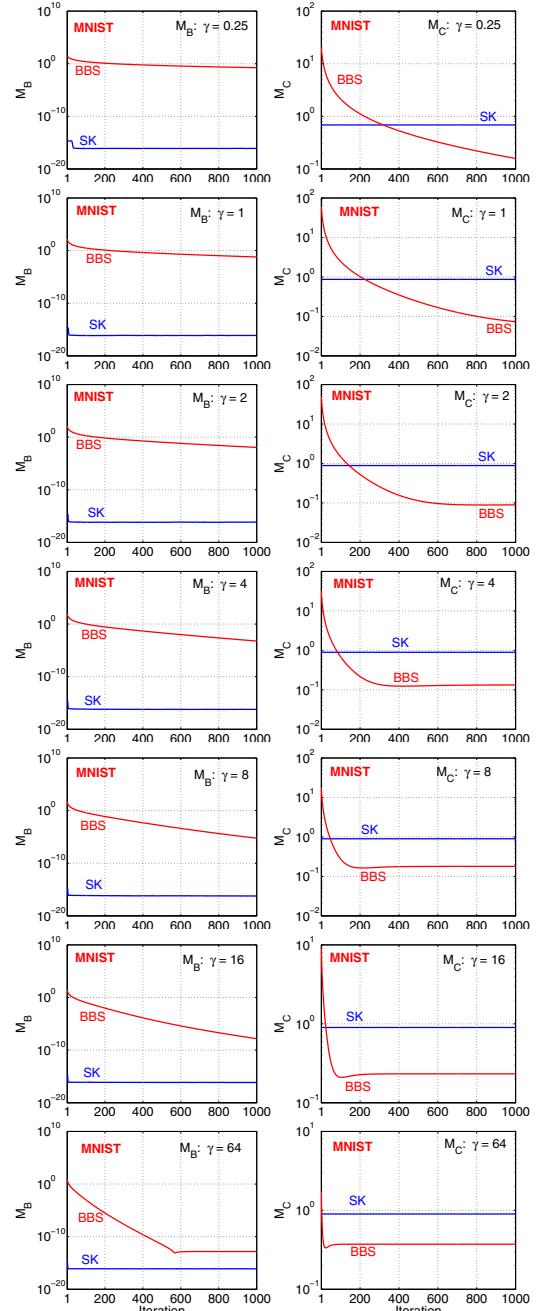


Figure 1. Quality measurements  $M_B$  (32) and  $M_C$  (33) (lower is better), on MNIST data, for up to 1000 iterations.  $\gamma$  is the kernel tuning parameter in Eq. (4). “BBS” labels the curves produced by the BBS algorithm with  $\phi(x) = x^2/2$  (i.e., Alg. 1) and “SK” the curves by the SK algorithm.

Fig. 1 presents the quality measurements on the MNIST data set, for a wide range of  $\gamma$  values. Fig. 2 presents the measurements on a variety of data sets for  $\gamma = 1$ :

- In terms of  $M_B$ , the SK algorithm performs well in producing good bi-stochastic matrices.
- In terms of  $M_C$ , the BBS algorithm using the Euclidian distance (Alg. 1) has noticeably better potential of producing good clustering results than the SK algorithm.

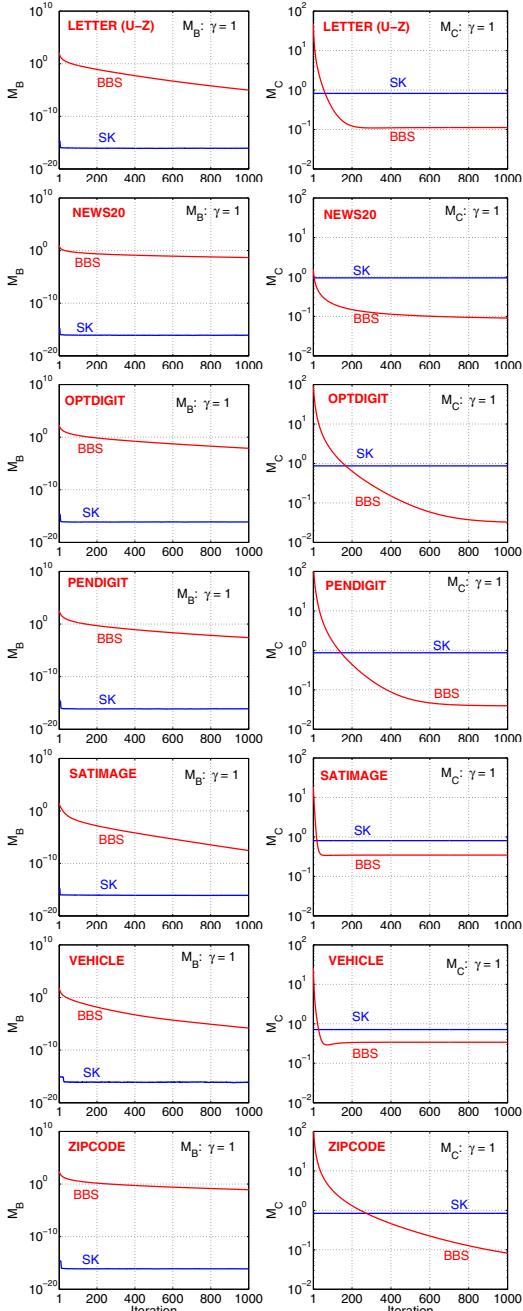


Figure 2. Quality measurements,  $M_B$ ,  $M_C$ , on a variety of data sets.

#### D. Comparing Clustering Results

We ultimately rely on the standard clustering procedure, e.g., [15], to assess clustering quality. Tables II to V provide the results for **BBS** (Alg. 1), **SK**, and three other methods:

- **K-means:** We directly used the original data sets (after normalizing each data point to have a unit  $l_2$  norm) and ran Matlab *kmeans* 100 times.
- **RA:** We ran spectral clustering directly on the similarity matrix  $\mathbf{K}$  (4). It is called *Ratio Association* [6].
- **Neut:** We ran spectral clustering on the normalized similarity matrix  $\tilde{\mathbf{K}} = \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2}$ , as in Eq. (5).

We report the clustering results on two metrics:

#### 1) Clustering Accuracy:

$$Accuracy = \frac{1}{n} \max \left( \sum_{\pi_i, \hat{\pi}_j} |\pi_i \cap \hat{\pi}_j| \right), \quad (34)$$

where  $\hat{\pi}_j$  denotes the  $j$ -th cluster in the output,  $\pi_i$  is the true  $i$ -th class, and  $|\pi_i \cap \hat{\pi}_j|$  is the number of data points from the  $i$ -th class are assigned to  $j$ -th cluster.

#### 2) Normalized Mutual Information (NMI) [21]:

$$NMI = \frac{\sum_{i=1}^k \sum_{j=1}^k |\pi_i \cap \hat{\pi}_j| \log \left( \frac{n \cdot |\pi_i \cap \hat{\pi}_j|}{|\pi_i| \cdot |\hat{\pi}_j|} \right)}{\sqrt{\left( \sum_{i=1}^k |\pi_i| \log \frac{|\pi_i|}{n} \right) \left( \sum_{j=1}^k |\hat{\pi}_j| \log \frac{|\hat{\pi}_j|}{n} \right)}} \quad (35)$$

We still need to address two more issues:

- For each case, we always ran *kmeans* 100 times. We report both the *average* and *maximum* measures of clustering quality (*Accuracy* and *NMI*). In practice, the *maximum* clustering performance may be quite attainable by tuning and running *kmeans* many times with different (random) initial starts.
- For **RA**, **Neut**, **SK** and **BBS**, we experimented with the similarity matrices  $\mathbf{K}$  (4) generated from a series of  $\gamma$  values (from 0.25 to 1024). Tables II to V report the best results among all  $\gamma$ 's. Again, the rationale is that, in practice, the best performance may be attainable by careful tuning. In addition, we believe it is also informative to present the results for all  $\gamma$  values, as in the Appendix; although due to the space limit, we could not present the results for all data sets.

Tables II to V demonstrate that, for many data sets, **BBS** (Alg. 1) can achieve considerably better clustering results than other methods, especially when evaluated using *maximum accuracy* and *maximum NMI*.

Table II  
AVERAGE ACCURACY

Data	K-means	RA	Neut	SK	BBS
MNIST	0.536	0.552	0.545	0.542	<b>0.633</b>
MNIST (0-4)	0.744	0.722	0.722	0.721	<b>0.805</b>
ISOLET (A-I)	0.621	<b>0.737</b>	0.735	0.709	0.713
ISOLET (J-R)	0.662	0.706	0.705	0.702	<b>0.708</b>
ISOLET (S-Z)	0.703	<b>0.787</b>	0.739	0.742	0.773
LETTER (A-E)	0.462	0.516	0.516	0.513	<b>0.539</b>
LETTER (F-J)	0.514	0.490	0.492	0.495	<b>0.619</b>
LETTER (K-O)	0.390	0.474	0.473	0.470	<b>0.502</b>
LETTER (P-T)	0.426	0.554	0.554	<b>0.555</b>	0.554
LETTER (U-Z)	0.467	0.511	<b>0.517</b>	0.512	0.505
NEWS20	0.277	0.244	0.245	0.244	<b>0.378</b>
OPTDIGIT	0.750	0.791	0.767	0.762	<b>0.848</b>
PENDIGIT	0.703	0.730	0.732	0.733	<b>0.756</b>
SATIMAGE	0.607	0.565	0.569	0.573	<b>0.617</b>
SHUTTLE	0.464	0.384	0.448	0.453	<b>0.647</b>
VEHICLE	0.366	0.389	0.371	0.374	<b>0.409</b>
ZIPCODE	0.650	0.678	0.678	0.674	<b>0.747</b>
ZIPCODE (0-4)	0.760	0.686	0.680	0.684	<b>0.908</b>

Table III  
AVERAGE NMI

Data	K-means	RA	Ncut	SK	BBS
MNIST	0.523	0.517	0.524	0.507	<b>0.711</b>
MNIST (0-4)	0.670	0.638	0.652	0.667	<b>0.850</b>
ISOLET (A-I)	0.711	0.756	0.755	0.746	<b>0.808</b>
ISOLET (J-R)	0.762	0.760	0.760	0.745	<b>0.832</b>
ISOLET (S-Z)	0.790	0.803	0.788	0.781	<b>0.843</b>
LETTER (A-E)	0.320	0.348	0.350	0.347	<b>0.397</b>
LETTER (F-J)	0.379	0.337	0.342	0.352	<b>0.469</b>
LETTER (K-O)	0.265	0.260	0.262	0.254	<b>0.379</b>
LETTER (P-T)	0.263	0.371	0.372	0.373	<b>0.417</b>
LETTER (U-Z)	0.344	0.397	0.403	0.399	<b>0.437</b>
NEWS20	0.283	0.241	0.241	0.238	<b>0.422</b>
OPTDIGIT	0.728	0.748	0.725	0.709	<b>0.874</b>
PENDIGIT	0.691	0.693	0.693	0.705	<b>0.776</b>
SATIMAGE	0.549	0.473	0.491	0.494	<b>0.603</b>
SHUTTLE	0.496	0.396	0.429	0.448	<b>0.542</b>
VEHICLE	0.116	0.108	0.124	0.123	<b>0.168</b>
ZIPCODE	0.631	0.625	0.625	0.624	<b>0.815</b>
ZIPCODE (0-4)	0.716	0.703	0.712	0.700	<b>0.913</b>

Table IV  
MAXIMUM ACCURACY

Data	K-means	RA	Ncut	SK	BBS
MNIST	0.588	0.608	0.603	0.603	<b>0.738</b>
MNIST (0-4)	0.853	0.756	0.790	0.827	<b>0.960</b>
ISOLET (A-I)	0.738	0.798	0.798	0.798	<b>0.819</b>
ISOLET (J-R)	0.760	0.750	0.746	0.746	<b>0.781</b>
ISOLET (S-Z)	0.861	0.846	0.870	0.794	<b>0.933</b>
LETTER (A-E)	0.518	0.589	0.589	0.589	<b>0.595</b>
LETTER (F-J)	0.590	0.584	0.584	0.546	<b>0.649</b>
LETTER (K-O)	0.463	0.487	0.500	0.510	<b>0.560</b>
LETTER (P-T)	0.496	0.604	0.556	0.556	<b>0.621</b>
LETTER (U-Z)	0.532	0.567	0.558	0.558	<b>0.585</b>
NEWS20	0.294	0.264	0.262	0.259	<b>0.419</b>
OPTDIGIT	0.875	0.814	0.824	0.801	<b>0.911</b>
PENDIGIT	0.778	0.795	0.794	0.820	<b>0.857</b>
SATIMAGE	0.632	0.582	0.588	0.590	<b>0.639</b>
SHUTTLE	0.598	0.474	0.506	0.510	<b>0.861</b>
VEHICLE	0.402	0.395	0.382	0.382	<b>0.479</b>
ZIPCODE	0.756	0.740	0.739	0.731	<b>0.897</b>
ZIPCODE (0-4)	0.891	0.813	0.808	0.809	<b>0.991</b>

Table V  
MAXIMUM NMI

Data	K-means	RA	Ncut	SK	BBS
MNIST	0.567	0.542	0.567	0.541	<b>0.741</b>
MNIST (0-4)	0.696	0.641	0.659	0.680	<b>0.897</b>
ISOLET (A-I)	0.788	0.779	0.781	0.770	<b>0.862</b>
ISOLET (J-R)	0.812	0.800	0.792	0.786	<b>0.883</b>
ISOLET (S-Z)	0.874	0.843	0.878	0.806	<b>0.897</b>
LETTER (A-E)	0.392	0.422	0.422	0.422	<b>0.468</b>
LETTER (F-J)	0.435	0.412	0.412	0.406	<b>0.524</b>
LETTER (K-O)	0.306	0.288	0.301	0.298	<b>0.430</b>
LETTER (P-T)	0.378	0.422	0.375	0.377	<b>0.499</b>
LETTER (U-Z)	0.395	0.445	0.437	0.437	<b>0.502</b>
NEWS20	0.310	0.254	0.254	0.252	<b>0.433</b>
OPTDIGIT	0.786	0.758	0.755	0.750	<b>0.897</b>
PENDIGIT	0.718	0.717	0.719	0.734	<b>0.826</b>
SATIMAGE	0.627	0.483	0.511	0.511	<b>0.623</b>
SHUTTLE	0.563	0.516	0.507	0.489	<b>0.705</b>
VEHICLE	0.172	0.125	0.150	0.144	<b>0.234</b>
ZIPCODE	0.665	0.652	0.649	0.645	<b>0.871</b>
ZIPCODE (0-4)	0.755	0.705	0.712	0.706	<b>0.964</b>

#### IV. RELATIONSHIP TO $k$ -MEANS

It is beneficial to gain some intuitive understanding on why the BBS algorithm with  $\phi(x) = x^2/2$  (i.e., Alg. 1) can perform well in clustering. We show that it is closely related to various relaxed  $k$ -means algorithms.

The  $k$ -means clustering aims to minimize the objective

$$J_1 = \sum_{c=1}^k \sum_{x_i \in \pi_c} \|x_i - \mu_c\|^2 \quad (36)$$

where  $\mu_c$  is the mean of cluster  $\pi_c$ . Some algebra can show that the minimizing  $J_1$  is equivalent to minimizing  $J_2$ :

$$J_2 = -\text{tr}(\tilde{\mathbf{F}}^\top \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{F}}) \quad (37)$$

where  $\tilde{\mathbf{F}}$  is the scaled partition matrix introduced at the beginning of the paper and  $\mathbf{X} = [x_1, x_2, \dots, x_n]^\top \in \mathbb{R}^{n \times d}$  is the data matrix. Let  $\mathbf{G} = \tilde{\mathbf{F}} \tilde{\mathbf{F}}^\top$  and  $\mathbf{K} = \mathbf{X} \mathbf{X}^\top$ . Then  $J_2 = -\text{tr}(\mathbf{KG})$ , which in fact can be viewed as a special case of the objective of BBS defined in Eq. (10):  $\text{tr}(\mathbf{G}^\top \mathbf{G} - 2\mathbf{K}^\top \mathbf{G})$ , because the term  $\text{tr}(\mathbf{G}^\top \mathbf{G})$  can be treated as a constant in this case:

$$\text{tr}(\mathbf{G}^\top \mathbf{G}) = \text{tr}(\tilde{\mathbf{F}} \tilde{\mathbf{F}}^\top \tilde{\mathbf{F}} \tilde{\mathbf{F}}^\top) = \text{tr}(\tilde{\mathbf{F}} \tilde{\mathbf{F}}^\top) = \text{tr}(\mathbf{G}) = k$$

In addition,  $\mathbf{K} = \mathbf{X} \mathbf{X}^\top$  is the linear kernel, which may be replaced by more flexible kernels, e.g., Eq. (4) as we use.

There are more than one way to formulate the relaxed  $k$ -means algorithm. For example,

$$\min_{\mathbf{G}} D_\phi(\mathbf{G}, \mathbf{K}), \quad (\text{where } \phi(x) = x^2) \quad (38)$$

$$\text{s.t. } \mathbf{G} \geq 0, \quad \mathbf{G} = \mathbf{G}^\top, \quad \mathbf{G}\mathbf{1} = \mathbf{1},$$

$$\mathbf{G}^2 = \mathbf{G}, \quad \text{tr}(\mathbf{G}) = k, \quad (39)$$

which is quite similar to our formulation of the BBS problem with the Euclidian distance. Our formulation discards the constraints (39) and hence its optimization task is easier.

#### V. EXTENSION: MULTIPLE BBS (MBBS)

Our detailed experiments reported in the Appendix illustrate that the clustering performance of the BBS algorithm (as well as other algorithms), to an extent, depends on the initial similarity matrix  $\mathbf{K}$ . This section extends BBS to combine the power of multiple input similarity matrices, e.g., a series of kernel matrices (4) using different  $\gamma$  values, to boost the performance. We name this scheme *Multiple BBS* or *MBBS*. This is in spirit related to *cluster ensemble* [21] and *Generalized Cluster Aggregation* [22].

Suppose we have  $m$  similarity matrices  $\{\mathbf{K}_{(i)}\}_{i=1}^m$ . We would like to obtain a bi-stochastic similarity matrix  $\mathbf{G}$  by solving the following optimization problem:

$$\min_{\mathbf{G}, \alpha} \sum_{i=1}^m \alpha_i D_\phi(\mathbf{G}, \mathbf{K}_{(i)}) + \lambda \Omega(\alpha) \quad (40)$$

$$\text{s.t. } \mathbf{G} \geq 0, \quad \mathbf{G} = \mathbf{G}^\top, \quad \mathbf{G}\mathbf{1} = \mathbf{1},$$

$$\forall i, \alpha_i \geq 0, \quad \sum_{i=1}^m \alpha_i = 1$$

We constrain the weight coefficients  $\alpha = \{\alpha_i\}_{i=1}^m$  to be in a simplex.  $\Omega(\alpha)$  is some regularizer to avoid trivial solutions.

There are two groups of variables  $\alpha$  and  $\mathbf{G}$ . Although the problem (40) is not jointly convex, it is convex with respect to one group of variables with the other group being fixed. Thus, it is reasonable to apply *block coordinate descent* [1].

### A. Fix $\alpha$ , Solve $\mathbf{G}$

At the  $t$ -th iteration, if  $\alpha$  is fixed to be  $\alpha = \alpha^{(t-1)}$ , the problem (40) becomes

$$\begin{aligned} \min_{\mathbf{G}} \quad & \sum_{i=1}^m \alpha_i^{(t-1)} D_\phi(\mathbf{G}, \mathbf{K}_{(i)}) \\ \text{s.t.} \quad & \mathbf{G} \geq 0, \quad \mathbf{G} = \mathbf{G}^\top, \quad \mathbf{G}\mathbf{1} = \mathbf{1}. \end{aligned} \quad (41)$$

Note that  $\Omega(\alpha)$  is irrelevant at this point. This is similar to problem (7) except for the summation form in the objective. The solution procedures are consequently also similar.

Here we assume  $\phi(x) = x^2/2$  for the illustration purpose.

$$\begin{aligned} & \sum_{i=1}^m \alpha_i^{(t-1)} D_\phi(\mathbf{G}, \mathbf{K}_{(i)}) \\ &= \frac{1}{2} \operatorname{tr} \left( \sum_{i=1}^m \alpha_i^{(t-1)} \mathbf{K}_{(i)}^\top \mathbf{K}_{(i)} + \mathbf{G}^\top \mathbf{G} - 2 \sum_{i=1}^m \alpha_i^{(t-1)} \mathbf{K}_{(i)}^\top \mathbf{G} \right) \end{aligned}$$

where we use the fact  $\sum_{i=1}^m \alpha_i^{(t-1)} = 1$ . As the term  $\sum_{i=1}^m \alpha_i^{(t-1)} \mathbf{K}_{(i)}^\top \mathbf{K}_{(i)}$  is irrelevant, the problem becomes

$$\begin{aligned} \min_{\mathbf{G}} \quad & \operatorname{tr} \left( \mathbf{G}^\top \mathbf{G} - 2 \left( \sum_{i=1}^m \alpha_i \mathbf{K}_{(i)} \right)^\top \mathbf{G} \right) \\ \text{s.t.} \quad & \mathbf{G} \geq 0, \quad \mathbf{G} = \mathbf{G}^\top, \quad \mathbf{G}\mathbf{1} = \mathbf{1} \end{aligned} \quad (42)$$

which is the same as Problem (10) if we make

$$\mathbf{K} = \sum_{i=1}^m \alpha_i^{(t-1)} \mathbf{K}_{(i)}. \quad (43)$$

### B. Fix $\mathbf{G}$ , Solve $\alpha$

When  $\mathbf{G}$  is fixed with  $\mathbf{G} = \mathbf{G}^{(t)}$  and for simplicity we only consider  $\Omega(\alpha) = \|\alpha\|^2 = \alpha^\top \alpha$ , the problem becomes

$$\begin{aligned} \min_{\mathbf{G}, \alpha} \quad & \sum_{i=1}^m \alpha_i D_\phi(\mathbf{G}^{(t)}, \mathbf{K}_i) + \lambda \|\alpha\|^2 \\ \text{s.t.} \quad & \forall i, \alpha_i \geq 0, \quad \sum_{i=1}^m \alpha_i = 1, \end{aligned} \quad (44)$$

which is a standard *Quadratic Programming* (QP) problem.

Here we will reformulate this problem to facilitate more efficient solutions. For the notational convenience, we denote  $\mathbf{g}^{(t)} = (g_1^{(t)}, g_2^{(t)}, \dots, g_m^{(t)})^\top$  with

$$g_i^{(t)} = D_\phi(\mathbf{G}^{(t)}, \mathbf{K}_{(i)}) \quad (45)$$

We first rewrite the objective of Problem (44) as

$$\alpha^\top \mathbf{g}^{(t)} + \lambda \|\alpha\|^2 = \left\| \sqrt{\lambda} \alpha - \frac{1}{\sqrt{2\lambda}} \mathbf{g}^{(t)} \right\|^2 + \frac{1}{2\lambda} (\mathbf{g}^{(t)})^\top \mathbf{g}^{(t)}$$

As  $\frac{1}{2\lambda} (\mathbf{g}^{(t)})^\top \mathbf{g}^{(t)}$  is irrelevant, (44) can be rewritten to be

$$\min_{\alpha} \left\| \alpha - \frac{1}{\sqrt{2\lambda}} \mathbf{g}^{(t)} \right\|^2, \quad \text{s.t. } \alpha \geq 0, \quad \alpha^\top \mathbf{1} = 1, \quad (46)$$

which is an Euclidian projection problem under the simplex constraint and can be solved efficiently, e.g., [9][14].

We will report extensive experiment results of Multiple BBS in a more comprehensive technical report.

## VI. CONCLUSIONS

We present **BBS** (*Bregmanian Bi-Stochastication*), a general framework for learning a bi-stochastic data similarity matrix from an initial similarity matrix, by minimizing the Bregmanian divergences such as the Euclidian distance or the KL divergence. The resultant bi-stochastic matrix can be used as input to clustering algorithms. The BBS framework is closely related to the relaxed  $k$ -means algorithms. Our extensive experiments on a wide range of public data sets demonstrate that the BBS algorithm using the Euclidian distance can often produce noticeably superior clustering results than other well-known algorithms including the **SK** algorithm and the **Ncut** algorithm.

## ACKNOWLEDGEMENT

This work is partially supported by NSF (DMS-0808864), ONR (YIP-N000140910911), and a grant from Microsoft.

## REFERENCES

- [1] D. P. Bertsekas. *Nonlinear Programming: 2nd Edition*. Athena Scientific, 1999.
- [2] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK., 2004.
- [3] P. K. Chan, D. F. Schlag, and J. Y. Zien. Spectral k-way ratio-cut partitioning and clustering. *IEEE Trans. Computer-Aided Design*, 13:1088–1096, 1994.
- [4] J. N. Darroch and D. Ratcliff. Generalized iterative scaling for log-linear models. *The Annals of Mathematical Statistics*, 43(5):1470–1480, 1972.
- [5] W. E. Deming and F. F. Stephan. On a least squares adjustment of a sampled frequency table when the expected marginal totals are known. *The Annals of Mathematical Statistics*, 11(4):427–444, 1940.
- [6] I. S. Dhillon, Y. Guan, and B. Kulis. A unified view of kernel k-means, spectral clustering and graph cuts. Technical report, Department of Computer Science, University of Texas at Austin. TR-04-25, 2004.
- [7] I. S. Dhillon and J. A. Tropp. Matrix nearness problems with bregman divergences. In *SIAM Journal on Matrix Analysis and Applications*, volume 29, pages 1120–1146, 2008.
- [8] C. Ding, X. He, H. Zha, M. Gu, and H. D. Simon. A min-max cut algorithm for graph partitioning and data clustering. In *Proceedings of the 1st International Conference on Data Mining*, pages 107–114, 2001.
- [9] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the L1-ball for learning in high dimensions. In *Proceedings of the 25th international conference on Machine learning*, pages 272–279, 2008.
- [10] R. Escalante and M. Raydan. Dykstra’s algorithm for a constrained least-squares matrix problem. *Numerical Linear Algebra with Applications*, 3(6):459–471, 1998.

- [11] W. W. Hager. Updating the inverse of a matrix. *SIAM Review*, 31(2):221–239, 1989.
- [12] A. Horn. Doubly stochastic matrices and the diagonal of a rotation matrix. *The American Journal of Mathematics*, 76:620–630, 1954.
- [13] A. K. Jain and R. C. Dubes. *Algorithms for Clustering Data*. Prentice Hall, 1988.
- [14] J. Liu and J. Ye. Efficient Euclidean projections in linear time. In *International Conference on Machine Learning*, pages 657–664, 2009.
- [15] A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: analysis and an algorithm. In *Advances in Neural Information Processing Systems 14*, pages 849–856, 2001.
- [16] J. Nocedal and S. J. Wright. *Numerical Optimization (2nd ed.)*. Springer-Verlag, Berlin, New York, 2006.
- [17] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- [18] R. Sinkhorn and P. Knopp. Concerning nonnegative matrices and doubly stochastic matrices. *Pacific J. Math.*, 21:343–348, 1967.
- [19] G. W. Soules. The rate of convergence of Sinkhorn balancing. *Linear Algebra and its Applications*, 150:3 – 40, 1991.
- [20] F. F. Stephan. An iterative method of adjusting sample frequency tables when expected marginal totals are known. *The Annals of Mathematical Statistics*, 13(2):166–178, 1942.
- [21] A. Strehl and J. Ghosh. Cluster ensembles - a knowledge reuse framework for combining multiple partitions. *Journal of Machine Learning Research*, 3:583–617, 2002.
- [22] F. Wang, X. Wang, and T. Li. Generalized cluster aggregation. In *Proceedings of The 21st International Joint Conference on Artificial Intelligence*, pages 1279–1284, 2009.
- [23] R. Zass and A. Shashua. A unifying approach to hard and probabilistic clustering. In *Proceedings of International Conference on Computer Vision*, pages 294–301, 2005.
- [24] H. Zha, X. He, C. Ding, M. Gu, and H. Simon. Spectral relaxation for k-means clustering. In *Advances in Neural Information Processing Systems 14*, 2001.

## APPENDIX

We generated the base similarity matrix  $\mathbf{K}$  using the Gaussian kernel (4) which has a tuning parameter  $\gamma > 0$ . The clustering performance can be, to an extent, sensitive to  $\gamma$ ; and hence we would like to present the clustering results for  $\gamma$  values ranging from  $2^{-2} = 0.25$  to  $2^{10} = 1024$ , for four algorithms: **RA**, **Neut**, **SK**, and **BBS** (using Euclidian distance), and two performance measures: *Accuracy* and *NMI*, as defined in Eq. (34) and Eq (35), respectively.

In the tables, each entry contains the *average* and *maximum* (in parentheses) clustering results from 100 runs of the Matlab *kmeans* program. Due to the space limit, we could not present the experiments for all the data sets.

		MNIST:		Accuracy
$\gamma$	RA	Ncut	SK	BBS
1024	0.532 (0.595)	0.533 (0.596)	0.536 (0.595)	0.553 (0.606)
256	0.536 (0.596)	0.537 (0.594)	0.539 (0.594)	0.566 (0.617)
64	0.541 (0.608)	0.540 (0.596)	0.530 (0.600)	0.581 (0.642)
32	0.538 (0.596)	0.537 (0.596)	0.537 (0.596)	0.612 (0.665)
16	0.537 (0.606)	0.532 (0.598)	0.538 (0.600)	0.630 (0.690)
8	0.539 (0.594)	0.536 (0.594)	0.541 (0.598)	0.622 (0.702)
4	0.537 (0.600)	0.542 (0.596)	0.538 (0.598)	0.630 (0.715)
2	0.540 (0.596)	0.542 (0.599)	0.542 (0.603)	0.633 (0.720)
1	0.548 (0.597)	0.545 (0.603)	0.536 (0.597)	0.621 (0.738)
0.5	0.552 (0.601)	0.535 (0.599)	0.521 (0.597)	0.463 (0.519)
0.25	0.527 (0.590)	0.514 (0.595)	0.493 (0.558)	0.451 (0.472)
		MNIST:		NMI
$\gamma$	RA	Ncut	SK	BBS
1024	0.473 (0.512)	0.475 (0.511)	0.476 (0.511)	0.491 (0.522)
256	0.475 (0.512)	0.475 (0.510)	0.476 (0.510)	0.530 (0.564)
64	0.479 (0.512)	0.478 (0.510)	0.474 (0.510)	0.593 (0.620)
32	0.478 (0.512)	0.476 (0.506)	0.476 (0.511)	0.626 (0.648)
16	0.475 (0.504)	0.475 (0.511)	0.476 (0.511)	0.651 (0.675)
8	0.479 (0.511)	0.476 (0.511)	0.477 (0.511)	0.668 (0.692)
4	0.475 (0.507)	0.479 (0.512)	0.479 (0.513)	0.686 (0.705)
2	0.481 (0.516)	0.483 (0.516)	0.485 (0.518)	0.700 (0.724)
1	0.494 (0.524)	0.493 (0.521)	0.492 (0.511)	0.711 (0.741)
0.5	0.508 (0.526)	0.510 (0.528)	0.507 (0.541)	0.467 (0.489)
0.25	0.517 (0.542)	0.524 (0.567)	0.500 (0.525)	0.386 (0.400)
		MNIST 0-4:		Accuracy
$\gamma$	RA	Ncut	SK	BBS
1024	0.721 (0.721)	0.719 (0.721)	0.720 (0.721)	0.722 (0.724)
256	0.720 (0.721)	0.719 (0.721)	0.719 (0.721)	0.673 (0.780)
64	0.720 (0.722)	0.722 (0.722)	0.721 (0.722)	0.762 (0.896)
32	0.721 (0.722)	0.721 (0.722)	0.721 (0.721)	0.762 (0.879)
16	0.722 (0.722)	0.721 (0.721)	0.720 (0.721)	0.757 (0.885)
8	0.721 (0.721)	0.718 (0.721)	0.715 (0.719)	0.760 (0.881)
4	0.718 (0.718)	0.716 (0.716)	0.714 (0.715)	0.777 (0.907)
2	0.703 (0.707)	0.698 (0.703)	0.692 (0.692)	0.762 (0.921)
1	0.675 (0.705)	0.675 (0.708)	0.675 (0.709)	0.805 (0.960)
0.5	0.661 (0.756)	0.658 (0.757)	0.645 (0.742)	0.438 (0.519)
0.25	0.506 (0.584)	0.708 (0.790)	0.721 (0.827)	0.532 (0.565)
		MNIST 0-4:		NMI
$\gamma$	RA	Ncut	SK	BBS
1024	0.572 (0.575)	0.571 (0.575)	0.572 (0.575)	0.577 (0.579)
256	0.572 (0.576)	0.571 (0.575)	0.571 (0.575)	0.635 (0.637)
64	0.573 (0.576)	0.575 (0.576)	0.574 (0.576)	0.741 (0.779)
32	0.576 (0.576)	0.575 (0.576)	0.575 (0.576)	0.764 (0.789)
16	0.577 (0.577)	0.577 (0.577)	0.576 (0.576)	0.785 (0.811)
8	0.579 (0.579)	0.578 (0.579)	0.577 (0.577)	0.800 (0.823)
4	0.585 (0.585)	0.583 (0.583)	0.580 (0.582)	0.815 (0.837)
2	0.589 (0.590)	0.587 (0.589)	0.590 (0.590)	0.813 (0.853)
1	0.614 (0.614)	0.609 (0.611)	0.602 (0.603)	0.850 (0.897)
0.5	0.638 (0.641)	0.632 (0.640)	0.636 (0.648)	0.346 (0.376)
0.25	0.560 (0.570)	0.652 (0.659)	0.667 (0.680)	0.418 (0.454)
		ISOLET A-I:		Accuracy
$\gamma$	RA	Ncut	SK	BBS
1024	0.695 (0.769)	0.688 (0.768)	0.692 (0.769)	0.691 (0.768)
256	0.694 (0.769)	0.697 (0.769)	0.695 (0.769)	0.681 (0.745)
64	0.692 (0.769)	0.688 (0.769)	0.692 (0.769)	0.657 (0.743)
32	0.699 (0.769)	0.702 (0.769)	0.693 (0.769)	0.658 (0.740)
16	0.687 (0.769)	0.684 (0.769)	0.689 (0.769)	0.682 (0.813)
8	0.697 (0.769)	0.700 (0.769)	0.694 (0.769)	0.679 (0.812)
4	0.695 (0.769)	0.692 (0.769)	0.684 (0.768)	0.674 (0.812)
2	0.732 (0.781)	0.735 (0.777)	0.694 (0.769)	0.678 (0.819)
1	0.737 (0.798)	0.726 (0.798)	0.709 (0.798)	0.678 (0.779)
0.5	0.715 (0.783)	0.705 (0.743)	0.700 (0.748)	0.713 (0.787)
0.25	0.677 (0.747)	0.687 (0.750)	0.552 (0.629)	0.590 (0.639)
		ISOLET A-I:		NMI
$\gamma$	RA	Ncut	SK	BBS
1024	0.677 (0.746)	0.675 (0.746)	0.673 (0.745)	0.680 (0.745)
256	0.681 (0.745)	0.683 (0.745)	0.682 (0.746)	0.727 (0.762)
64	0.679 (0.745)	0.677 (0.745)	0.675 (0.745)	0.720 (0.773)
32	0.681 (0.746)	0.682 (0.745)	0.677 (0.746)	0.731 (0.785)
16	0.678 (0.749)	0.675 (0.745)	0.680 (0.745)	0.762 (0.822)
8	0.686 (0.748)	0.684 (0.748)	0.678 (0.746)	0.773 (0.826)
4	0.690 (0.754)	0.684 (0.747)	0.675 (0.742)	0.773 (0.832)
2	0.732 (0.756)	0.723 (0.749)	0.696 (0.742)	0.771 (0.832)
1	0.756 (0.779)	0.751 (0.777)	0.745 (0.769)	0.762 (0.791)
0.5	0.749 (0.777)	0.755 (0.775)	0.746 (0.770)	0.808 (0.862)
0.25	0.722 (0.770)	0.743 (0.781)	0.624 (0.655)	0.586 (0.614)



