

Technical Report No. 177

Approximation Algorithms for
Bregman Clustering
Co-clustering and Tensor
Clustering

Suvrit Sra¹, Stefanie Jegelka¹
Arindam Banerjee²

17 Oct 2008

¹ MPI für biologische Kybernetik, AGBS; ² University of Minnesota, MN, USA.

Approximation Algorithms for Bregman Clustering Co-clustering and Tensor Clustering

Suvrit Sra, Stefanie Jegelka, and Arindam Banerjee

Abstract. The Euclidean K-means problem is fundamental to clustering and over the years it has been intensely investigated. More recently, generalizations such as Bregman k-means [8], co-clustering [10], and tensor (multi-way) clustering [40] have also gained prominence. A well-known computational difficulty encountered by these clustering problems is the NP-Hardness of the associated optimization task, and commonly used methods guarantee at most local optimality. Consequently, approximation algorithms of varying degrees of sophistication have been developed, though largely for the basic Euclidean K-means (or ℓ_1 -norm K-median) problem. In this paper we present approximation algorithms for several Bregman clustering problems by building upon the recent paper of Arthur and Vassilvitskii [5]. Our algorithms obtain objective values within a factor $O(\log K)$ for Bregman k-means, Bregman co-clustering, Bregman tensor clustering, and weighted kernel k-means. To our knowledge, except for some special cases, approximation algorithms have not been considered for these general clustering problems. There are several important implications of our work: (i) under the same assumptions as Ackermann et al. [2] it yields a much faster algorithm (non-exponential in K , unlike [2]) for information-theoretic clustering, (ii) it answers several open problems posed by [4], including generalizations to Bregman co-clustering, and tensor clustering, (iii) it provides practical and easy to implement methods—in contrast to several other common approximation approaches.

1 Introduction

Partitioning data points into clusters is a fundamentally hard problem. The well-known Euclidean k-means problem that seeks to partition the input data into K clusters, so that the sum of squared distances of the input points to their corresponding cluster centroids is minimized, is an NP-Hard problem [22]. Simple and frequently used procedures that rapidly obtain local minima exist since a long time [26, 32]. For example, Lloyd’s algorithm [32], which is commonly referred to as the K-means algorithm is arguably the most popular approach to solving Euclidean k-means. Here, one begins with K centers (usually chosen randomly) and assigns points to their closest centers. Each cluster center is then recomputed as the mean of the points assigned to it, and these two steps are repeated until the procedure converges. A similar greedy procedure also exists for the Bregman k-means problem, as shown in [8]. Despite enjoying properties such as monotonic descent in the objective function value and utter simplicity of implementation, these simplistic iterative approaches can often get stuck in poor local-optima. Therefore, heuristic local search strategies (e.g., [21]), or even guaranteed approximation algorithms have been designed for it (e.g., [31] or references therein). Heuristic strategies can be quite effective but are not accompanied by better than local optimality guarantees, while standard approximation algorithms quickly sacrifice the simplicity, and thereby the efficiency of the K-means algorithm.

Fortunately, in a recent paper Arthur and Vassilvitskii [5] presented a simple initialization scheme for Euclidean k-means along with an elegant analysis guaranteeing an $O(\log K)$ approximation to the globally optimal objective function value. The greatest advantage of their scheme is that it retains the simplicity and efficiency of the K-means algorithm, while still maintaining theoretical guarantees. This

paper is directly motivated by their work, which we greatly extend to obtain approximation algorithms for several Bregman clustering problems. We summarize our main results below.

1.1 Results.

We present approximation algorithms for the following Bregman divergence based clustering problems:

1. Bregman k-means [8] (§2),
2. Bregman co-clustering [10] (§5),
3. Bregman tensor clustering [9] (§5).

Additionally as an easy generalization of [5] we also obtain an approximation algorithm for weighted kernel k-means [20] (§4).

Implications. Our results have several important implications. Under assumptions similar to that of (Ackermann et al. [2], 2008), we obtain a much faster approximation algorithm for information-theoretic clustering as a special case of our approximation for Bregman k-means (§3.1). Ackermann et al. [2] require time exponential (or worse) in K , while our methods run in time linear in K . In fact, while preparing this paper we became aware of a very recent SODA 2009 paper of Ackermann and Blömer [1] (yet to appear in print), who provide new approximation algorithms for Bregman k-means. However, their new algorithms, while faster than those in [2], are *still* exponential (or worse) in K —our algorithm operates under the *same* assumptions on the Bregman divergences as made by [1], and is much faster (non-exponential in K).

Our results for Bregman co-clustering and Bregman tensor clustering answer two open problems posed by [4], and yield the first (to our knowledge) known approximation algorithms for these problems. Finally, using our $O(\log K)$ approximation for weighted kernel K-means, one can obtain potentially better algorithms for graph-cut objectives and certain semi-supervised clustering problems by exploiting the equivalences described by [20, 30].

1.2 Related work

There exist several books and a vast array of papers dealing with the problem of clustering. However, as our focus is on approximation algorithms for Bregman divergence based clustering problems, we summarize below only work dealing with approximation algorithms for clustering. Graph partitioning also forms a large class of clustering problems and algorithms. However, it lies outside the scope of this paper, apart from the connection via weighted kernel k-means as mentioned above.

1.2.1 Clustering

The most directly related work is the paper [5] that has motivated our algorithms for clustering. If one fixes the number of clusters K , and the data dimensionality d , then Euclidean k-means can be solved exactly, in time $O(n^{Kd})$ [28]. We remark that using the Bregman Voronoi ideas of Nielsen et al. [35], it might be possible to generalize the work of Inaba et al. [28] to Bregman k-means.

Several other polynomial time approximation algorithms for K-means have been proposed in the literature, for example, [17, 25, 31] (also see the references therein). All of the algorithms proposed in these papers suffer from a common problem, namely exponential (or poly-exponential) dependence on K , rendering them impractical despite their theoretical pleasantness.

Of particular interest is the paper of Ackermann et al. [2], who extended clustering guarantees of Kumar et al. [31] to generic divergence measures (including Bregman divergences). Under the *same* assumptions on the underlying Bregman divergence as Ackermann et al. [2] we obtain a much faster and practical approximation algorithm for Bregman k-means than their methods. Their approximation factor is $(1 + \epsilon)$ with a running time of $O(dn2^{(K/\epsilon)^{O(1)}})$, while our factor is $O(\log K)$ with a running time of

$O(dnK)$. In an even more recent paper that will appear in SODA 2009, Ackermann and Blömer [1] (preprint available from the authors’ website) have introduced a new $O(dKn + d^2 2^{(K/\epsilon)^{\Theta(1)}} \log^{K+2} n)$ approximation algorithm for Bregman k-means that again achieves $(1 + \epsilon)$ approximation, with the same assumptions as [2] on the underlying Bregman divergences. Our approximation algorithms for Bregman k-means are much more practical than theirs, both because of the lower running time as well as the implementational simplicity.

Related to Bregman k-means is the important special case of information theoretic clustering, wherein one minimizes the sum of KL divergences of input data points to their cluster centroids. In addition to the generic methods already summarized above, recent important work worth mentioning here is the paper of Chaudhuri and McGregor [14], who present a KL-divergence clustering algorithm that does not make *any* assumptions on the input data, but yields a non-constant $O(\log n)$ approximation.

Other relevant works such as [29, 34, 36] are summarized in [1, 2, 5], and we refer the reader to those papers for additional information.

1.2.2 Co-clustering and Tensor clustering.

For a detailed discussion of co-clustering and several relevant references we refer the reader to [10], while for the lesser known problem of tensor clustering we refer the reader to [3, 9, 11, 23, 33, 40].

Approximation algorithms for co-clustering are much less well-studied. We are aware of only two very recent attempts (both papers are from 2008), namely, [38] and [4]—and both of the papers follow similar approaches to obtain their approximation guarantees. In this paper, we build upon [4] and obtain approximation algorithms for Bregman co-clustering as well as Bregman tensor clustering. We therefore answer *two open problems* posed by [4], namely, whether their methods for Euclidean co-clustering could be generalized to Bregman co-clustering, and more importantly, whether generalizations to tensors could be found. Our approximation results for co-clustering and tensor clustering may be viewed independently of our results for Bregman clustering, because they are based on being able to solve the 1-dimensional (standard) clustering problem with *any* guaranteed approximation method. One can, and we do, however, invoke our Bregman clustering results to obtain actual efficient algorithms.

Now we are ready to discuss details of our methods and we begin with Bregman clustering below.

2 Bregman Clustering

Bregman k-means (BREGM) was introduced by Banerjee et al. [8], and it can be viewed as a generalization of Euclidean k-means and information theoretic clustering (ITC) [19]. Below we derive a randomized algorithm called BREG++ for Bregman k-means and prove it to be within $O(\log K)$ of the optimal. In §3.1 we discuss the particularly interesting special case of ITC in further detail, especially because ITC has only recently (in 2008) witnessed some progress in terms of approximation algorithms [2, 14]. We also discuss implications of our BREG++ for mixture-modeling in §3.2.

2.1 Setup and Algorithm.

Let $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ be the input data, and $\{w_1, \dots, w_n\}$ corresponding non-negative weights. For a strictly convex function f , let B_f denote a Bregman divergence defined as [13]

$$B_f(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y}). \tag{2.1}$$

Given B_f , Bregman k-means seeks a partition $\mathcal{C} = \{\mathcal{C}_1, \dots, \mathcal{C}_K\}$ of \mathcal{X} , such that the following objective is minimized:

$$J(\mathcal{C}) = \sum_{h=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_h} w_i B_f(\mathbf{x}_i, \boldsymbol{\mu}_h), \tag{2.2}$$

where $\boldsymbol{\mu}_h$ is the weighted mean of cluster \mathcal{C}_h , i.e.,

$$\boldsymbol{\mu}_h = \frac{\sum_{\mathbf{x}_i \in \mathcal{C}_h} w_i \mathbf{x}_i}{\sum_{\mathbf{x}_i \in \mathcal{C}_h} w_i}. \quad (2.3)$$

The means given by (2.3) are optimal for a given clustering, as shown formally below.

Lemma 2.1. *Let \mathcal{A} be a set of points with weighted mean $\boldsymbol{\mu}_\mathcal{A}$, and let \mathbf{z} be an arbitrary point. Then,*

$$\sum_{\mathbf{x}_i \in \mathcal{A}} w_i B_f(\mathbf{x}_i, \mathbf{z}) = \sum_{\mathbf{x}_i \in \mathcal{A}} w_i B_f(\mathbf{x}_i, \boldsymbol{\mu}_\mathcal{A}) + W_\mathcal{A} B_f(\boldsymbol{\mu}_\mathcal{A}, \mathbf{z}),$$

where $W_\mathcal{A} = \sum_{\mathbf{x}_i \in \mathcal{A}} w_i$.

Proof. Follows directly from the definition (2.1) of B_f and equation (2.3)—for details see [8]. \square

Algorithm. Given K initial means, the Bregman k-means (BREGM) algorithm [8] follows the outline:

1. For each $i = 1 \dots N$, assign point \mathbf{x}_i to its nearest (in B_f) mean updating cluster \mathcal{C}_h .
2. For each $h = 1 \dots K$, compute $\boldsymbol{\mu}_h$ using (2.3).
3. Repeat steps 1 and 2 to convergence.

This simple k-means type approach monotonically decreases the objective function, finally stopping once the clusters stabilize. To obtain approximation guarantees, we must modify this basic algorithm a little. To that end, we generalize the careful initialization technique of [5], as shown below.

As its initialization, BREG++ selects cluster centers from \mathcal{X} sequentially following a weighted farthest-first scheme. The first center $\boldsymbol{\mu}_1$ is chosen with probability proportional to its weight, i.e., for some $\mathbf{x}_i \in \mathcal{X}$

$$P(\boldsymbol{\mu}_1 = \mathbf{x}_i) = \frac{w_i}{\sum_{j=1}^n w_j}.$$

The remaining centers are chosen from \mathcal{X} with a different weighting. At a given stage in the initialization, let C be the set of centers already chosen. Let $D(\mathbf{x})$ denote the smallest Bregman divergence of a point \mathbf{x} in \mathcal{X} to an already chosen center, i.e.,

$$D(\mathbf{x}) = \min_{\boldsymbol{\mu} \in C} B_f(\mathbf{x}, \boldsymbol{\mu}). \quad (2.4)$$

Then BREG++ chooses the next center $\boldsymbol{\mu}_h$ by letting $\boldsymbol{\mu}_h = \mathbf{x}_i \in \mathcal{X}$, with probability

$$P(\boldsymbol{\mu}_h = \mathbf{x}_i) = \frac{w_i D(\mathbf{x}_i)}{\sum_{j=1}^n w_j D(\mathbf{x}_j)}. \quad (2.5)$$

The initialization steps (2.4) and (2.5) are repeated until we have chosen K centers, which are then used to initialize the standard BREGM algorithm. Interestingly, this weighted farthest-first initialization alone is sufficient to bring BREG++ within a factor $O(\log K)$ of the optimal, as the analysis below shows.

2.2 Analysis.

Arthur and Vassilvitskii's [2007] analysis does not directly generalize to Bregman k-means. Some additional details must be developed as outlined in this section. We assume that the Bregman divergence being minimized has bounded curvature, i.e., $\exists \sigma_1, \sigma_2$ with $0 < \sigma_1 \leq \sigma_2 < \infty$, such that

$$\sigma_1 \|\mathbf{x} - \mathbf{y}\|^2 \leq B_f(\mathbf{x}, \mathbf{y}) \leq \sigma_2 \|\mathbf{x} - \mathbf{y}\|^2. \quad (2.6)$$

Note that the bounds (2.6) need not hold over the entire domain of f —they can be limited to convex hulls of the input data points. Specifically, we can select

$$\sigma_1 = \inf_{\substack{\mathbf{x} \in \mathcal{X} \\ \mathbf{y} \in \text{conv}(\mathcal{X})}} \frac{B_f(\mathbf{x}, \mathbf{y})}{\|\mathbf{x} - \mathbf{y}\|^2}, \quad \sigma_2 = \sup_{\substack{\mathbf{x} \in \text{conv}(\mathcal{X}) \\ \mathbf{y} \in \mathcal{X}}} \frac{B_f(\mathbf{x}, \mathbf{y})}{\|\mathbf{x} - \mathbf{y}\|^2}, \quad (2.7)$$

where $\text{conv}(\mathcal{X})$ denotes the convex hull of \mathcal{X} , i.e., the set of points that can be expressed as $\mathbf{y} = \sum_{\mathbf{x} \in \mathcal{X}} \alpha_{\mathbf{x}} \mathbf{x}$ with $\alpha_{\mathbf{x}} \geq 0$ and $\sum_{\mathbf{x} \in \mathcal{X}} \alpha_{\mathbf{x}} = 1$. Though these bounds might appear restrictive, they are in fact not that limiting, as our treatment of information theoretic clustering in Section 3.1 shows. In fact, it turns out the similar bounds were assumed in the very recent work of Ackermann et al. [2], and [1]. In the language of convex-optimization, these bounds are nothing but bounds on curvature (Hessian) of the convex function B_f (e.g., in the context of *strong convexity* [12, §9.1.2]).

On a more intriguing note, it seems that without such assumptions on the curvature of B_f , one might not be able to obtain constant approximation ratios; this intuition is reinforced by the recent results of [14], who avoided making such assumptions, but ended up with an $O(\log n)$ approximation.

We now prove the approximation in three steps (following [5]). First we show that BREG++ is competitive in those clusters out of the optimal clustering \mathcal{C}_{OPT} from which it happens to sample a center.

Lemma 2.2. *Let \mathcal{A} be an arbitrary cluster in \mathcal{C}_{OPT} and let \mathcal{C} be the clustering with just one center that was chosen with probability proportional to the weight of points in \mathcal{A} . Then, if $J(\mathcal{A})$ is the contribution of points in \mathcal{A} to the final objective, we have*

$$E[J(\mathcal{A})] \leq \left(1 + \frac{\sigma_2}{\sigma_1}\right) J_{\text{OPT}}(\mathcal{A}).$$

Proof. Let $\mu_{\mathcal{A}}$ denote the weighted mean of cluster \mathcal{A} . Since \mathcal{C}_{OPT} is optimal, it must be using $\mu_{\mathcal{A}}$ as its center. Let $W_{\mathcal{A}} = \sum_{x_i \in \mathcal{A}} w_i$. Now invoking Lemma 2.1 we note that $E[J(\mathcal{A})]$ is given by

$$\begin{aligned} & \sum_{\mu_0 \in \mathcal{A}} \frac{w_o}{W_{\mathcal{A}}} \left(\sum_{\mathbf{x}_i \in \mathcal{A}} w_i B_f(\mathbf{x}_i; \mu_0) \right) \\ &= \sum_{\mu_0 \in \mathcal{A}} \frac{w_o}{W_{\mathcal{A}}} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i B_f(\mathbf{x}_i; \mu_{\mathcal{A}}) + \sum_{\mu_0} W_{\mathcal{A}} B_f(\mu_{\mathcal{A}}; \mu_0) \\ &= J_{\text{OPT}}(\mathcal{A}) + \sum_{\mu_0 \in \mathcal{A}} W_{\mathcal{A}} \frac{B_f(\mu_{\mathcal{A}}; \mu_0)}{B_f(\mu_0; \mu_{\mathcal{A}})} B_f(\mu_0; \mu_{\mathcal{A}}) \\ &\leq J_{\text{OPT}}(\mathcal{A}) + \frac{\sigma_2}{\sigma_1} J_{\text{OPT}}(\mathcal{A}) = \left(1 + \frac{\sigma_2}{\sigma_1}\right) J_{\text{OPT}}(\mathcal{A}), \end{aligned}$$

where the last inequality follows from (2.7). □

The second step consists of showing how the algorithm behaves for the remaining centers that are chosen with the weighted farthest-first sampling.

Lemma 2.3. *Let \mathcal{A} be an arbitrary cluster in \mathcal{C}_{OPT} , and let \mathcal{C} be an arbitrary clustering. If we add a random point of \mathcal{A} as a center to \mathcal{C} using (2.4) and (2.5), then*

$$E[J(\mathcal{A})] \leq 4 \frac{\sigma_2}{\sigma_1} \left(1 + \frac{\sigma_2}{\sigma_1}\right) J_{\text{OPT}}(\mathcal{A}).$$

Proof. After choosing a center \mathbf{x}_0 from \mathcal{A} , any point $\mathbf{x}_i \in \mathcal{A}$ will contribute $w_i \min(D(\mathbf{x}_i), B_f(\mathbf{x}_i, \mathbf{x}_0))$ to the objective. Since we sample according to (2.5), the expected value of the objective $E[J(\mathcal{A})]$ is

$$\sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0 D(\mathbf{x}_0)}{\sum_{\mathbf{x} \in \mathcal{A}} w_{\mathbf{x}} D(\mathbf{x})} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \min(D(\mathbf{x}_i), B_f(\mathbf{x}_i, \mathbf{x}_0)).$$

Let \mathbf{c}_0 and \mathbf{c}_i be the centers closest to \mathbf{x}_0 and \mathbf{x}_i , respectively, From the triangle inequality we have

$$\|\mathbf{x}_0 - \mathbf{c}_i\| \leq \|\mathbf{x}_0 - \mathbf{x}_i\| + \|\mathbf{x}_i - \mathbf{c}_i\|.$$

Then using (2.7) we can bound the divergence

$$\begin{aligned} B_f(\mathbf{x}_0, \mathbf{c}_i) &\leq \sigma_2 \|\mathbf{x}_0 - \mathbf{c}_i\|^2 \leq \sigma_2 (\|\mathbf{x}_0 - \mathbf{x}_i\| + \|\mathbf{x}_i - \mathbf{c}_i\|)^2 \\ &\leq 2\sigma_2 \|\mathbf{x}_0 - \mathbf{x}_i\|^2 + 2\sigma_2 \|\mathbf{x}_i - \mathbf{c}_i\|^2 \\ &\leq 2\frac{\sigma_2}{\sigma_1} B_f(\mathbf{x}_i, \mathbf{x}_0) + 2\frac{\sigma_2}{\sigma_1} B_f(\mathbf{x}_i, \mathbf{c}_i). \end{aligned}$$

Noting that $D(\mathbf{x}_0) = B_f(\mathbf{x}_0, \mathbf{c}_0) \leq B_f(\mathbf{x}_0, \mathbf{c}_i)$ and $D(\mathbf{x}_i) = B_f(\mathbf{x}_i, \mathbf{c}_i)$, we have the bound

$$D(\mathbf{x}_0) \leq 2\frac{\sigma_2}{\sigma_1} B_f(\mathbf{x}_i, \mathbf{x}_0) + 2\frac{\sigma_2}{\sigma_1} D(\mathbf{x}_i).$$

Multiplying both sides by w_i and summing over all $\mathbf{x}_i \in \mathcal{A}$, we have (for $W_{\mathcal{A}} = \sum_{\mathbf{x}_i \in \mathcal{A}} w_i$)

$$\begin{aligned} W_{\mathcal{A}} D(\mathbf{x}_0) &\leq 2\frac{\sigma_2}{\sigma_1} \left(\sum_{\mathbf{x}_i \in \mathcal{A}} w_i B_f(\mathbf{x}_i, \mathbf{x}_0) + w_i D(\mathbf{x}_i) \right), \text{ i.e.,} \\ w_0 D(\mathbf{x}_0) &\leq 2\frac{\sigma_2}{\sigma_1} \frac{w_0}{W_{\mathcal{A}}} \left(\sum_{\mathbf{x}_i \in \mathcal{A}} w_i B_f(\mathbf{x}_i, \mathbf{x}_0) + w_i D(\mathbf{x}_i) \right). \end{aligned}$$

Now letting $R(\mathcal{A}) = \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \min(D(\mathbf{x}_i), B_f(\mathbf{x}_i, \mathbf{x}_0))$, we see that $E[J(\mathcal{A})]$ is upper bounded by

$$\begin{aligned} &2\frac{\sigma_2}{\sigma_1} \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \left(\frac{\sum_{\mathbf{x}_i \in \mathcal{A}} w_i (B_f(\mathbf{x}_i, \mathbf{x}_0) + D(\mathbf{x}_i))}{\sum_{\mathbf{x} \in \mathcal{A}} w_{\mathbf{x}} D(\mathbf{x})} R(\mathcal{A}) \right) \\ &\leq 2\frac{\sigma_2}{\sigma_1} \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i B_f(\mathbf{x}_i, \mathbf{x}_0) \\ &+ 2\frac{\sigma_2}{\sigma_1} \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i B_f(\mathbf{x}_i, \mathbf{x}_0) \\ &= 4\frac{\sigma_2}{\sigma_1} \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i B_f(\mathbf{x}_i, \mathbf{x}_0) \\ &\leq 4\frac{\sigma_2}{\sigma_1} \left(1 + \frac{\sigma_2}{\sigma_1}\right) J_{\text{OPT}}(\mathcal{A}), \end{aligned}$$

where in the second line we simplified $R(\mathcal{A})$ by using $\min(D(\mathbf{x}_i), B_f(\mathbf{x}_i, \mathbf{x}_0)) \leq D(\mathbf{x}_i)$, while for the third line we used $\min(D(\mathbf{x}_i), B_f(\mathbf{x}_i, \mathbf{x}_0)) \leq B_f(\mathbf{x}_i, \mathbf{x}_0)$. The last inequality follows from Lemma 2.2. \square

Remark 2.3 (K-means). For $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{x}$ we have $\sigma_1 = \sigma_2$, and Bregman k-means is reduces to Euclidean k-means, and our analysis reduces to that of [5].

Remark 2.4 (Mahalanobis). For $f = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x}$, where \mathbf{A} is a positive-definite matrix, the resulting Bregman divergence is a Mahalanobis distance. Here, one has $\sigma_1 = \lambda_{\min}(\mathbf{A})$, and $\sigma_2 = \lambda_{\max}(\mathbf{A})$, independent of the input data (the λ s denote eigenvalues of \mathbf{A}). The approximation ratio depends naturally on the condition number $\lambda_{\max}(\mathbf{A})/\lambda_{\min}(\mathbf{A})$ of \mathbf{A} .

Given Lemmas 2.2 and 2.3, the third step of our proof is simple as we can essentially invoke Lemma 3.3 of [5] to show that the total error incurred via the weighted farthest-first sampling is within a factor $O(\log K)$ of the optimal.

Lemma 2.4. *Let \mathcal{C} be an arbitrary clustering. Choose $u > 0$ “uncovered”¹ clusters from \mathcal{C}_{OPT} , and let \mathcal{X}_u denote the set of points in these clusters. Also let $\mathcal{X}_c = X \setminus \mathcal{X}_u$. Now suppose we add $t \leq u$ random centers to \mathcal{C} , chosen with the weighted farthest-first sampling. Let \mathcal{C}' denote the resulting clustering, and let J' denote the corresponding objective. Then, $E[J']$ is at most*

$$\left(J(\mathcal{X}_c) + 4 \frac{\sigma_2}{\sigma_1} \left(1 + \frac{\sigma_2}{\sigma_1} \right) J_{OPT}(\mathcal{X}_u) \right) \cdot (1 + H_t) + \frac{u-t}{u} \cdot J(\mathcal{X}_u),$$

where H_t denotes the Harmonic number $1 + \frac{1}{2} + \dots + \frac{1}{t}$.

Proof. Direct from the proof of Lemma 3.3 of [5]. □

Finally, we have the main approximation theorem.

Theorem 2.5. *A clustering \mathcal{C} obtained via BREG++ satisfies*

$$E[J(\mathcal{C})] \leq 4 \frac{\sigma_2}{\sigma_1} \left(1 + \frac{\sigma_2}{\sigma_1} \right) (\log K + 2) J_{OPT}.$$

Proof. Immediate from Theorem 3.1 of [5] by a direct application of Lemma 2.4. □

3 Implications of BREG++

We now describe some important implications of our BREG++ method derived above.

3.1 Information Theoretic Clustering.

With $f(\mathbf{x}) = \sum_j x_j \log x_j$, the Bregman divergence B_f becomes the (un-normalized) Kullback-Leibler divergence, and Bregman k-means reduces to information theoretic clustering (ITC). Even though local and greedy methods for ITC have been well studied [6, 19, 37], approximation algorithms for it have only been developed very recently [2, 14] (both papers are from 2008).

For ITC, our assumptions on bounded σ_1 and σ_2 are equivalent to those of [2] as mentioned previously. Under these assumptions we obtain an efficient k-means type $O(\log K)$ approximation algorithm, while Ackermann et al. [2] obtain an $O(1 + \epsilon)$ approximation algorithm, with an impractical running time of $O(dn2^{(\frac{K}{\epsilon})^{O(1)}})$. The even more recent paper of Ackermann and Blömer [1] yields an ITC algorithm faster than that of [2], but still has an impractical running time of $O(dKn + d^2 2^{(K/\epsilon)^{\Theta(1)}} \log^{K+2} n)$.

Chaudhuri and McGregor [14] develop an approximation algorithm ITC that does not make any assumptions on the data. They first lower bound the KL-divergence using a Hellinger distance, then cluster approximately using this distance, before recovering clusters for KL. Their method is clever, but leads to a non-constant approximation ratio of $O(\log n)$ (n is the number of input points). Obtaining an $O(\log K)$ algorithm for ITC without any assumptions on the data therefore remains an open problem—though we suspect that without additional assumptions, ITC might be inapproximable to better than a polylog factor.

Details. Observe that given the definition the Bregman divergence (2.1), using Taylor series expansion of f around \mathbf{x} we immediately have

$$B_f(\mathbf{x}, \mathbf{y}) = \frac{1}{2} (\mathbf{x} - \mathbf{y})^T \nabla^2 f(\xi_{\mathbf{x}, \mathbf{y}}) (\mathbf{x} - \mathbf{y}),$$

where $\xi_{\mathbf{x}, \mathbf{y}}$ is some point between \mathbf{x} and \mathbf{y} . Since $\nabla^2 f$ is positive definite, the constants σ_1 and σ_2 can be obtained by bounding the minimum and maximum eigenvalues of $\nabla^2 f$. For ITC, one assumes the input data to be normalized, i.e., for each $\mathbf{x} \in \mathcal{X}$, $\sum_j x_j = 1$. Since the Hessian for the KL-divergence is $\nabla^2 f(\mathbf{x}) = \text{Diag}(x_1^{-1}, \dots, x_d^{-1})$ at a point $\mathbf{x} = (x_1, \dots, x_d)$, its maximum eigenvalue is $(1 - \|\mathbf{x} - \mathbf{x}\|_\infty)^{-1}$ and the minimum eigenvalue is $\|\mathbf{x}\|_\infty^{-1}$.

¹An “uncovered” cluster is one from which a center has not been chosen by the weighted farthest-first sampling procedure.

Since the interpolating point $\xi_{\mathbf{x}, \mathbf{y}}$ lies in $\text{conv}(\mathcal{X})$ for $\mathbf{x}, \mathbf{y} \in \text{conv}(\mathcal{X})$, we can select

$$\sigma_1^{-1} = \max_{\mathbf{y} \in \text{conv}(\mathcal{X})} \|\mathbf{y}\|_\infty = \max_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_\infty.$$

Analogously we have,

$$\sigma_2^{-1} = \min_{\mathbf{y} \in \text{conv}(\mathcal{X})} (1 - \|\mathbf{1} - \mathbf{y}\|_\infty) = \min_{\mathbf{x} \in \mathcal{X}} (1 - \|\mathbf{1} - \mathbf{x}\|_\infty).$$

The reduction from $\text{conv}(\mathcal{X})$ to \mathcal{X} results from each $\mathbf{y} \in \text{conv}(\mathcal{X})$ being a convex combination of the data points, whereby its coordinates are a convex combination of the data point coordinates. This convex combination is maximized by putting all weight on the maximum component. Thus, $\sigma_1 \geq 1$ and σ_2 corresponds to inverse of the minimum coordinate entry γ in the data set, so $\sigma_2/\sigma_1 \leq \gamma^{-1}$. Ackermann et al. [2] and Ackermann and Blömer [1] will also have a similar dependence on γ .

3.2 Mixture Modeling on Exponential Families

The parametric mixture modeling problem entails fitting a mixture of K distributions from a pre-defined family to a set of observations. Let \mathbf{x} denote an observation, $\boldsymbol{\pi}$ a prior over the mixture components, and $\boldsymbol{\theta}_h$ the parameters corresponding to the z^{th} mixture component. Then, a mixture model assumes the following generative process: (i) sample $z \sim \boldsymbol{\pi}$, and (ii) sample $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta}_z)$. Given a set of observations, the basic mixture modeling problem is that of finding the parameters $\Theta = (\boldsymbol{\pi}, \boldsymbol{\theta}_z, \{z\}_1^K)$ such that $\log p(\mathbf{x}|\Theta)$ is maximized.² More formally, if X denotes the random variable corresponding to the observations and Z denotes the one corresponding to the mixture components, a direct calculation [7] shows that the problem is equivalent to maximizing

$$J_{MM}(\Theta) = E_{Z|X} [\log p(X, Z|\Theta)] + H(Z|X) \quad (3.1)$$

over Θ , where

$$p(z|\mathbf{x}) = \frac{\pi_z p(\mathbf{x}|\boldsymbol{\theta}_z)}{p(\mathbf{x})},$$

and $H(\cdot)$ denotes the Shannon entropy of $Z|X$. For the purposes of analysis one can focus on the expected log-likelihood of the data, i.e., the first term in (3.1). In practice, for several real datasets, the distribution $p(z|\mathbf{x})$ is typically skewed in that it has a high value ≈ 1 for some z^* , and low values ≈ 0 for other z , so that the entropy $H(Z|X)$ is rather small. As a result, ignoring the entropy term may be reasonable in an application. A more theoretically well motivated justification can be given by considering “hard clustering” for the mixture modeling problem, where we focus on the family of posterior distributions

$$q(z|\mathbf{x}) = \begin{cases} 1, & \text{if } p(z|\mathbf{x}) > p(z'|\mathbf{x}), \forall z' \neq z, \\ 0, & \text{otherwise.} \end{cases}$$

For simplicity, we let $z_i^* = z$ such that $p(z|\mathbf{x}_i) > p(z'|\mathbf{x}_i), \forall z' \neq z$. If $J_Q(\Theta)$ is the corresponding objective, then following [7] we have

$$J_{MM}(\Theta) - H(Z|X) \leq J_Q(\Theta) \leq J_{MM}(\Theta).$$

Thus, $J_Q(\Theta)$ forms a tight lower bound to the original objective $J_{MM}(\Theta)$, especially when entropy of $Z|X$ is small, which is true for several real world problems.

²Note that the objective $\log p(\mathbf{x}|\Theta)$ is always non-positive since $p(\mathbf{x}|\Theta) \leq 1$. All objective functions in this section share the same property.

We focus on mixture models over exponential family distributions, whose density functions can be written as

$$p(\mathbf{x}|\theta) = \exp(\langle \mathbf{x}, \theta \rangle - \psi(\theta_z))p_0(\mathbf{x}),$$

where ψ is a convex function of Legendre [39] type known as the cumulant, θ is the natural parameter, and $p_0(\mathbf{x})$ is a base measure. A particular choice of ψ determines a family, such as Gaussian or Poisson, while a particular choice of θ determines a specific distribution in the family. The expectation parameter $\boldsymbol{\mu} = E[X]$ of an exponential family distribution is uniquely tied to the natural parameter through a Legendre transform $\boldsymbol{\mu} = \nabla\psi(\theta)$, and $\theta = \nabla\phi(\boldsymbol{\mu})$ where ϕ is the Legendre conjugate of ψ [39].

Our results below rely on the following key connection between exponential family distributions and Bregman divergences [8]. The density function $p(\mathbf{x}|\theta_z)$ of an exponential family distribution can be uniquely written as

$$p(\mathbf{x}|\theta_z) = \exp(-B_\phi(\mathbf{x}, \boldsymbol{\mu}))f_0(\mathbf{x}), \quad (3.2)$$

where ϕ is the Legendre conjugate of ψ , and $\boldsymbol{\mu} = E[X] = \nabla\psi(\theta)$ is the expectation parameter.

Now we use BREG++ to optimize $J_Q(\Theta)$ based on the expectation parameter $\boldsymbol{\mu}_z$ of each component $z = 1, \dots, K$. The only additional step is to set (for each i) $z_i^* = z$ if \mathbf{x}_i is assigned to cluster z .

Lemma 3.1. *Let Θ_{MM} denote the natural parameters corresponding to the final mean parameters after convergence, and let $\pi_z = |C_z|/n$, where $|C_z|$ denotes the number of elements in the z -th cluster. Then, (recall J_Q is negative)*

$$E[J_Q(\Theta_{MM})] \geq 4\frac{\sigma_2}{\sigma_1} \left(1 + \frac{\sigma_2}{\sigma_1}\right) J_Q(\Theta^*),$$

where Θ^* denotes an optimum set of parameters.

Proof. By definition,

$$\begin{aligned} \max_{\Theta} J_Q(\Theta) &= \max_{\Theta} E_{Z|X \sim Q}[\log p(X, Z|\Theta)] + H_Q(Z|X) \\ &= \max_{\Theta} \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{x}_i|\theta_{z_i^*}) \\ &= \max_{\boldsymbol{\mu}} -\frac{1}{n} \sum_{i=1}^n B_\phi(\mathbf{x}_i, \boldsymbol{\mu}_{z_i^*}) = \min_{\boldsymbol{\mu}} \frac{1}{n} \sum_{i=1}^n B_\phi(\mathbf{x}_i, \boldsymbol{\mu}_{z_i^*}), \end{aligned}$$

which is precisely the objective function for Bregman k-means. Thus, the result follows from Lemma 2.3, and a change in the direction of the inequality due to conversion of the minimization problem to a maximization problem by multiplying both sides with -1. \square

4 Weighted Kernel K-means

In this section we present WKKM++, an $O(\log K)$ approximation algorithm for the weighted kernel k-means (WKKM) problem [20].

Let $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ denote the set of input data points (which may or may not be available explicitly), and let $\phi : \mathbf{x} \rightarrow \phi(\mathbf{x}) \in \mathcal{H}$, denote the feature map that takes \mathbf{x} to its corresponding point in an RKHS \mathcal{H} . Further, let w_1, w_2, \dots, w_n denote non-negative weights corresponding to each input point.

WKKM seeks a clustering $\mathcal{C} = \{C_1, C_2, \dots, C_K\}$ such that the following objective is minimized

$$J(\mathcal{C}) = \sum_{h=1}^K \sum_{\mathbf{x}_i \in C_h} w_i \|\phi(\mathbf{x}_i) - \boldsymbol{\mu}_h\|^2, \quad (4.1)$$

where $\boldsymbol{\mu}_h$ is the weighted mean of cluster \mathcal{C}_h , i.e.,

$$\boldsymbol{\mu}_h = \frac{\sum_{\mathbf{x}_i \in \mathcal{C}_h} w_i \phi(\mathbf{x}_i)}{\sum_{\mathbf{x}_i \in \mathcal{C}_h} w_i}. \quad (4.2)$$

For a given clustering, the weighted centroid $\boldsymbol{\mu}_h$ (4.2) is optimal; formally stated

Lemma 4.1 (Optimality of Mean). *Let \mathcal{A} be a set of points with weighted mean $\boldsymbol{\mu}_{\mathcal{A}}$, and let $\phi(\mathbf{z})$ be an arbitrary point. Then,*

$$\begin{aligned} & \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \phi(\mathbf{z})\|^2 \\ &= \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \boldsymbol{\mu}_{\mathcal{A}}\|^2 + W \|\boldsymbol{\mu}_{\mathcal{A}} - \phi(\mathbf{z})\|^2, \end{aligned}$$

where $W = \sum_{\mathbf{x}_i \in \mathcal{A}} w_i$.

Proof. Elementary; similar to Lemma 2.2. □

The WKKM++ algorithm proceeds exactly like the BREG++ algorithm of Section 2. Specifically, WKKM++ selects K initial means from amongst the data points using a particular sampling procedure. These K means are then used as an initialization for the WKKM algorithm:

1. For each $i = 1..N$ assign point \mathbf{x}_i to its nearest mean, update corresponding cluster \mathcal{C}_h
2. For each $h = 1..K$ update $\boldsymbol{\mu}_h$ using (4.2).
3. Repeat steps 1 and 2 until convergence.

This standard approach is guaranteed to only monotonically decrease the objective function. However, the crux of the analysis is in showing that after just the weighted farthest-first initialization based on (4.3) and (4.4), WKKM++ comes to within a factor $O(\log K)$ of the optimal. We elaborate on this below.

The first mean $\boldsymbol{\mu}_1 = \phi(\mathbf{x}_i)$ is chosen uniformly at random from the data points. The remaining means are chosen with a weighted farthest-first sampling procedure outlined below.

First, we define the weighting function

$$D(\mathbf{x}) = \min_{\boldsymbol{\mu} \in \mathcal{C}} \|\phi(\mathbf{x}) - \boldsymbol{\mu}\|^2, \quad (4.3)$$

which is easily computable using dot-products only because

$$\|\phi(\mathbf{x}) - \boldsymbol{\mu}\|^2 = \phi(\mathbf{x})^T \phi(\mathbf{x}) + \boldsymbol{\mu}^T \boldsymbol{\mu} - 2\phi(\mathbf{x})^T \boldsymbol{\mu},$$

and $\boldsymbol{\mu}$ is just one of the points \mathbf{x}_i during the initialization.

At a given stage in the algorithm, suppose we wish to select the next mean $\boldsymbol{\mu}_h$. The probability that a given point $\phi(\mathbf{x}_i)$ is chosen to be $\boldsymbol{\mu}_h$ is set to

$$P(\boldsymbol{\mu}_h = \phi(\mathbf{x}_i)) = \frac{w_i D(\mathbf{x}_i)^2}{\sum_{j=1}^n w_j D(\mathbf{x}_j)^2}. \quad (4.4)$$

The probability (4.4) is also computable using dot-products only, as it involves only the weighting function $D(\mathbf{x})$, which itself is so computable. We repeatedly select means using (4.3) and (4.4) until we have selected K different means.

Now we proceed to show that the weighted farthest-first initialization as described above is sufficient to guarantee an $O(\log K)$ factor. First we show that WKKM++ is competitive in those clusters out of the optimal clustering \mathcal{C}_{OPT} from which it happens to sample a center.

Lemma 4.2 (First mean). *Let \mathcal{A} be an arbitrary cluster in \mathcal{C}_{OPT} , and let \mathcal{C} be the clustering with just one center, which is chosen with probability proportional to the weight of points in \mathcal{A} . Then, if $J(\mathcal{A})$ is the contribution of points in \mathcal{A} to the final objective, $E[J(\mathcal{A})] \leq 2J_{OPT}(\mathcal{A})$.*

Proof. Let $\mu_{\mathcal{A}}$ denote the weighted mean of cluster \mathcal{A} . Since \mathcal{C}_{OPT} is optimal, it must be using $\mu_{\mathcal{A}}$ as its center. Let $W_{\mathcal{A}} = \sum_{\mathbf{x}_i \in \mathcal{A}} w_i$. With the random initialization, assuming all points in \mathcal{A} stay assigned to the cluster \mathcal{A} till the end, using Lemma 4.1 we see that $E[J(\mathcal{A})]$ is given by

$$\begin{aligned} & \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \left(\sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|^2 \right) \\ &= \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \left(\sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \mu_{\mathcal{A}}\|^2 + W_{\mathcal{A}} \|\phi(\mathbf{x}_0) - \mu_{\mathcal{A}}\|^2 \right) \\ &= \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \mu_{\mathcal{A}}\|^2 + \sum_{\mathbf{x}_0 \in \mathcal{A}} w_0 \|\phi(\mathbf{x}_0) - \mu_{\mathcal{A}}\|^2 \\ &= 2J_{OPT}(\mathcal{A}). \end{aligned}$$

Since the contribution of each point can only decrease in subsequent WKKM++ iterations, we have $E[J(\mathcal{A})] \leq 2J_{OPT}(\mathcal{A})$. \square

Next we show how WKKM++ behaves for the remaining centers that it picks.

Lemma 4.3 (Other means). *Let \mathcal{A} be an arbitrary cluster in \mathcal{C}_{OPT} , and let \mathcal{C} be an arbitrary clustering. If we add a random point from \mathcal{A} as a center to \mathcal{C} using the farthest-first sampling, then $E[J(\mathcal{A})] \leq 8J_{OPT}(\mathcal{A})$.*

Proof. After choosing a center $\phi(\mathbf{x}_0)$, any point $\mathbf{x}_i \in \mathcal{A}$ will contribute $w_i \min(D(\mathbf{x}_i), \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|)^2$ to the objective. Since subsequent assignments can only decrease the contribution, we have

$$E[J(\mathcal{A})] \leq \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0 D(\mathbf{x}_0)^2}{\sum_{\mathbf{x}' \in \mathcal{A}} w_{\mathbf{x}'} D(\mathbf{x}')^2} \times \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \min(D(\mathbf{x}_i), \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|)^2.$$

Let $\mathbf{c}_0, \mathbf{c}_i$ be the closest centers to $\phi(\mathbf{x}_0), \phi(\mathbf{x}_i)$, respectively. Then, from the triangle inequality we have

$$\|\phi(\mathbf{x}_0) - \mathbf{c}_i\| \leq \|\phi(\mathbf{x}_0) - \phi(\mathbf{x}_i)\| + \|\phi(\mathbf{x}_i) - \mathbf{c}_i\|.$$

Further, from the Cauchy-Schwartz inequality we have

$$\|\phi(\mathbf{x}_0) - \mathbf{c}_i\|^2 \leq 2\|\phi(\mathbf{x}_0) - \phi(\mathbf{x}_i)\|^2 + 2\|\phi(\mathbf{x}_i) - \mathbf{c}_i\|^2.$$

Noting that $D(\mathbf{x}_0) = \|\phi(\mathbf{x}_0) - \mathbf{c}_0\| \leq \|\phi(\mathbf{x}_0) - \mathbf{c}_i\|$, we hence have

$$D(\mathbf{x}_0)^2 \leq 2\|\phi(\mathbf{x}_0) - \phi(\mathbf{x}_i)\|^2 + 2D(\mathbf{x}_i)^2.$$

Multiplying both sides by w_i and summing over $\mathbf{x}_i \in \mathcal{A}$ we obtain

$$\begin{aligned} W_{\mathcal{A}} D(\mathbf{x}_0)^2 &\leq 2 \left(\sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|^2 + w_i D(\mathbf{x}_i)^2 \right) \\ \Rightarrow w_0 D(\mathbf{x}_0)^2 &\leq 2 \frac{w_0}{W_{\mathcal{A}}} \left(\sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|^2 + w_i D(\mathbf{x}_i)^2 \right). \end{aligned}$$

Hence, $E[J(\mathcal{A})]$ is upper bounded by

$$\begin{aligned}
& 2 \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \left(\frac{\sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|^2}{\sum_{\mathbf{x}' \in \mathcal{A}} w_{\mathbf{x}'} D(\mathbf{x}')^2} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \min(D(\mathbf{x}_i), \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|)^2 \right) \\
& + 2 \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \left(\frac{\sum_{\mathbf{x}_i \in \mathcal{A}} w_i D(\mathbf{x}_i)^2}{\sum_{\mathbf{x}' \in \mathcal{A}} w_{\mathbf{x}'} D(\mathbf{x}')^2} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \min(D(\mathbf{x}_i), \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|)^2 \right) \\
& \leq 2 \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|^2 + 2 \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|^2 \\
& = 4 \sum_{\mathbf{x}_0 \in \mathcal{A}} \frac{w_0}{W_{\mathcal{A}}} \sum_{\mathbf{x}_i \in \mathcal{A}} w_i \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|^2 \\
& \leq 8J_{\text{OPT}}(\mathcal{A}),
\end{aligned}$$

where in the first line we used $\min(D(\mathbf{x}_i), \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|)^2 \leq D(\mathbf{x}_i)^2$ in the first expression, and $\min(D(\mathbf{x}_i), \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|)^2 \leq \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_0)\|^2$ in the second expression, and the last line follows from Lemma 4.2. \square

Theorem 4.4 (Approximation Ratio). *A clustering \mathcal{C} obtained via WKKM++ satisfies*

$$E[J(\mathcal{C})] \leq 8(\log K + 2)J_{\text{OPT}}.$$

Proof. Exactly follows proof structure in Section 2. \square

Remark 4.1 (Implications). WKKM++ has two important implications. First, by exploiting the equivalence between several graph-cut criteria and WKKM [20], one can hope to obtain better graph-cuts using WKKM++. A formal proof of this observation however remains an open problem. Second, the connection of WKKM to semi-supervised graph-clustering [30] leads to a potentially improved algorithm for semi-supervised clustering.

5 Approximation Algorithms for Tensor Clustering and Co-clustering

In its simplest formulation, co-clustering refers to the simultaneous partitioning of the rows and columns of the input data matrix into $K \times L$ co-clusters (sub-matrices). Co-clustering, also called bi-clustering [15, 27] has witnessed increasing interest over the years (see [10] and references therein). Anagnostopoulos et al. [4] seem to be the first to present an approximation algorithm for co-clustering based on a minimum-sum squared residue criterion of [16]. They (i.e., [4]) posed two open questions:

- Could one extend their ideas to obtain approximation algorithms for Bregman co-clustering?
- Could one design approximation algorithms for 3-way co-clustering of a tensor in $\mathbb{R}^{n_1 \times n_2 \times n_3}$?

Below we answer both these questions in the affirmative, leading to the first (to our knowledge) approximation algorithms for Bregman matrix and tensor co-clustering. In fact, our results hold for arbitrary m -way tensor co-clustering, not just the 3-way case.

We directly develop an approximation algorithm for m -way tensor co-clustering (hereafter *clustering*) that yields an approximation algorithm for Bregman co-clustering, which is nothing but tensor clustering for $m = 2$. Specifically, we show a competitiveness of $O(m \log K)$ for m -way Bregman co-clustering.

Tensors are well-studied in multilinear algebra [24], but they are not so widespread in the machine learning community. Therefore, to facilitate an easier understanding of our proofs, we briefly summarize some important tensor notation below for those unfamiliar to it.

5.1 Background on Tensors

Most of the material in this section is taken from the well-written paper of de Silva and Lim [18], whose notation turns out to be particularly suitable for our analysis. An order- m tensor \mathbf{A} may be viewed as an element of the vector space $\mathbb{R}^{n_1 \times \dots \times n_m}$. A particular component of the tensor \mathbf{A} is represented by the multiply-indexed value $a_{i_1 i_2 \dots i_m}$, where $i_j = 1 \dots n_j$ for $1 \leq j \leq m$.

Multilinear matrix multiplication. The most important operation that we do with tensors is that of multilinear matrix multiplication, which is a generalization of the familiar concept of matrix multiplication. Matrices *act* on other matrices by either left or right multiplications. For an order-3 tensor, there are three dimensions along which a matrix may act via matrix multiplication. For example, given an order-3 tensor $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, and three matrices $\mathbf{P} \in \mathbb{R}^{p_1 \times n_1}$, $\mathbf{Q} \in \mathbb{R}^{p_2 \times n_2}$, and $\mathbf{R} \in \mathbb{R}^{p_3 \times n_3}$, *multilinear matrix multiplication* is the operation defined by the action of these three matrices on the different dimensions of \mathbf{A} that yields the tensor $\mathbf{A}' \in \mathbb{R}^{p_1 \times p_2 \times p_3}$. Formally, the entries of the tensor \mathbf{A}' are given by

$$a'_{lmn} = \sum_{i,j,k=1}^{n_1, n_2, n_3} p_{li} q_{mj} r_{nk} a_{ijk}, \quad (5.1)$$

and this operation is written compactly as

$$\mathbf{A}' = (\mathbf{P}, \mathbf{Q}, \mathbf{R}) \cdot \mathbf{A}. \quad (5.2)$$

The notation (5.2) is particularly nice and may be viewed as the *group action* of $(\mathbf{P}, \mathbf{Q}, \mathbf{R})$ on \mathbf{A} (group-action refers to the situation when a group with a particular algebraic structure “acts” on another set; for the multilinear multiplication notation one can view it as the set $G = \mathbb{R}^{p_1 \times n_1} \times \mathbb{R}^{p_2 \times n_2} \times \mathbb{R}^{p_3 \times n_3}$ acting on the set $X = \mathbb{R}^{n_1 \times n_2 \times n_3}$). Addition in G is defined entry-wise:

$$(\mathbf{P}_1, \mathbf{Q}_1, \mathbf{R}_1) + (\mathbf{P}_2, \mathbf{Q}_2, \mathbf{R}_2) = (\mathbf{P}_1 + \mathbf{P}_2, \mathbf{Q}_1 + \mathbf{Q}_2, \mathbf{R}_1 + \mathbf{R}_2).$$

Multilinear multiplication extends naturally to tensors of arbitrary (finite) order. If $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_m}$, and $\mathbf{P}_1 \in \mathbb{R}^{p_1 \times n_1}, \dots, \mathbf{P}_m \in \mathbb{R}^{p_m \times n_m}$, then $\mathbf{A}' = (\mathbf{P}_1, \dots, \mathbf{P}_m) \cdot \mathbf{A}$ has entries

$$a'_{i_1 i_2 \dots i_m} = \sum_{j_1, \dots, j_m=1}^{n_1, \dots, n_m} p_{i_1 j_1}^{(1)} \cdots p_{i_m j_m}^{(m)} a_{j_1 \dots j_m}, \quad (5.3)$$

where $p_{ij}^{(k)}$ denotes the ij -entry of matrix \mathbf{P}_k .

Example 5.1 (Matrix Multiplication). *Let $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$, $\mathbf{P} \in \mathbb{R}^{p \times n_1}$, and $\mathbf{Q} \in \mathbb{R}^{q \times n_2}$ be given. The matrix product $\mathbf{P}\mathbf{A}\mathbf{Q}^T$ can be written as the multilinear multiplication $(\mathbf{P}, \mathbf{Q}) \cdot \mathbf{A}$.*

Example 5.2 (Basic Properties). *The following properties of multilinear multiplication are easily verified (and generalized to tensors of arbitrary order):*

1. **Linearity:** *Let $\alpha, \beta \in \mathbb{R}$, and \mathbf{A} and \mathbf{B} be tensors with same dimensions, then*

$$(\mathbf{P}, \mathbf{Q}) \cdot (\alpha\mathbf{A} + \beta\mathbf{B}) = \alpha(\mathbf{P}, \mathbf{Q}) \cdot \mathbf{A} + \beta(\mathbf{P}, \mathbf{Q}) \cdot \mathbf{B}$$

2. **Product rule:** *For matrices $\mathbf{P}_1, \mathbf{P}_2, \mathbf{Q}_1, \mathbf{Q}_2$ of appropriate dimensions, and a tensor \mathbf{A} we have*

$$(\mathbf{P}_1, \mathbf{P}_2) \cdot ((\mathbf{Q}_1, \mathbf{Q}_2) \cdot \mathbf{A}) = (\mathbf{P}_1\mathbf{Q}_1, \mathbf{P}_2\mathbf{Q}_2) \cdot \mathbf{A}$$

3. **Multilinearity:** *Let $\alpha, \beta \in \mathbb{R}$, and \mathbf{P}, \mathbf{Q} , and \mathbf{R} be matrices of appropriate dimensions. Then, for a tensor \mathbf{A} the following holds*

$$(\mathbf{P}, \alpha\mathbf{Q} + \beta\mathbf{R}) \cdot \mathbf{A} = \alpha(\mathbf{P}, \mathbf{Q}) \cdot \mathbf{A} + \beta(\mathbf{P}, \mathbf{R}) \cdot \mathbf{A}$$

Inner Product: The Frobenius norm induces an inner-product that can be defined as

$$\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i_1, \dots, i_m} a_{i_1 \dots i_m} b_{i_1 \dots i_m}, \quad (5.4)$$

so that $\|\mathbf{A}\|_F^2 = \langle \mathbf{A}, \mathbf{A} \rangle$ holds as usual. The following property of this inner product is easily verified

$$\langle (\mathbf{P}_1, \dots, \mathbf{P}_m) \cdot \mathbf{A}, (\mathbf{Q}_1, \dots, \mathbf{Q}_m) \cdot \mathbf{B} \rangle = \langle \mathbf{A}, (\mathbf{P}_1^T \mathbf{Q}_1, \dots, \mathbf{P}_m^T \mathbf{Q}_m) \cdot \mathbf{B} \rangle. \quad (5.5)$$

Proof: Using definition (5.3) along with (5.4) we have

$$\begin{aligned} \langle (\mathbf{P}_1, \dots, \mathbf{P}_m) \cdot \mathbf{A}, (\mathbf{Q}_1, \dots, \mathbf{Q}_m) \cdot \mathbf{B} \rangle &= \sum_{i_1, \dots, i_m} \sum_{\substack{j_1, \dots, j_m \\ k_1, \dots, k_m}} p_{i_1 j_1}^{(1)} q_{i_1 k_1}^{(1)} \cdots p_{i_m j_m}^{(m)} q_{i_m k_m}^{(m)} a_{j_1 \dots j_m} b_{k_1 \dots k_m}, \\ &= \sum_{\substack{j_1, \dots, j_m \\ k_1, \dots, k_m}} \left(\sum_{i_1} p_{i_1 j_1}^{(1)} q_{i_1 k_1}^{(1)} \right) \cdots \left(\sum_{i_m} p_{i_m j_m}^{(m)} q_{i_m k_m}^{(m)} \right) a_{j_1 \dots j_m} b_{k_1 \dots k_m} \\ &= \sum_{\substack{j_1, \dots, j_m \\ k_1, \dots, k_m}} (\mathbf{P}_1^T \mathbf{Q}_1)_{j_1 k_1} \cdots (\mathbf{P}_m^T \mathbf{Q}_m)_{j_m k_m} a_{j_1 \dots j_m} b_{k_1 \dots k_m} = \sum_{j_1 \dots j_m} a_{j_1 \dots j_m} b'_{j_1 \dots j_m} = \langle \mathbf{A}, \mathbf{B}' \rangle, \end{aligned}$$

where $\mathbf{B}' = (\mathbf{P}_1^T \mathbf{Q}_1, \dots, \mathbf{P}_m^T \mathbf{Q}_m) \cdot \mathbf{B}$.

5.2 Tensor clustering

Given the background above, we are now ready to formally state the Bregman tensor clustering problem.

Let $\mathbf{A} \in \mathbb{R}^{n_1 \times \dots \times n_m}$ be an order- m tensor. Tensor clustering refers to a partitioning of \mathbf{A} into sub-tensors or simply clusters, so that the entries of each cluster are as coherent as possible. The goal of (one simple version) of Bregman tensor clustering is to partition \mathbf{A} into sub-tensors so that the sum of the Bregman divergences of individual elements in the sub-tensor to their corresponding cluster representatives is minimized. The cluster representatives turn out to be simply the means of the associated sub-tensors because we are minimizing Bregman divergences.

A cluster (sub-tensor) is indexed by subsets of indices along each dimension. Let $I_j \subseteq \{1, \dots, n_j\}$ denote such an index subset for dimension j . Then the cluster representative corresponding to a sub-tensor is simply its mean, i.e.

$$M_{I_1 \dots I_m} = \frac{1}{|I_1| \cdots |I_m|} \sum_{i_1 \in I_1, \dots, i_m \in I_m} a_{i_1 \dots i_m}. \quad (5.6)$$

Assuming that each dimension j is partitioned into k_j clusters, we can collect all the different representatives (each of which can be written in the form (5.6)) into a *means tensor* $\mathbf{M} \in \mathbb{R}^{k_1 \times \dots \times k_m}$. Thus, we have a total of $\prod_j k_j$ tensor clusters. Let $\bar{\mathbf{C}}_j \in \{0, 1\}^{n_j \times k_j}$ denote the cluster indicator matrix for dimension j . In such a matrix, entry i of column k is one if and only if i is in the k th index set for tensor dimension j .

Given this notation, we can now formally state the Bregman tensor clustering problem:

$$\underset{\bar{\mathbf{C}}_1, \dots, \bar{\mathbf{C}}_m}{\text{minimize}} \quad B_f(\mathbf{A}, (\bar{\mathbf{C}}_1, \dots, \bar{\mathbf{C}}_m) \cdot \mathbf{M}), \quad \text{s.t. } \bar{\mathbf{C}}_j \in \{0, 1\}^{n_j \times k_j}. \quad (5.7)$$

Problem (5.7) can be rewritten in a more useful form. To that end, let \mathbf{C}_j be the normalized cluster indicator matrix obtained from $\bar{\mathbf{C}}_j$ by normalizing the columns to have unit-norm (i.e., $\mathbf{C}_j^T \mathbf{C}_j = \mathbf{I}_{k_j}$). Then (5.7) may be rewritten as

$$\underset{\mathbf{C}_1, \dots, \mathbf{C}_m}{\text{minimize}} \quad J(\mathbf{C}) = B_f(\mathbf{A}, (\mathbf{P}_1, \dots, \mathbf{P}_m) \cdot \mathbf{A}), \quad \text{where } \mathbf{P}_j = \mathbf{C}_j \mathbf{C}_j^T, \quad (5.8)$$

In the sequel, we will refer to a clustering by its parametrizations via both indicator and projection matrices. A little remark on the side: since $C_j^T C_j = \mathbf{I}_{k_j}$, one can relax the “hard”-clustering constraints on C_j to just orthogonality constraints. Indeed, such relaxations form the basis of spectral-relaxations for Euclidean K-means as well as co-clustering [16]. However, we will not use such relaxations in this paper.

Summary of the Algorithm: Broadly speaking, our tensor clustering approximation algorithm is based on clustering along subsets of dimensions using a guaranteed approximation algorithm, and then combining the resulting clusterings to obtain tensor clusters. A particular example of such a scheme would be to cluster along single dimensions using a method such as BREG++ clustering. Note that clustering along a single dimension in a tensor is a generalization of clustering the one-dimensional sub-tensors, i.e. vectors, in a matrix. In a tensor, we form groups of $m - 1$ -way tensors, for instance, grouping matrices in a 3-way tensor. Thanks to the separability of our Bregman divergences, BREG++ directly extends to sub-tensor objects. Taking the 3-way example with sub-matrices \mathbf{A}, \mathbf{B} , recall that $B_f(\mathbf{A}, \mathbf{B}) = B_f(\text{vec}(\mathbf{A}), \text{vec}(\mathbf{B}))$. So we simply treat the sub-tensors as vectors.

Our analysis below establishes that given a clustering algorithm that clusters along t of the dimensions at a time with an approximation factor of α_t , we can achieve an objective within $O(\lceil m/t \rceil \frac{\sigma_2}{\sigma_1} \alpha_t)$ of the optimal; the scaling factors σ_1 and σ_2 are defined as

$$\sigma_1 = \inf_{\substack{x \in \{a_{ij}\} \\ y \in \text{conv}(\{a_{ij}\})}} \frac{B_f(x, y)}{(x - y)^2}, \quad \sigma_2 = \sup_{\substack{x \in \{a_{ij}\} \\ y \in \text{conv}(\{a_{ij}\})}} \frac{B_f(x, y)}{(x - y)^2}. \quad (5.9)$$

Note: For simplicity of exposition we assume that we cluster an order- m tensor along t dimensions at a time and to eventually combine the resulting m/t sub-clusterings. Our analysis can be generalized (at the expense of laborious algebra) to the case where we cluster along *partitions* of varying sizes, say $\{t_1, \dots, t_r\}$, where $t_1 + \dots + t_r = m$.

5.2.1 Analysis

In this section we prove our tensor clustering approximation theorem, which yields as corollaries efficient approximation algorithms based on BREG++ for both Bregman co-clustering, and Bregman tensor clustering.

Theorem 5.3 (Approximation guarantee). *Let \mathbf{A} be the input order- m tensor, and let \mathcal{C}_j denote the clustering of \mathbf{A} along the j th subset of t dimensions ($1 \leq j \leq m/t$), as obtained by a multiway clustering algorithm with guarantee³ α_t . Let $\mathcal{C} = (\mathcal{C}_1, \dots, \mathcal{C}_{m/t})$ denote the induced tensor clustering. Then⁴*

$$J(\mathcal{C}) \leq 2^{\log(m/t)} \frac{\sigma_2}{\sigma_1} \alpha_t J_{OPT}(m)$$

Corollary 5.4 (Approximation with BREG++). *Let \mathbf{A} be the input order- m tensor, and let \mathcal{C}_j denote the clustering of \mathbf{A} along dimension j ($1 \leq j \leq m$), as obtained via BREG++. Let $\mathcal{C} = (\mathcal{C}_1, \dots, \mathcal{C}_m)$ denote the induced tensor clustering. Then,*

$$E[J(\mathcal{C})] \leq 4m \frac{\sigma_2^2}{\sigma_1^2} \left(1 + \frac{\sigma_2}{\sigma_1}\right) (\log K^* + 2) J_{OPT}(m),$$

where $K^* = \max_{1 \leq j \leq m} k_j$ is the maximum number of clusters across all dimensions.

³By “guarantee α ”, we mean that the algorithm yields a solution that is guaranteed to have an objective value within a factor of $O(\alpha)$ of the optimum.

⁴Here and in the sequel, the argument m to J_{OPT} denotes the best m -way clustering to avoid confusions about dimensions.

To establish Theorem 5.3, we will first bound the quality of a combination of dimension-wise clusterings for the Frobenius norm, with the help of the Pythagorean Property (Lemma 5.5). It is clear that compressing along only a subset of dimensions achieves lower divergence than clustering along all dimensions. Generalizing an idea of [4], we upper bound the full combined clustering in terms of the (approximately) optimal clustering along a subset of dimensions (Prop. 5.6). Finally, we extend this upper bound to general Bregman divergences and relate it to the optimal tensor clustering.

In the analysis below we assume without loss of generality that $m = 2^h t$ for an integer h (otherwise, pad in empty dimensions). We assume that we have access to an algorithm that can cluster along a subset of t dimensions, while achieving an objective function within a factor α_t of the optimal (for those t dimensions), i.e., $\alpha_t J_{\text{OPT}}(t)$. For example, when $t = 1$ we can use BREG++ (or in theory, even the approximation algorithms of [1]).

Lemma 5.5 (Pythagorean Property). *Let $\mathcal{P} = (\mathbf{P}_1, \dots, \mathbf{P}_t)$, $\mathcal{Q} = (\mathbf{P}_{t+1}, \dots, \mathbf{P}_m)$, and $\mathcal{P}^\perp = (\mathbf{I} - \mathbf{P}_1, \dots, \mathbf{I} - \mathbf{P}_t)$ be combinations of projection matrices \mathbf{P}_j . Then*

$$\|(\mathcal{P}, \mathcal{Q}) \cdot \mathbf{A} + (\mathcal{P}^\perp, \mathcal{R}) \cdot \mathbf{B}\|^2 = \|(\mathcal{P}, \mathcal{Q}) \cdot \mathbf{A}\|^2 + \|(\mathcal{P}^\perp, \mathcal{R}) \cdot \mathbf{B}\|^2, \quad (5.10)$$

where \mathcal{R} is some arbitrary combination of $m - t$ projection matrices.

Proof. Using the inner-product (5.4) we can rewrite (5.10) as

$$\|(\mathcal{P}, \mathcal{Q}) \cdot \mathbf{A} + (\mathcal{P}^\perp, \mathcal{R}) \cdot \mathbf{B}\|^2 = \|(\mathcal{P}, \mathcal{Q}) \cdot \mathbf{A}\|^2 + \|(\mathcal{P}^\perp, \mathcal{R}) \cdot \mathbf{B}\|^2 + 2 \langle (\mathcal{P}, \mathcal{Q}) \cdot \mathbf{A}, (\mathcal{P}^\perp, \mathcal{R}) \cdot \mathbf{B} \rangle.$$

With (5.5) the latter term simplifies to

$$\langle (\mathcal{P}, \mathcal{Q}) \cdot \mathbf{A}, (\mathcal{P}^\perp, \mathcal{R}) \cdot \mathbf{B} \rangle = \langle \mathbf{A}, (\mathcal{P}^T \mathcal{P}^\perp, \mathcal{Q}^T \mathcal{R}) \cdot \mathbf{B} \rangle = 0,$$

thus yielding the claim (5.10). \square

Some more notation. Before diving into the proofs, we outline some more useful notation. Since we can only cluster along t dimensions at a time, we recursively half the initial set of m dimensions until, after $\log(m/t) + 1$ recursions, the sets have length t . Let l denote the level of recursion, starting at $l = \log(m/t) = h$ down to $l = 0$, where the sets have length t . At level l , the sets will have length $2^l t$. Each clustering along a subset of $2^l t$ dimensions is represented by the corresponding $2^l t$ projection matrices. We denote their combination by \mathcal{P}_i^l . At level l , i ranges from 1 to 2^{h-l} .

For illustration, consider an order-8 tensor, and $t = 2$. Then $h = \log(m/t) = 2$, so we will need 3 levels. For simplicity, we always partition the set of dimensions in the middle, i.e. $\{1, \dots, 8\}$ into $\{1, \dots, 4\}$ and $\{5, \dots, 8\}$ and so on, ending with $\{\{1, 2\}, \{3, 4\}, \{5, 6\}, \{7, 8\}\}$. The projection matrix for dimension i is \mathbf{P}_i . The full tensor clustering is $(\mathbf{P}_1, \dots, \mathbf{P}_8)$. So here we get

$$\begin{aligned} \mathcal{P}_1^2 &= (\mathbf{P}_1, \mathbf{P}_2, \mathbf{P}_3, \mathbf{P}_4, \mathbf{P}_5, \mathbf{P}_6, \mathbf{P}_7, \mathbf{P}_8) \\ \mathcal{P}_1^1 &= (\mathbf{P}_1, \mathbf{P}_2, \mathbf{P}_3, \mathbf{P}_4), & \mathcal{P}_2^1 &= (\mathbf{P}_5, \mathbf{P}_6, \mathbf{P}_7, \mathbf{P}_8) \\ \mathcal{P}_1^0 &= (\mathbf{P}_1, \mathbf{P}_2), & \mathcal{P}_2^0 &= (\mathbf{P}_3, \mathbf{P}_4), & \mathcal{P}_3^0 &= (\mathbf{P}_5, \mathbf{P}_6), & \mathcal{P}_4^0 &= (\mathbf{P}_7, \mathbf{P}_8) \end{aligned}$$

To represent a clustering of the tensor along only a subset of dimensions, we pad the corresponding \mathcal{P}_i^l with $m - 2^l t$ identity matrices for the non-clustered dimensions. We refer to this padded collection as \mathcal{Q}_i^l . In the example above, e.g. $\mathcal{Q}_1^0 = (\mathbf{P}_1, \mathbf{P}_2, \mathbf{I}, \mathbf{I}, \mathbf{I}, \mathbf{I}, \mathbf{I}, \mathbf{I})$, $\mathcal{Q}_2^0 = (\mathbf{I}, \mathbf{I}, \mathbf{I}, \mathbf{I}, \mathcal{P}_2^0)$, and $\mathcal{Q}_1^1 = \mathcal{P}_1^1$. With recursive partitions of the dimensions, \mathcal{Q}_i^l subsumes \mathcal{Q}_j^0 for $2^l(i-1) < j \leq 2^l i$: $\mathcal{Q}_i^l = \sum_{j=2^l(i-1)}^{2^l i} \mathcal{Q}_j^0$. The algorithm for the subsets of dimensions will yield the \mathcal{Q}_i^0 and \mathcal{P}_i^0 . The remaining clusterings are simply combinations of those level-0 clusterings. Finally, we refer to the collection of $m - 2^l t$ identity matrices (for simplicity, we assume that they have the correct dimensionalities) as \mathcal{I}^l , so, for instance, $\mathcal{Q}_1^1 = (\mathcal{P}_1^1, \mathcal{I}^1)$.

Note that the order of the dimensions is arbitrary, as long as the index sets remain the same and we reorder the dimensions of all tensors and matrices correspondingly. Hence, we always shift the identity matrices to the back for “ease” of notation. Furnished with this notation, we can now turn towards the details of the proofs. We start with the relation of the combined clustering to a subclustering with the Frobenius norm objective function.

Proposition 5.6. *Let \mathbf{A} be an order- m tensor and $m \geq 2^l t$. The objective function for any $2^l t$ -way clustering $\mathcal{P}_1^l = (\mathcal{P}_1^0, \dots, \mathcal{P}_{2^l}^0)$ can be bounded via the subclusterings along only one set of dimensions of size t :*

$$\|\mathbf{A} - \mathcal{Q}_1^l \cdot \mathbf{A}\|^2 = \|\mathbf{A} - (\mathcal{P}_1^l, \mathcal{I}_1^l) \cdot \mathbf{A}\|^2 \leq \max_{1 \leq j \leq 2^l} 2^l \|\mathbf{A} - \mathcal{Q}_j^0 \cdot \mathbf{A}\|^2. \quad (5.11)$$

Not that the Proposition actually holds for any set of 2^l sub-clusterings by permuting dimensions accordingly:

$$\|\mathbf{A} - \mathcal{Q}_i^l \cdot \mathbf{A}\|^2 \leq \max_{2^{l(i-1)} < j \leq 2^l i} 2^l \|\mathbf{A} - \mathcal{Q}_j^0 \cdot \mathbf{A}\|^2 \quad (5.12)$$

We will prove the proposition for $i = 1$ for ease of notation. If $m = 2^h t$ then the factor is $2^h = m/t$. (If m/t is not a power of 2, then we get the factor with $h = \lceil \log(m/t) \rceil$).

Proof. We prove the proposition by induction on l .

Base case: Let $l = 0$. Then $\mathcal{Q}_j^0 = \mathcal{P}_1^0 = \mathcal{P}_1^l$ and the claim holds trivially.

Induction step: Assume the claim holds for $l \geq 0$. Let us look at a clustering $\mathcal{P}_1^{l+1} = (\mathcal{P}_1^l, \mathcal{P}_2^l)$. Then \mathbf{A} decomposes as

$$\mathbf{A} = \left((\mathcal{P}_1^l, \mathcal{P}_2^l, \mathcal{I}^{l+1}) + ((\mathcal{P}_1^l)^\perp, \mathcal{P}_2^l, \mathcal{I}^{l+1}) + (\mathcal{P}_1^l, (\mathcal{P}_2^l)^\perp, \mathcal{I}^{l+1}) + ((\mathcal{P}_1^l)^\perp, (\mathcal{P}_2^l)^\perp, \mathcal{I}^{l+1}) \right) \cdot \mathbf{A},$$

whereby the Pythagorean Property 5.5 yields

$$\|\mathbf{A} - (\mathcal{P}_1^{l+1}, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2 = \|\mathbf{A} - (\mathcal{P}_1^l, \mathcal{P}_2^l, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2 \quad (5.13)$$

$$= \|((\mathcal{P}_1^l)^\perp, \mathcal{P}_2^l, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2 + \|(\mathcal{P}_1^l, (\mathcal{P}_2^l)^\perp, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2 + \|((\mathcal{P}_1^l)^\perp, (\mathcal{P}_2^l)^\perp, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2, \quad (5.14)$$

Assuming without loss of generality that

$$\|((\mathcal{P}_1^l)^\perp, \mathcal{P}_2^l, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2 \geq \|\mathcal{P}_1^l, (\mathcal{P}_2^l)^\perp, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|,$$

we have from (5.13) and (5.14) the following inequalities

$$\begin{aligned} & \|\mathbf{A} - (\mathcal{P}_1^l, \mathcal{P}_2^l, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2 \\ & \leq 2(\|((\mathcal{P}_1^l)^\perp, \mathcal{P}_2^l, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2 + \|((\mathcal{P}_1^l)^\perp, (\mathcal{P}_2^l)^\perp, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2) \\ & = 2\|((\mathcal{P}_1^l)^\perp, \mathcal{P}_2^l, \mathcal{I}^{l+1}) \cdot \mathbf{A} + ((\mathcal{P}_1^l)^\perp, (\mathcal{P}_2^l)^\perp, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2 \\ & = 2\|((\mathcal{P}_1^l)^\perp, \mathcal{I}^l, \mathcal{I}^{l+1}) \cdot \mathbf{A}\|^2 \\ & = 2\|\mathbf{A} - \mathcal{Q}_1^l \cdot \mathbf{A}\|^2 \\ & \leq 2 \max_{1 \leq j \leq 2} \|\mathbf{A} - \mathcal{Q}_j^l \cdot \mathbf{A}\|^2 \\ & \leq 2 \cdot 2^l \max_{1 \leq j \leq 2^{l+1}} \|\mathbf{A} - \mathcal{Q}_j^0 \cdot \mathbf{A}\|^2, \end{aligned}$$

where the last step follows from the induction hypothesis and (5.12), and the two norm terms are combined using the Pythagorean Property. \square

The proof of Theorem 5.3 generalizes Proposition 5.6 to arbitrary Bregman divergences, and then relates the objective for the sub-clustering to the objective for the full m -way clustering.

Proof. (Theorem 5.3.) Let $m = 2^h t$. Via the approximation algorithm with guarantee α_t , we cluster the subsets of size t to obtain \mathcal{P}_i^0 . Let $\check{\mathcal{Q}}_i^0$ be the optimal subclustering of dimension set i , i.e. the result that \mathcal{Q}_i^0 would be if $\alpha_t = 1$. With σ_1 and σ_2 from (5.9), we can bound the combination \mathcal{P}_1^h of the sub-clusterings \mathcal{P}_i^0 :

$$\begin{aligned} B_f(\mathbf{A}, \mathcal{P}_1^h \cdot \mathbf{A}) &= \sum_{i_1, \dots, i_m} B_f(a_{i_1, \dots, i_m}, \mu_{I_1(i_1), \dots, I_m(i_m)}) \\ &\leq \sigma_2 \sum_{i_1, \dots, i_m} (a_{i_1, \dots, i_m} - \mu_{I_1(i_1), \dots, I_m(i_m)})^2 = \sigma_2 \|\mathbf{A} - \mathcal{P}_1^h \cdot \mathbf{A}\|^2 \\ &\leq 2^h \sigma_2 \max_j \|\mathbf{A} - \mathcal{Q}_j^0 \cdot \mathbf{A}\|^2 \end{aligned} \quad (5.15)$$

$$\leq 2^h \sigma_2 \sigma_1^{-1} \max_j B_f(\mathbf{A}, \mathcal{Q}_j^0 \cdot \mathbf{A}) \leq 2^h \sigma_2 \sigma_1^{-1} \alpha_t \max_j B_f(\mathbf{A}, \check{\mathcal{Q}}_j^0 \cdot \mathbf{A}). \quad (5.16)$$

Inequality (5.15) follows from Proposition 5.6, and (5.16) from the guarantee of the algorithm we used to get the separate sub-clustering \mathcal{Q}_j^0 .

Let us look at the relation between an optimal clustering $\check{\mathcal{Q}}^0$ of an arbitrary subset of t dimensions and the optimal full tensor clustering $\check{\mathcal{Q}}^0$ of all $2^h t$ dimensions. Let $\check{\mathbf{C}}_j$ be the cluster indicator matrices of the clustered dimensions in $\check{\mathcal{Q}}^0$, and $\check{\mathbf{C}}_j$ their normalized versions such that $\check{\mathcal{Q}}^0 = (\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t, \mathcal{I}^0)(\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t, \mathcal{I}^0)^\top$ ⁵. By definition, $\check{\mathcal{Q}}^0$ solves

$$\underset{\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t, \mathbf{M}}{\text{minimize}} \quad B_f(\mathbf{A}, (\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t, \mathcal{I}^0) \cdot \mathbf{M}), \quad \text{s.t. } \check{\mathbf{C}}_j \in \{0, 1\}^{n_j \times k_j},$$

with $(\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t, \mathcal{I}^0) \cdot \check{\mathbf{M}} = \check{\mathcal{Q}}^0 \cdot \mathbf{A}$. In that respect, $\check{\mathcal{Q}}^0$ even beats the sub-clustering $(\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t)$ taken from the optimal full m -way clustering $\check{\mathcal{Q}}_1^h = \check{\mathcal{P}}_1^h$, i.e.

$$\begin{aligned} B_f(\mathbf{A}, (\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t, \mathcal{I}^0) \cdot ((\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t, \mathcal{I}^0)^\top \cdot \mathbf{A})) &\leq \min_{\mathbf{B}} B_f(\mathbf{A}, (\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t, \mathcal{I}^0) \cdot \mathbf{B}) \\ &\leq B_f(\mathbf{A}, (\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_t, \mathcal{I}^0)(I, \dots, I, \check{\mathbf{C}}_{t+1}, \dots, \check{\mathbf{C}}_m) \cdot \check{\mathbf{M}}) = B_f(\mathbf{A}, \check{\mathcal{P}}_1^h \cdot \mathbf{A}), \end{aligned}$$

where $\check{\mathbf{M}}$ is the tensor of means for the optimal m -way clustering, $(\check{\mathbf{C}}_1, \dots, \check{\mathbf{C}}_m) \cdot \check{\mathbf{M}} = \check{\mathcal{P}}_1^h \cdot \mathbf{A}$. Combining this bound with (5.16) yields the final bound for the combined clustering \mathcal{P}_1^h ,

$$B_f(\mathbf{A}, \mathcal{P}_1^h \cdot \mathbf{A}) \leq 2^h \sigma_2 \sigma_1^{-1} \alpha_t B_f(\mathbf{A}, \check{\mathcal{P}}_1^h \cdot \mathbf{A}) = 2^h \sigma_2 \sigma_1^{-1} \alpha_t J_{OPT}(m),$$

and completes the proof of the theorem. \square

Special cases of Theorem 5.3 are Euclidean m -way co-clustering and Bregman co-clustering.

Corollary 5.7 (m -way Euclidean tensor clustering). *For $f = \frac{1}{2}x^2$, we have $\sigma_1 = \sigma_2 = 1$. Thus, using KMEANS++ as the base algorithm for obtaining the combined clustering \mathcal{P}_1^h leads to the approximation guarantee:*

$$E[J(\mathcal{P}_1^h)] \leq 8m(\log K^* + 2)J_{OPT}(m).$$

Corollary 5.8 (Co-clustering). *Let $\mathcal{C}_R, \mathcal{C}_C$ be the clusterings of the rows and columns of \mathbf{A} obtained via BREG++, and $(\mathcal{C}_R, \mathcal{C}_C)$ the induced co-clustering. Then the expected objective value J is bounded as*

$$E[J(\mathcal{C}_R, \mathcal{C}_C)] \leq 8 \frac{\sigma_2^2}{\sigma_1^2} \left(1 + \frac{\sigma_2}{\sigma_1}\right) (\log K^* + 2) J_{OPT}(2),$$

where σ_1 and σ_2 are as defined in Equation 5.9.

⁵Without loss of generality, we again assume that the clustered dimensions are the first t . Otherwise, permute the dimensions

References

- [1] M. R. Ackermann and Johannes Blömer. Coresets and Approximate Clustering for Bregman Divergences. In *Proc. 20th ACM-SIAM Symposium on Discrete Algorithms (SODA '09)*, 2009. To appear.
- [2] M. R. Ackermann, J. Blömer, and C. Sohler. Clustering for metric and non-metric distance measures. In *ACM-SIAM SODA*, April 2008.
- [3] S. Agarwal, J. Lim, L. Zelnik-Manor, P. Perona, D. Kriegman, and S. Belongie. Beyond pairwise clustering. In *IEEE CVPR*, 2005.
- [4] A. Anagnostopoulos, A. Dasgupta, and R. Kumar. Approximation algorithms for co-clustering. In *PODS*, 2008.
- [5] D. Arthur and S. Vassilvitskii. k -means++: The Advantages of Careful Seeding. In *SODA*, pages 1027–1035, 2007.
- [6] L. D. Baker and A. McCallum. Distributional clustering for text classification. In *Proc. SIGIR*, 1998.
- [7] A. Banerjee, I. S. Dhillon, J. Ghosh, and S. Sra. Clustering on the Unit Hypersphere using von Mises-Fisher Distributions. *JMLR*, 6:1345–1382, Sep 2005.
- [8] A. Banerjee, S. Merugu, I. S. Dhillon, and J. Ghosh. Clustering with Bregman Divergences. *JMLR*, 6(6):1705–1749, October 2005.
- [9] A. Banerjee, S. Basu, and S. Merugu. Multi-way Clustering on Relation Graphs. In *SIAM Data Mining*, 2007.
- [10] A. Banerjee, I. S. Dhillon, J. Ghosh, S. Merugu, and D. S. Modha. A Generalized Maximum Entropy Approach to Bregman Co-clustering and Matrix Approximation. *JMLR*, 8:1919–1986, 2007.
- [11] R. Bekkerman, R. El-Yaniv, and A. McCallum. Multi-way distributional clustering via pairwise interactions. In *ICML*, 2005.
- [12] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, March 2004.
- [13] Y. Censor and S. A. Zenios. *Parallel Optimization: Theory, Algorithms, and Applications*. Oxford University Press, 1997.
- [14] K. Chaudhuri and A. McGregor. Finding metric structure in information theoretic clustering. In *Conf. on Learning Theory, COLT*, July 2008.
- [15] Y. Cheng and G. Church. Biclustering of expression data. In *Proceedings ISMB*, pages 93–103. AAAI Press, 2000.
- [16] H. Cho, I. S. Dhillon, Y. Guan, and S. Sra. Minimum Sum Squared Residue based Co-clustering of Gene Expression data. In *Proc. 4th SIAM International Conference on Data Mining (SDM)*, pages 114–125, Florida, 2004. SIAM.
- [17] W. Fernandez de la Vega, M. Karpinski, C. Kenyon, and Y. Rabani. Approximation schemes for clustering problems. In *STOC'03*, 2003.
- [18] V. de Silva and L.-H. Lim. Tensor Rank and the Ill-Posedness of the Best Low-Rank Approximation Problem. *SIAM J. on Matrix Analysis and Applications*, 30(3):1084–1127, 2008.
- [19] I. S. Dhillon, S. Mallela, and R. Kumar. A divisive information-theoretic feature clustering algorithm for text classification. *JMLR*, 3:1265–1287, March 2003.

- [20] I. S. Dhillon, Y. Guan, and B. Kulis. Weighted Graph Cuts without Eigenvectors A Multilevel Approach. *IEEE Trans. Pattern Anal. Mach. Intell.*, 29(11):1944–1957, 2007.
- [21] Inderjit S. Dhillon, Yuqiang Guan, and J. Kogan. Iterative Clustering of High Dimensional Text Data Augmented by Local Search. *IEEE International Conference on Data Mining*, page 131, 2002.
- [22] P. Drineas, A. Frieze, R. Kannan, S. Vempala, and V. Vinay. Clustering large graphs via the singular value decomposition. *Machine Learning*, 56:9–33, 2004.
- [23] V. M. Govindu. A tensor decomposition for geometric grouping and segmentation. In *IEEE CVPR*, 2005.
- [24] W. H. Greub. *Multilinear Algebra*. Springer, 1967.
- [25] S. Har-Peled and S. Mazumdar. On coresets for k -means and k -median clustering. In *STOC'04*, 2004.
- [26] J. A. Hartigan. *Clustering Algorithms*. Wiley, 1975.
- [27] J. A. Hartigan. Direct clustering of a data matrix. *Journal of the American Statistical Association*, 67(337):123–129, March 1972.
- [28] M. Inaba, N. Katoh, and H. Imai. Applications of weighted voronoi diagrams and randomization to variance based k -clustering. In *SCG'94: Proc. 10th Symp. on Comp. Geo.*, 1994.
- [29] Tapas Kanungo, David M. Mount, Nathan S. Netanyahu, Christine D. Piatko, Ruth Silverman, and Angela Y. Wu. A local search approximation algorithm for k -means clustering. *Comput. Geom. Theory Appl.*, 28(2-3):89–112, 2004.
- [30] B. Kulis, S. Basu, I. Dhillon, and R. Mooney. Semi-supervised graph clustering: a kernel approach. *Machine Learning*, 2008. To appear.
- [31] A. Kumar, Y. Sabharwal, and S. Sen. A simple linear time $(1 + \epsilon)$ -approximation algorithms for k -means clustering in any dimensions. In *IEEE Symp. on Foundations of Comp. Sci.*, 2004.
- [32] S. P. Lloyd. Least squares quantization in PCM. *IEEE Tran. on Inf. Theory*, 28(2):129–136, 1982.
- [33] B. Long, X. Wu, and Z. Zhang. Unsupervised learning on k -partite graphs. In *SIGKDD*, 2006.
- [34] R. R. Mettu and C. G. Plaxton. Optimal Time Bounds for Approximate Clustering. *Mach. Learn.*, 56(1-3):35–60, 2004.
- [35] F. Nielsen, J.-D. Boissonnat, and R. Nock. On Bregman Voronoi Diagrams. In *ACM-SIAM Symposium on Discrete Algorithms*, pages 746–755, 2007.
- [36] R. Ostrovsky, Y. Rabani, L. Schulman, and C. Swamy. The effectiveness of Lloyd-type methods for the k -means problem. In *FOCS*, 2006.
- [37] F. C. Pereira, N. Tishby, and L. Lee. Distributional clustering of english words. In *Meeting of the Assoc. for Comp. Linguistics*, 1993.
- [38] K. Puolamäki, S. Hanhijärvi, and G. C. Garriga. An approximation ratio for biclustering. *Inf. Process. Lett.*, 108(2):45–49, 2008.
- [39] R. T. Rockafellar. *Convex Analysis*. Princeton Univ. Press, 1970.
- [40] A. Shashua, R. Zass, and T. Hazan. Multi-way Clustering Using Super-Symmetric Non-negative Tensor Factorization. *LNCS*, 3954:595–608, 2006.