

SVD, discrepancy, and regular structure of contingency tables

Marianna Bolla *

Institute of Mathematics, Budapest University of Technology and Economics

Abstract

We will use the factors obtained by correspondence analysis to find biclustering of a contingency table such that the row-column cluster pairs are regular, i.e., they have small discrepancy. In our main theorem, the constant of the so-called volume-regularity is related to the SVD of the normalized contingency table. Our result is applicable to two-way cuts when both the rows and columns are divided into the same number of clusters, thus extending partly the result of [7] estimating the discrepancy of a contingency table by the second largest singular value of the normalized table (one-cluster, rectangular case), and partly the result of [5] for estimating the constant of volume-regularity by the structural eigenvalues and the distances of the corresponding eigen-subspaces of the normalized modularity matrix of an edge-weighted graph (several clusters, symmetric case).

Key words: Normalized contingency table, Regular row-column pairs, Biclustering, Discrepancy, Cluster variances, Directed graphs

1 Introduction

A typical problem of contemporary cluster analysis is to find relatively small number of groups of objects, belonging to rows and columns of a contingency table which exhibit homogeneous behavior with respect to each other and do not differ significantly in size. To make inferences on the separation that can be achieved for a given number of clusters, minimum normalized two-way cuts are investigated and related to the SVD of the correspondence matrix.

* Research supported in part by the Hungarian National Research Grants OTKA 76481 and OTKA-KTIA 77778; further, by the TÁMOP-4.2.2.B-10/1-2010-0009 project.

Email address: marib@math.bme.hu (Marianna Bolla).

Contingency tables are rectangular arrays with nonnegative, real entries. One example is the keyword–document matrix. Here the entries are associations between documents and words. Based on network data, the entry in the i th row and j th column is the relative frequency of word j in document i . Latent semantic indexing looks for real scores of the documents and keywords such that the score of any document be proportional to the total scores of the keywords occurring in it, and vice versa, the score of any keyword be proportional to the total scores of the documents containing it. Not surprisingly, the solution is given by the SVD of the binary table, where the document- and keyword-scores are the coordinates of the left and right singular vectors corresponding to its largest non-trivial singular value which gives the constant of proportionality.

This idea is generalized in [10] in the following way. We can think of the above relation between keywords and documents as the relation with respect to the most important topic (or context, or factor). After this, we are looking for another scoring with respect to the second topic, up to k (where k is a positive integer not exceeding the rank of the table). The solution is given by the singular vector pairs corresponding to the k largest singular values of the table.

If a scoring system is endowed with the marginal measures, the problem can be formulated in terms of correspondence analysis and correlation maximization. The problem is solved by the SVD of the correspondence matrix (normalized contingency table), where the singular vector pairs are also transformed, see [4]. In this way, instead of scores, the documents and keywords have k -dimensional representatives, based of which further investigations, spacial representation, or biclustering can be performed that finds simultaneous clustering of the rows and columns of the table with densities as homogeneous as possible between the keyword–document cluster pairs.

The problem is also related to the Pagerank (see [11]) and to microarray analysis (see [12]) when we want to find clusters of the rows and columns of a microarray, simultaneously. Here rows correspond to genes and columns to different conditions, whereas the entries are expression levels of genes under specific conditions. We also look for a bipartition of the genes and conditions such that genes in the same cluster equally (not necessarily weakly or strongly) influence conditions of the same cluster.

In Section 2 we deal with the singular value decomposition (SVD) of a correspondence matrix. In Section 3 we relate it to normalized two-way cuts of the contingency table, while in Section 4 the constant of volume-regularity of row–column clusters pairs is estimated by means of the SVD. Section 5 is devoted to discussion, application and possible extension to directed graphs.

2 SVD of contingency tables and correspondence matrices

Let \mathbf{C} be a contingency table on row set $Row = \{1, \dots, n\}$ and column set $Col = \{1, \dots, m\}$, where \mathbf{C} is $n \times m$ matrix of entries $c_{ij} \geq 0$. Without loss of generality, we suppose that there are not identically zero rows or columns. Here c_{ij} is some kind of association between the objects behind row i and column j , where 0 means no interaction at all.

Let the row- and column-sums of \mathbf{C} be

$$d_{row,i} = \sum_{j=1}^m c_{ij} \quad (i = 1, \dots, n) \quad \text{and} \quad d_{col,j} = \sum_{i=1}^n c_{ij} \quad (j = 1, \dots, m)$$

which are collected in the main diagonals of the $n \times n$ and $m \times m$ diagonal matrices \mathbf{D}_{row} and \mathbf{D}_{col} , respectively.

For a given integer $1 \leq k \leq \min\{n, m\}$, we are looking for k -dimensional representatives $\mathbf{r}_1, \dots, \mathbf{r}_n$ of the rows and $\mathbf{c}_1, \dots, \mathbf{c}_m$ of the columns such that they minimize the objective function

$$Q_k = \sum_{i=1}^n \sum_{j=1}^m c_{ij} \|\mathbf{r}_i - \mathbf{c}_j\|^2 \quad (1)$$

subject to

$$\sum_{i=1}^n d_{row,i} \mathbf{r}_i \mathbf{r}_i^T = \mathbf{I}_k, \quad \sum_{j=1}^m d_{col,j} \mathbf{c}_j \mathbf{c}_j^T = \mathbf{I}_k. \quad (2)$$

When minimized, the objective function Q_k favors k -dimensional placement of the rows and columns such that representatives of highly associated rows and columns are forced to be close to each other. As we will see, this is equivalent to the problem of correspondence analysis.

Indeed, let us put both the objective function and the constraints in a more favorable form. Let \mathbf{X} be the $n \times k$ matrix of rows $\mathbf{r}_1^T, \dots, \mathbf{r}_n^T$; let $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$ denote the columns of \mathbf{X} , for which fact we use the notation $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_k)$. Similarly, let \mathbf{Y} be the $m \times k$ matrix of rows $\mathbf{c}_1^T, \dots, \mathbf{c}_m^T$; let $\mathbf{y}_1, \dots, \mathbf{y}_k \in \mathbb{R}^m$ denote the columns of \mathbf{Y} , i.e., $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_k)$. Hence, the constraints (2) can be formulated like

$$\mathbf{X}^T \mathbf{D}_{row} \mathbf{X} = \mathbf{I}_k, \quad \mathbf{Y}^T \mathbf{D}_{col} \mathbf{Y} = \mathbf{I}_k.$$

With this notation, the objective function (1) is

$$\begin{aligned} Q_k &= \sum_{i=1}^n \sum_{j=1}^m c_{ij} \|\mathbf{r}_i - \mathbf{c}_j\|^2 = \sum_{i=1}^n d_{row,i} \|\mathbf{r}_i\|^2 + \sum_{j=1}^m d_{col,j} \|\mathbf{c}_j\|^2 - \sum_{i=1}^n \sum_{j=1}^m c_{ij} \mathbf{r}_i^T \mathbf{c}_j \\ &= 2k - \text{tr} \mathbf{X}^T \mathbf{C} \mathbf{Y} = 2k - \text{tr} (\mathbf{D}_{row}^{1/2} \mathbf{X})^T (\mathbf{D}_{row}^{-1/2} \mathbf{C} \mathbf{D}_{col}^{-1/2}) (\mathbf{D}_{col}^{1/2} \mathbf{Y}), \end{aligned} \quad (3)$$

where the matrix $\mathbf{C}_{corr} = \mathbf{D}_{row}^{-1/2} \mathbf{C} \mathbf{D}_{col}^{-1/2}$ is the *correspondence matrix* (*normalized contingency table*) belonging to the table \mathbf{C} , see [4]. If we multiply all the entries of \mathbf{C} with the same positive constant, the correspondence matrix \mathbf{C}_{corr} will not change. Therefore, without the loss of generality, $\sum_{i=1}^n \sum_{j=1}^m c_{ij} = 1$ will be supposed in the sequel. The correspondence matrix has SVD

$$\mathbf{C}_{corr} = \sum_{i=1}^r s_i \mathbf{v}_i \mathbf{u}_i^T, \quad (4)$$

where $r \leq \min\{n, m\}$ is the rank of \mathbf{C}_{corr} , or equivalently (since there are not identically zero rows or columns), the rank of \mathbf{C} . Here $1 = s_1 \geq s_2 \geq \dots \geq s_r > 0$ are the non-zero singular values of \mathbf{C}_{corr} , and 1 is a single singular value if \mathbf{C}_{corr} , or equivalently, \mathbf{C} is non-decomposable ($\mathbf{C}\mathbf{C}^T$ is irreducible). In this case $\mathbf{v}_1 = (\sqrt{d_{row,1}}, \dots, \sqrt{d_{row,n}})^T$ and $\mathbf{u}_1 = (\sqrt{d_{col,1}}, \dots, \sqrt{d_{col,m}})^T$.

Note that the singular spectrum of a decomposable contingency table can be composed from the singular spectra of its non-decomposable parts, as well as their singular vector pairs. Therefore, in the future, the non-decomposability of the underlying contingency table will be supposed. In this way, the following representation theorem for contingency tables can be formulated.

Theorem 1 *Let \mathbf{C} be a non-decomposable contingency table with SVD (4) of its correspondence matrix \mathbf{C}_{corr} . Let $k \leq r$ be a positive integer such that $s_k > s_{k+1}$. Then the minimum of (1) subject to (2) is $2k - \sum_{i=1}^k s_i$ and it is attained with the optimum row representatives $\mathbf{r}_1^*, \dots, \mathbf{r}_n^*$ and column representatives $\mathbf{c}_1^*, \dots, \mathbf{c}_m^*$, the transposes of which are row vectors of $\mathbf{X}^* = \mathbf{D}_{row}^{-1/2}(\mathbf{v}_1, \dots, \mathbf{v}_k)$ and $\mathbf{Y}^* = \mathbf{D}_{col}^{-1/2}(\mathbf{u}_1, \dots, \mathbf{u}_k)$, respectively.*

PROOF. In view of (3), we have to maximize

$$\text{tr} (\mathbf{D}_{row}^{1/2} \mathbf{X})^T \mathbf{C}_{corr} (\mathbf{D}_{col}^{1/2} \mathbf{Y})$$

under the given constraints. Separation theorems for the singular value decomposition (see e.g., [1] and [13]) are applicable, yielding the required statement.

The vectors $\mathbf{r}_1^*, \dots, \mathbf{r}_n^*$ and $\mathbf{c}_1^*, \dots, \mathbf{c}_m^*$ giving the optimum in the above theorem are called *optimum k -dimensional representatives* of the rows and columns,

while the transformed singular vectors $\mathbf{D}_{row}^{-1/2}\mathbf{v}_1, \dots, \mathbf{D}_{row}^{-1/2}\mathbf{v}_k$ and $\mathbf{D}_{col}^{-1/2}\mathbf{u}_1, \dots, \mathbf{D}_{col}^{-1/2}\mathbf{u}_k$ are called *vector components* of the rows and columns taking part in the k -dimensional representation.

Observe that the dimension k does not play an important role here: the vector components can be included successively up to a k such that $s_k > s_{k+1}$. We remark that the singular vectors can arbitrarily be chosen in the isotropic subspaces corresponding to possible multiple singular values, under the orthogonality conditions. Further, provided that 1 is a single singular value, the first vector components are the constantly $\mathbf{1}$ vectors in \mathbb{R}^n and \mathbb{R}^m , respectively, and hence, the k -dimensional representation is realized in a $(k-1)$ -dimensional hyperplane of \mathbb{R}^k .

A symmetric contingency table corresponds to a weighted graph, and our correspondence matrix is the identity minus the normalized Laplacian, called normalized modularity matrix in [5]. In another view, a contingency table can be considered as part of the weight matrix of a bipartite graph on vertex set $Row \cup Col$. However, it would be tedious to always distinguish between these two types of vertices, we rather use the framework of correspondence analysis, and formulate our statements in terms of rows and columns.

3 Normalized two-way cuts of contingency tables

Given the $n \times m$ contingency table \mathbf{C} on row set Row and column set Col , further, an integer k ($0 < k \leq r$), we want to simultaneously partition its rows and columns into disjoint, nonempty subsets

$$Row = R_1 \cup \dots \cup R_k, \quad Col = C_1 \cup \dots \cup C_k$$

such that the *cuts* $c(R_a, C_b) = \sum_{i \in R_a} \sum_{j \in C_b} c_{ij}$ ($a, b = 1, \dots, k$) between the row-column cluster pairs be as homogeneous as possible. For this requirement, the following so-called *normalized two-way cut* of the contingency table with respect to the above k -partitions $P_{row} = (R_1, \dots, R_k)$ and $P_{col} = (C_1, \dots, C_k)$ of its rows and columns and the collection of signs σ is defined as follows:

$$\nu_k(P_{row}, P_{col}, \sigma) = \sum_{a=1}^k \sum_{b=1}^k \left(\frac{1}{\text{Vol}(R_a)} + \frac{1}{\text{Vol}(C_b)} + \frac{2\sigma_{ab}\delta_{ab}}{\sqrt{\text{Vol}(R_a)\text{Vol}(C_b)}} \right) c(R_a, C_b),$$

where

$$\text{Vol}(R_a) = \sum_{i \in R_a} d_{row,i} = \sum_{i \in R_a} \sum_{j=1}^m c_{ij}, \quad \text{Vol}(C_b) = \sum_{j \in C_b} d_{col,j} = \sum_{j \in C_b} \sum_{i=1}^n c_{ij}$$

are volumes of the clusters, δ_{ab} is the Kronecker delta, and the sign σ_{ab} is equal to 1 or -1 (it only has relevance in the $a = b$ case, when it helps balancing between the volumes of the same index row and column clusters), $\sigma := (\sigma_{11}, \dots, \sigma_{kk})$. We want to minimize the above normalized two-way cut with respect to all possible k -partitions $\mathcal{P}_{row,k}$ and $\mathcal{P}_{col,k}$ of the rows and columns, further, to σ , simultaneously. The objective function penalizes row- and column clusters of extremely different volumes in the $a \neq b$ case, whereas in the $a = b$ case σ_{aa} moderates the balance between $\text{Vol}(R_a)$ and $\text{Vol}(C_a)$.

Definition 2 *The normalized two-way cut of the contingency table \mathbf{C} is*

$$\nu_k(\mathbf{C}) = \min_{P_{row}, P_{col}, \sigma} \nu_k(P_{row}, P_{col}, \sigma).$$

Theorem 3 *Let $1 = s_1 > s_2 \dots \geq s_r$ be the positive singular values of the correspondence matrix belonging to the non-decomposable contingency table \mathbf{C} of rank r , and $k \leq r$ be a positive integer. Then*

$$\nu_k(\mathbf{C}) \geq 2k - \sum_{i=1}^k s_i.$$

PROOF. We will show that $\nu_k(P_{row}, P_{col}, \sigma)$ is Q_k in the special representation, where the column vectors of \mathbf{X} and \mathbf{Y} are partition vectors belonging to P_{row} and P_{col} , respectively. Therefore, the statement follows, as the overall minimum is $2k - \sum_{i=1}^k s_i$. Indeed, let the i th coordinate of the left vector component \mathbf{x}_a be

$$x_{ia} := \frac{1}{\sqrt{\text{Vol}(R_a)}} \quad \text{if } i \in R_a, a = 1, \dots, k;$$

similarly, let the j th coordinate of the right vector component \mathbf{y}_b be

$$y_{jb} = \sigma_{bb} \frac{1}{\sqrt{\text{Vol}(C_b)}} \quad \text{if } j \in C_b, b = 1, \dots, k,$$

otherwise the coordinates are zeros. With this, the matrices \mathbf{X} and \mathbf{Y} satisfy the conditions imposed on the representatives, further

$$\|\mathbf{r}_i - \mathbf{c}_j\|^2 = \frac{1}{\text{Vol}(R_a)} + \frac{1}{\text{Vol}(C_b)} + \frac{2\sigma_{bb}\delta_{ab}}{\sqrt{\text{Vol}(R_a)\text{Vol}(C_b)}}, \quad \text{if } i \in R_a, j \in C_b.$$

In case of a symmetric contingency table (weight matrix \mathbf{W} of an edge-weighted graph), we get the same result with the representation based on the eigenvectors belonging to the largest absolute value eigenvalues of the matrix $\mathbf{D}^{-1/2}\mathbf{W}\mathbf{D}^{-1/2}$, where $\mathbf{D} = \mathbf{D}_{row} = \mathbf{D}_{col}$, see [5]. However, $\nu_k(P_{row}, P_{col}, \sigma)$

cannot always be directly related to the normalized cut, except the following two special cases.

- When the $k - 1$ largest absolute value eigenvalues of the normalized modularity matrix are all positive, or equivalently, if the k smallest eigenvalues (including the zero) of the normalized Laplacian matrix are farther from 1 than any other eigenvalue which is greater than 1. In this case the $k - 1$ largest singular values (apart from the 1) of the correspondence matrix are identical to the $k - 1$ largest eigenvalues of the normalized modularity matrix, and the left and right singular vectors are identical to the corresponding eigenvector with the same orientation. Consequently, for the k -dimensional row- and column-representatives $\mathbf{r}_i = \mathbf{c}_i$ ($i = 1, \dots, n = m$) holds. With the choice $\sigma_{bb} = 1$ ($b = 1, \dots, k$), the corresponding $\nu_k(\mathbf{C})$ is twice the normalized cut of our weighted graph in which weights of edges within the clusters do not count. In this special situation, the normalized two-way cut also favors k -partitions with low inter-cluster edge-densities (consequently, intra-cluster densities tend to be large, as they do not count in the objective function).
- When the $k - 1$ largest absolute value eigenvalues of the normalized modularity matrix are all negative, then $\mathbf{r}_i = -\mathbf{c}_i$ for all $(k - 1)$ -dimensional row and column representatives, and any (but only one) of them can be the corresponding vertex representative. Now $\nu_k(\mathbf{C})$, which is attained with the choice $\sigma_{bb} = -1$ ($b = 1, \dots, k$), differs from the normalized cut in that it also counts the edge-weights within the clusters. Indeed, in the $a = b$, $R_a = C_a = V_a$ case

$$\|\mathbf{r}_i - \mathbf{c}_j\|^2 = \frac{1}{\text{Vol}(V_a)} + \frac{1}{\text{Vol}(V_b)} + \frac{2}{\sqrt{\text{Vol}(V_a)\text{Vol}(V_b)}} = \frac{4}{\text{Vol}(V_a)}$$

if $i, j \in V_a$. Here, by minimizing the normalized k -way cut, rather a so-called anti-community structure is detected in that $c(R_a, C_a) = c(V_a, V_a)$ is suppressed to compensate for the term $\frac{4}{\text{Vol}(V_a)}$.

We remark that Ding et al. [9] treat this problem for two row- and column-clusters and minimize another objective function such that it favors 2-partitions where $c(R_1, C_2)$ and $c(R_2, C_1)$ are small compared to $c(R_1, C_1)$ and $c(R_2, C_2)$. The solution is also given by the transformed $\mathbf{v}_2, \mathbf{u}_2$ pair. However, it is the objective function Q_k which best complies with the SVD of the correspondence matrix, and hence, gives the continuous relaxation of the normalized cut minimization problem. The idea of Ding et al. could be naturally extended to the case of several, but the same number of row and column clusters, and it may work well in the keyword-document classification problem. Though, in some real-life problems, e.g., clustering genes and conditions of microarrays, we rather want to find clusters of similarly functioning genes that equally (not especially weakly or strongly) influence conditions of the same cluster.

Dhillon [8] also suggests a multipartition algorithm that runs the k-means algorithm simultaneously for the row and column representatives.

4 Regular row-column cluster pairs

Let us start with the one-cluster case. Let \mathbf{C} be an $n \times m$ contingency table and let \mathbf{C}_{corr} be the correspondence matrix belonging to it. The Expander Mixing Lemma for edge-weighted graphs naturally extends to this situation, see the following result of [7].

Proposition 4 *Let \mathbf{C} be a non-decomposable contingency table (i.e., $\mathbf{C}\mathbf{C}^T$ is irreducible) on row set Row and column set Col , and of total volume 1. Then for all $R \subset Row$ and $C \subset Col$*

$$|c(R, C) - \text{Vol}(R)\text{Vol}(C)| \leq s_2 \sqrt{\text{Vol}(R)\text{Vol}(C)},$$

where s_2 is the largest but 1 singular value of the normalized contingency table \mathbf{C}_{corr} .

Since the spectral gap of \mathbf{C}_{corr} is $1 - s_2$, in view of the above Expander Mixing Lemma, 'large' spectral gap is an indication that the weighted cut between any row and column subset of the contingency table is near to what is expected in a random table. The following notion of discrepancy is just measures the deviation from this random situation. The discrepancy (see [7]) of the contingency table \mathbf{C} of total volume 1 is the smallest $\alpha > 0$ such that for all $R \subset Row$ and $C \subset Col$

$$|c(R, C) - \text{Vol}(R)\text{Vol}(C)| \leq \alpha \sqrt{\text{Vol}(R)\text{Vol}(C)}.$$

In view of this, the result of Theorem 4 can be interpreted as follows: α singular value separation causes α discrepancy, where the singular value separation is the second largest singular value of the normalized contingency table, which is the smaller the bigger the separation between the largest singular value (the 1) of the normalized contingency table and the other singular values is. Based on the ideas of [2] and [6], Butler [7] proves the converse of the Expander Mixing Lemma for contingency tables, namely that

$$s_2 \leq 150\alpha(1 - 8 \log \alpha).$$

Now we extend the notion of discrepancy to volume-regular pairs.

Definition 5 *The row-column cluster pair $R \subset Row$, $C \subset Col$ of the contingency table \mathbf{C} of total volume 1 is γ -volume regular if for all $X \subset R$ and $Y \subset C$ the relation*

$$|c(X, Y) - \rho(R, C)\text{Vol}(X)\text{Vol}(Y)| \leq \gamma \sqrt{\text{Vol}(R)\text{Vol}(C)} \quad (5)$$

holds, where $\rho(R, C) = \frac{c(R, C)}{\text{vol}(R)\text{vol}(C)}$ is the relative inter-cluster density of the row-column pair R, C .

Now we will show that for given k , if the clusters are formed via applying the weighted k -means algorithm for the optimal row- and column representatives, respectively, then the so obtained row-column cluster pairs are homogeneous in the sense that they form equally dense parts of the contingency table. More precisely, the constant γ of the volume regularity of the pairs will be related to the SVD of \mathbf{C}_{corr} . To this end, we introduce the following notion.

The weighted k -variance of the k -dimensional row representatives is defined by

$$S_k^2(\mathbf{X}) = \min_{(R_1, \dots, R_k)} \sum_{a=1}^k \sum_{j \in R_a} d_{row,j} \|\mathbf{r}_j - \bar{\mathbf{r}}_a\|^2, \quad (6)$$

where $\bar{\mathbf{r}}_a = \frac{1}{\text{vol}(R_a)} \sum_{j \in R_a} d_{row,j} \mathbf{r}_j$ is the weighted center of cluster R_a ($a = 1, \dots, k$). Similarly, the weighted k -variance of the k -dimensional column representatives is

$$S_k^2(\mathbf{Y}) = \min_{(C_1, \dots, C_k)} \sum_{a=1}^k \sum_{j \in C_a} d_{col,j} \|\mathbf{c}_j - \bar{\mathbf{c}}_a\|^2, \quad (7)$$

where $\bar{\mathbf{c}}_a = \frac{1}{\text{vol}(C_a)} \sum_{j \in C_a} d_{col,j} \mathbf{c}_j$ is the weighted center of cluster C_a ($a = 1, \dots, k$). Observe, that the trivial vector components can be omitted, and the k -variance of the so obtained $(k-1)$ -dimensional representatives will be the same.

Definition 6 *The cut-norm of the rectangular real matrix \mathbf{A} with row-set Row and column-set Col is*

$$\|\mathbf{A}\|_{\square} = \max_{R \subset \text{Row}, C \subset \text{Col}} \left| \sum_{i \in R} \sum_{j \in C} a_{ij} \right|.$$

Lemma 7 *For the cut-norm of the $n \times m$ real matrix \mathbf{A}*

$$\|\mathbf{A}\|_{\square} \leq \sqrt{nm} \|\mathbf{A}\|$$

holds, where the right hand side contains its spectral norm, i.e., the largest singular value of \mathbf{A} .

PROOF.

$$\begin{aligned} \|\mathbf{A}\|_{\square} &= \max_{\mathbf{x} \in \{0,1\}^n, \mathbf{y} \in \{0,1\}^m} |\mathbf{x}^T \mathbf{A} \mathbf{y}| = \max_{\mathbf{x} \in \{0,1\}^n, \mathbf{y} \in \{0,1\}^m} \left| \left(\frac{\mathbf{x}}{\|\mathbf{x}\|} \right)^T \mathbf{A} \left(\frac{\mathbf{y}}{\|\mathbf{y}\|} \right) \right| \cdot \|\mathbf{x}\| \cdot \|\mathbf{y}\| \\ &\leq \sqrt{nm} \max_{\|\mathbf{x}\|=1, \|\mathbf{y}\|=1} |\mathbf{x}^T \mathbf{A} \mathbf{y}| = \sqrt{nm} \|\mathbf{A}\|, \end{aligned}$$

since for $\mathbf{x} \in \{0, 1\}^n$, $\|\mathbf{x}\| \leq \sqrt{n}$, and for $\mathbf{y} \in \{0, 1\}^m$, $\|\mathbf{y}\| \leq \sqrt{m}$.

Theorem 8 *Let \mathbf{C} be a non-decomposable contingency table of n -element row set Row and m -element column set Col , with row- and column sums $d_{row,1}, \dots, d_{row,n}$ and $d_{col,1}, \dots, d_{col,m}$, respectively. Suppose that $\sum_{i=1}^n \sum_{j=1}^m c_{ij} = 1$ and there are no dominant rows and columns: $d_{row,i} = \Theta(1/n)$, $(i = 1, \dots, n)$ and $d_{col,j} = \Theta(1/m)$, $(j = 1, \dots, m)$ as $n, m \rightarrow \infty$. Let the singular values of \mathbf{C}_{corr} be*

$$1 = s_1 > s_2 \geq \dots \geq s_k > \varepsilon \geq s_i, \quad i \geq k+1.$$

The partition (R_1, \dots, R_k) of Row and (C_1, \dots, C_k) of Col are defined so that they minimize the weighted k -variances $S_k^2(\mathbf{X})$ and $S_k^2(\mathbf{Y})$ of the row and column representatives defined in (6) and (7), respectively. Suppose that there are constants $0 < K_1, K_2 \leq \frac{1}{k}$ such that $|R_i| \geq K_1 n$ and $|C_i| \geq K_2 m$ ($i = 1, \dots, k$), respectively. Then the R_i, C_j pairs are $\mathcal{O}(\sqrt{2k}(S_k(\mathbf{X})S_k(\mathbf{Y}) + \varepsilon))$ -volume regular ($i, j = 1, \dots, k$).

PROOF. Recall that provided \mathbf{C} is non-decomposable, the largest singular value $s_1 = 1$ of \mathbf{C}_{corr} is single with corresponding singular vector pair $\mathbf{v}_1 = \mathbf{D}_{row}^{1/2} \mathbf{1}$ and $\mathbf{u}_1 = \mathbf{D}_{col}^{1/2} \mathbf{1}$ with the constantly $\mathbf{1}$ vectors of appropriate size. The optimal k -dimensional representatives of the rows and columns are row vectors of the matrices $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_k)$ and $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_k)$, where $\mathbf{x}_i = \mathbf{D}_{row}^{-1/2} \mathbf{v}_i$ and $\mathbf{y}_i = \mathbf{D}_{col}^{-1/2} \mathbf{u}_i$, respectively ($i = 1, \dots, k$). Suppose that the minimum k -variance is attained on the k -partition (R_1, \dots, R_k) of the rows and (C_1, \dots, C_k) of the columns. By an easy analysis of variance argument of [3] it follows that

$$S_k^2(\mathbf{X}) = \sum_{i=1}^k \text{dist}^2(\mathbf{v}_i, F), \quad S_k^2(\mathbf{Y}) = \sum_{i=1}^k \text{dist}^2(\mathbf{u}_i, G),$$

where $F = \text{Span}\{\mathbf{D}_{row}^{1/2} \mathbf{w}_1, \dots, \mathbf{D}_{row}^{1/2} \mathbf{w}_k\}$ and $G = \text{Span}\{\mathbf{D}_{col}^{1/2} \mathbf{z}_1, \dots, \mathbf{D}_{col}^{1/2} \mathbf{z}_k\}$ with the so-called normalized row partition vectors $\mathbf{w}_1, \dots, \mathbf{w}_k$ of coordinates $w_{ji} = \frac{1}{\sqrt{\text{vol}(R_i)}}$ if $j \in R_i$ and 0, otherwise, and column partition vectors $\mathbf{z}_1, \dots, \mathbf{z}_k$ of coordinates $z_{ji} = \frac{1}{\sqrt{\text{vol}(C_i)}}$ if $j \in C_i$ and 0, otherwise ($i = 1, \dots, k$).

Note that the vectors $\mathbf{D}_{row}^{1/2} \mathbf{w}_1, \dots, \mathbf{D}_{row}^{1/2} \mathbf{w}_k$ and $\mathbf{D}_{col}^{1/2} \mathbf{z}_1, \dots, \mathbf{D}_{col}^{1/2} \mathbf{z}_k$ form orthonormal systems in \mathbb{R}^n and \mathbb{R}^m , respectively (but they are, usually, not complete). By [3], we can find orthonormal systems $\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_k \in F$ and $\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_k \in G$ such that

$$S_k^2(\mathbf{X}) \leq \sum_{i=1}^k \|\mathbf{v}_i - \tilde{\mathbf{v}}_i\|^2 \leq 2S_k^2(\mathbf{X}), \quad S_k^2(\mathbf{Y}) \leq \sum_{i=1}^k \|\mathbf{u}_i - \tilde{\mathbf{u}}_i\|^2 \leq 2S_k^2(\mathbf{Y}).$$

We approximate the matrix $\mathbf{C}_{corr} = \sum_{i=1}^r s_i \mathbf{v}_i \mathbf{u}_i^T$ by the rank k matrix $\sum_{i=1}^k s_i \tilde{\mathbf{v}}_i \tilde{\mathbf{u}}_i^T$ with the following accuracy (in spectral norm):

$$\left\| \sum_{i=1}^r s_i \mathbf{v}_i \mathbf{u}_i^T - \sum_{i=1}^k s_i \tilde{\mathbf{v}}_i \tilde{\mathbf{u}}_i^T \right\| \leq \sum_{i=1}^k s_i \left\| \mathbf{v}_i \mathbf{u}_i^T - \tilde{\mathbf{v}}_i \tilde{\mathbf{u}}_i^T \right\| + \left\| \sum_{i=k+1}^r s_i \mathbf{v}_i \mathbf{u}_i^T \right\|, \quad (8)$$

where the spectral norm of the last term is at most ε , and the individual terms of the first one are estimated from above in the following way.

$$\begin{aligned} s_i \left\| \mathbf{v}_i \mathbf{u}_i^T - \tilde{\mathbf{v}}_i \tilde{\mathbf{u}}_i^T \right\| &\leq \left\| (\mathbf{v}_i \mathbf{u}_i^T - \tilde{\mathbf{v}}_i \mathbf{u}_i^T) + (\tilde{\mathbf{v}}_i \mathbf{u}_i^T - \tilde{\mathbf{v}}_i \tilde{\mathbf{u}}_i^T) \right\| \\ &\leq \left\| (\mathbf{v}_i - \tilde{\mathbf{v}}_i) \mathbf{u}_i^T \right\| + \left\| \tilde{\mathbf{v}}_i (\mathbf{u}_i - \tilde{\mathbf{u}}_i)^T \right\| \\ &= \sqrt{\left\| (\mathbf{v}_i - \tilde{\mathbf{v}}_i) \mathbf{u}_i^T \mathbf{u}_i (\mathbf{v}_i - \tilde{\mathbf{v}}_i)^T \right\|} + \sqrt{\left\| (\mathbf{u}_i - \tilde{\mathbf{u}}_i) \tilde{\mathbf{v}}_i^T \tilde{\mathbf{v}}_i (\mathbf{u}_i - \tilde{\mathbf{u}}_i)^T \right\|} \\ &= \sqrt{(\mathbf{v}_i - \tilde{\mathbf{v}}_i)^T (\mathbf{v}_i - \tilde{\mathbf{v}}_i)} + \sqrt{(\mathbf{u}_i - \tilde{\mathbf{u}}_i)^T (\mathbf{u}_i - \tilde{\mathbf{u}}_i)} \\ &= \left\| \mathbf{v}_i - \tilde{\mathbf{v}}_i \right\| + \left\| \mathbf{u}_i - \tilde{\mathbf{u}}_i \right\|, \end{aligned}$$

where we exploited that the spectral norm (i.e., the largest singular value) of an $n \times m$ matrix \mathbf{A} is equal to either the squareroot of the largest eigenvalue of the matrix $\mathbf{A} \mathbf{A}^T$ or equivalently, that of $\mathbf{A}^T \mathbf{A}$. In the above calculations all of these matrices are of rank 1, hence, the largest eigenvalue of the symmetric, positive semidefinite matrix under the squareroot is the only non-zero eigenvalue of it, therefore, it is equal to its trace; finally, we used the commutativity of the trace, and in the last line we have the usual vector norm.

Therefore the first term in (8) can be estimated from above by

$$\begin{aligned} \sum_{i=1}^k \left\| \mathbf{v}_i \mathbf{u}_i^T - \tilde{\mathbf{v}}_i \tilde{\mathbf{u}}_i^T \right\| &\leq \sqrt{k} \sqrt{\sum_{i=1}^k \left\| \mathbf{v}_i - \tilde{\mathbf{v}}_i \right\|^2} + \sqrt{k} \sqrt{\sum_{i=1}^k \left\| \mathbf{u}_i - \tilde{\mathbf{u}}_i \right\|^2} \\ &\leq \sqrt{k} (\sqrt{2S_k^2(\mathbf{X})} + \sqrt{2S_k^2(\mathbf{Y})}) = \sqrt{2k} (S_k(\mathbf{X}) + S_k(\mathbf{Y})). \end{aligned}$$

Based on these considerations and relation between the cut norm and the spectral norm (see Lemma 7), the densities to be estimated in the defining formula (5) of volume regularity can be written in terms of stepwise constant vectors in the following way. The vectors $\hat{\mathbf{v}}_i := \mathbf{D}_{row}^{-1/2} \tilde{\mathbf{v}}_i$ are stepwise constants on the partition (R_1, \dots, R_k) of the rows, whereas the vectors $\hat{\mathbf{u}}_i := \mathbf{D}_{col}^{-1/2} \tilde{\mathbf{u}}_i$ are stepwise constants on the partition (C_1, \dots, C_k) of the columns, $i = 1, \dots, k$. The matrix

$$\sum_{i=1}^k s_i \hat{\mathbf{v}}_i \hat{\mathbf{u}}_i^T$$

is therefore an $n \times m$ block-matrix on $k \times k$ blocks belonging to the above partition of the rows and columns. Let \hat{c}_{ab} denote its entries in the a, b block ($a, b = 1, \dots, k$). Using (8), the rank k approximation of the matrix \mathbf{C} is

performed with the following accuracy of the perturbation \mathbf{E} in spectral norm:

$$\|\mathbf{E}\| = \left\| \mathbf{C} - \mathbf{D}_{row} \left(\sum_{i=1}^k s_i \hat{\mathbf{v}}_i \hat{\mathbf{u}}_i^T \right) \mathbf{D}_{col} \right\| = \left\| \mathbf{D}_{row}^{1/2} (\mathbf{C}_{corr} - \sum_{i=1}^k s_i \mathbf{v}_i \mathbf{u}_i^T) \mathbf{D}_{col}^{1/2} \right\|.$$

Therefore, the entries of \mathbf{C} – for $i \in R_a, j \in C_b$ – can be decomposed as

$$c_{ij} = d_{row,i} d_{col,j} \hat{c}_{ab} + \eta_{ij},$$

where the cut norm of the $n \times m$ error matrix $\mathbf{E} = (\eta_{ij})$ restricted to $R_a \times C_b$ (otherwise it contains entries all zeroes) and denoted by \mathbf{E}_{ab} , is estimated as follows:

$$\begin{aligned} \|\mathbf{E}_{ab}\|_{\square} &\leq \sqrt{nm} \|\mathbf{E}_{ab}\| \leq \sqrt{nm} \cdot \|\mathbf{D}_{row,a}^{1/2}\| \cdot (\sqrt{2k}(S_k(\mathbf{X}) + S_k(\mathbf{Y})) + \varepsilon) \cdot \|\mathbf{D}_{col,b}^{1/2}\| \\ &\leq \sqrt{nm} \sqrt{c_1 \frac{\text{Vol}(R_a)}{|R_a|}} \cdot \sqrt{c_2 \frac{\text{Vol}(C_b)}{|C_b|}} (\sqrt{2k}(S_k(\mathbf{X}) + S_k(\mathbf{Y})) + \varepsilon) \\ &= \sqrt{c_1 c_2} \cdot \sqrt{\frac{n}{|R_a|}} \cdot \sqrt{\frac{m}{|C_b|}} \cdot \sqrt{\text{Vol}(R_a)} \sqrt{\text{Vol}(C_b)} (\sqrt{2k}(S_k(\mathbf{X}) + S_k(\mathbf{Y})) + \varepsilon) \\ &\leq \sqrt{\frac{c_1 c_2}{K_1 K_2}} \sqrt{\text{Vol}(R_a)} \sqrt{\text{Vol}(C_b)} (\sqrt{2k} s + \varepsilon) \\ &= c \sqrt{\text{Vol}(R_a)} \sqrt{\text{Vol}(C_b)} (\sqrt{2k}(S_k(\mathbf{X}) + S_k(\mathbf{Y})) + \varepsilon), \end{aligned}$$

where the $n \times n$ diagonal matrix $\mathbf{D}_{row,a}$ inherits \mathbf{D}_{row} 's diagonal entries over R_a , whereas the $m \times m$ diagonal matrix $\mathbf{D}_{col,b}$ inherits \mathbf{D}_{col} 's diagonal entries over C_b , otherwise they are zeros. Further, the constants c_1, c_2 are due to the fact that there are no dominant rows and columns, while K_1, K_2 are derived from the cluster size balancing conditions. Hence, the constant c does not depend on n and m . Consequently, for $a, b = 1, \dots, k$ and $X \subset R_a, Y \subset C_b$:

$$\begin{aligned} |c(X, Y) - \rho(R_a, C_b) \text{Vol}(X) \text{Vol}(Y)| &= \\ \left| \sum_{i \in X} \sum_{j \in Y} (d_{row,i} d_{col,j} \hat{c}_{ab} + \eta_{ij}^{ab}) - \frac{\text{Vol}(X) \text{Vol}(Y)}{\text{Vol}(R_a) \text{Vol}(C_b)} \sum_{i \in R_a} \sum_{j \in C_b} (d_{row,i} d_{col,j} \hat{c}_{ab} + \eta_{ij}^{ab}) \right| &= \\ \left| \sum_{i \in X} \sum_{j \in Y} \eta_{ij}^{ab} - \frac{\text{Vol}(X) \text{Vol}(Y)}{\text{Vol}(R_a) \text{Vol}(C_b)} \sum_{i \in R_a} \sum_{j \in C_b} \eta_{ij}^{ab} \right| &\leq 2 \|\mathbf{E}_{ab}\|_{\square} \\ &\leq 2c (\sqrt{2k}(S_k(\mathbf{X}) + S_k(\mathbf{Y})) + \varepsilon) \sqrt{\text{Vol}(R_a) \text{Vol}(C_b)}, \end{aligned}$$

that gives the required statement for $a, b = 1, \dots, k$.

Note that when we use Definition 5 of γ -volume regularity for the row-column cluster pairs R_i, C_j ($i, j = 1, \dots, k$), then we may say that the k -way *discrepancy* of the underlying contingency table is the minimum γ for which all the row-column cluster pairs are γ -volume regular. With this nomenclature,

Theorem 8 states that the k -way discrepancy of a contingency table can be estimated from above by the $(k + 1)$ th largest singular value of the correspondence matrix and the k -variance of the clusters obtained by the left and right singular vectors corresponding to the k largest singular values of this matrix. Hence, SVD based representation is applicable to find volume regular cluster pairs for given k , where k is the number of structural (protruding) singular values.

5 Discussion, application, and extension to directed graphs

In the ideal k -cluster case, we consider the following generalized random binary contingency table model: given the partition (R_1, \dots, R_k) of the rows and (C_1, \dots, C_k) of the columns, the entry in the row $i \in R_a$ and column $j \in C_b$ is 1 with probability p_{ab} , and 0 otherwise, independently of other rows of R_a and columns of C_b , $1 \leq a, b \leq k$. We can think of the probability p_{ab} as the inter-cluster density of the row-column cluster pair R_a, C_b . Since generalized contingency tables can be viewed as block-matrices (with $k \times k$ blocks) burdened with a general random noise, in [4], we gave the following spectral characterization of them. Fixing k , and tending with n and m to infinity in such a way that the cluster sizes grow at the same rate and also n and m subpolynomially, there exists a positive number $\theta \leq 1$, independent of n and m , such that for every $0 < \tau < 1/2$ there are exactly k singular values of \mathbf{C}_{corr} greater than $\theta - \max\{n^{-\tau}, m^{-\tau}\}$, while all the others are at most $\max\{n^{-\tau}, m^{-\tau}\}$; further, the weighted k -variance of the row and column representatives constructed by the k transformed structural left and right singular vectors is $\mathcal{O}(\max\{n^{-\tau}, m^{-\tau}\})$, respectively.

For general contingency tables, our result is that the existence of k singular values of \mathbf{C}_{corr} , separated from 0 by ε , is indication of a k -cluster structure, while the eigenvalues accumulating around 0 are responsible for the pairwise regularities. The clusters themselves can be recovered by applying the k -means algorithm for the row and column representatives obtained by the left and right singular vectors corresponding to the structural singular values.

We applied the biclustering algorithm to find simultaneously clusters of stores and products based on their consumption in TESCO stores. Figure 1 shows 3 clusters of the stores in which the consumption of the products belonging to the same cluster was homogeneous with consumption-density $\frac{c(R_a, C_b)}{\text{Vol}(R_a)\text{Vol}(C_b)}$ between store-cluster R_a and product-cluster C_b ($a, b = 1, \dots, 3$). After sorting the rows and columns according to their cluster memberships, we plotted the entries $\frac{c_{ij}}{d_{row,i}d_{col,j}}$ (there was one exceptional store-cluster which contained only 3 stores, but the others could be identified with groups of smaller and larger stores associated with product groups of high consumption-density

within them).

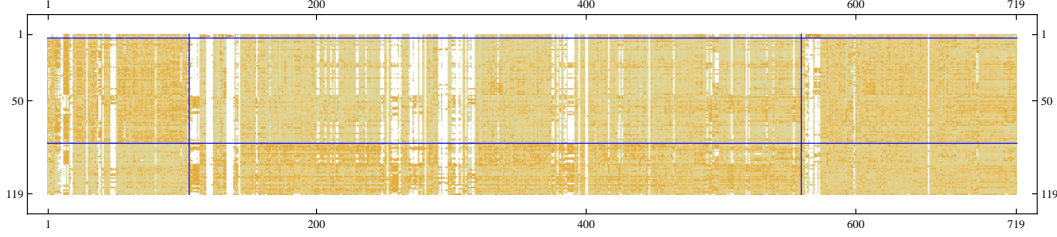


Fig. 1. Result of biclustering 119 stores and 719 products into 3 clusters

We can consider quadratic, but not symmetric contingency tables with zero diagonal as edge-weight matrices of directed graphs. The $n \times n$ edge-weight matrix \mathbf{W} of a directed graph has zero diagonal, but is usually not symmetric: w_{ij} is the weight of the $i \rightarrow j$ edge ($i, j = 1, \dots, n; i \neq j$). In this setup, the generalized in- and out-degrees are

$$d_{out,i} = \sum_{j=1}^n w_{ij} \quad (i = 1, \dots, n) \quad \text{and} \quad d_{in,j} = \sum_{i=1}^n w_{ij} \quad (j = 1, \dots, n);$$

further, $\mathbf{D}_{in} = \text{diag}(d_{in,1}, \dots, d_{in,n})$ and $\mathbf{D}_{out} = \text{diag}(d_{out,1}, \dots, d_{out,n})$ are the in- and out-degree matrices. Suppose that there are no sources and sinks (i.e. no zero out- and in-degrees), further, that \mathbf{W} is non-decomposable. Then the correspondence matrix belonging to \mathbf{W} is

$$\mathbf{W}_{corr} = \mathbf{D}_{out}^{-1/2} \mathbf{W} \mathbf{D}_{in}^{-1/2},$$

and its SVD is used to minimize the normalized two-way cut of \mathbf{W} as a contingency table, see Section 3. Butler [7] generalized the Expander Mixing Lemma for this situation. We can further generalize it to obtain regular in- and out-vertex cluster pairs, for a given k , in the following sense. The V_{in}, V_{out} in- and out-vertex cluster pair of the directed graph (with sum of the weights of directed edges 1) is γ -volume regular if for all $X \subset V_{out}$ and $Y \subset V_{in}$ the relation

$$|w(X, Y) - \rho(V_{out}, V_{in}) \text{Vol}_{out}(X) \text{Vol}_{in}(Y)| \leq \gamma \sqrt{\text{Vol}_{out}(V_{out}) \text{Vol}_{in}(V_{in})}$$

holds, where the *directed cut* $w(X, Y)$ is the sum the weights of the $X \rightarrow Y$ edges, $\text{Vol}_{out}(X) = \sum_{i \in X} d_{out,i}$, $\text{Vol}_{in}(Y) = \sum_{j \in Y} d_{in,j}$, and $\rho(V_{out}, V_{in}) = \frac{w(V_{out}, V_{in})}{\text{Vol}_{out}(V_{out}) \text{Vol}_{in}(V_{in})}$ is the relative inter-cluster density of the out-in cluster pair V_{out}, V_{in} . The clustering $(V_{in,1}, \dots, V_{in,k})$ and $(V_{out,1}, \dots, V_{out,k})$ of the columns and rows – guaranteed by Theorem 8 – corresponds to in- and out-clusters of the same vertex set such that the directed information flow $V_{out,a} \rightarrow V_{in,b}$ is as homogeneous as possible for all $a, b = 1, \dots, k$ pairs.

Acknowledgements

We are indebted to the Tesco Hungary for making their data available and Tamás Kóí for computer processing the data.

References

- [1] Bhatia, R., Matrix Analysis, Springer (1996).
- [2] Bilu, Y. and Linial, N., Lifts, discrepancy and nearly optimal spectral gap, *Combinatorica* **26** (2006), 495–519.
- [3] Bolla, M., Tusnády, G., Spectra and optimal partitions of weighted graphs, *Discrete Mathematics* **128** (1994), 1–20.
- [4] Bolla, M., Friedl, K., Krámlí, A., Singular value decomposition of large random matrices (for two-way classification of microarrays), *Journal of Multivariate Analysis* **101** (2010), 434–446.
- [5] Bolla, M., Spectra and structure of weighted graphs, *Electronic Notes in Discrete Mathematics* **38** (2011), 149–154.
- [6] Bollobás, B., Nikiforov, V., Hermitian matrices and graphs: singular values and discrepancy, *Discrete Mathematics* **285** (2004), 17–32.
- [7] Butler, S., Using discrepancy to control singular values for nonnegative matrices, *Lin. Alg. Appl.* **419** (2006), 486–493.
- [8] Dhillon, I. S., Co-clustering documents and words using bipartite spectral graph partitioning. In: Proc. ACM Int’l Conf. Knowledge Disc. Data Mining (KDD 2001), 2001.
- [9] Ding, C., He, X., Zha, H., Gu, M., Simon, H. D., A minmax cut spectral method for data clustering and data partitioning, Lawrence Berkeley National Laboratory Tech. Rep. 54111, 2003.
- [10] Frieze, A., Kannan, R., Vempala, S., Fast Monte-Carlo Algorithms for finding low-rank approximations. In: Proc. of the 39th Annual IEEE Symposium on Foundations of Computer Science (FOCS), pp. 370–386, 1998.
- [11] Kleinberg, J., Authoritative sources in hyperlinked environment, IBM Research Report RJ 10076 (91892), 1997.
- [12] Kluger, Y., Basri, R., Chang, J. T., Gerstein, M., Spectral biclustering of microarray data: clustering genes and conditions, *Genome Research* **13** (2003), 703–716.
- [13] Rao, C. R., Separation theorems for singular values of matrices and their applications in multivariate analysis, *J. Multivariate Analysis* **9** (1979), 362–377.