

Data Mining and Matrices

05 – Semi-Discrete Decomposition

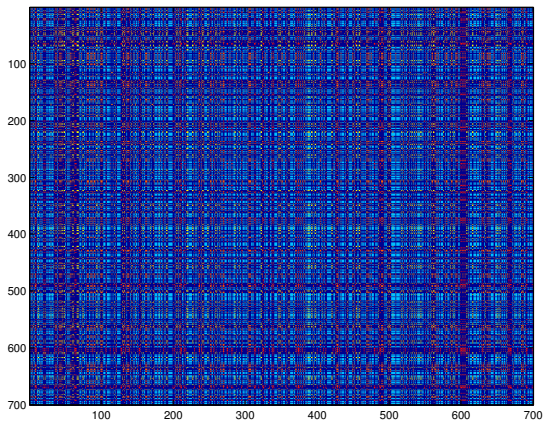
Rainer Gemulla, Pauli Miettinen

May 16, 2013

Outline

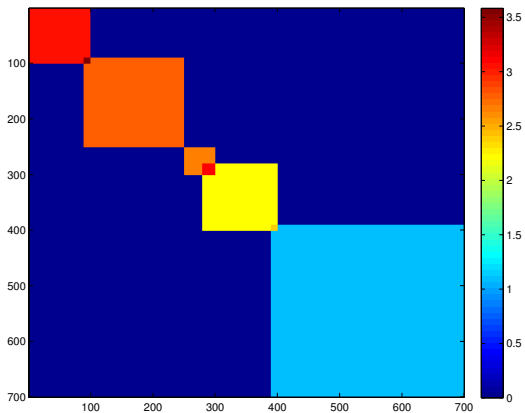
- 1 Hunting the Bump
- 2 Semi-Discrete Decomposition
- 3 The Algorithm
- 4 Applications
 - SDD alone
 - SVD + SDD
- 5 Wrap-Up

An example data



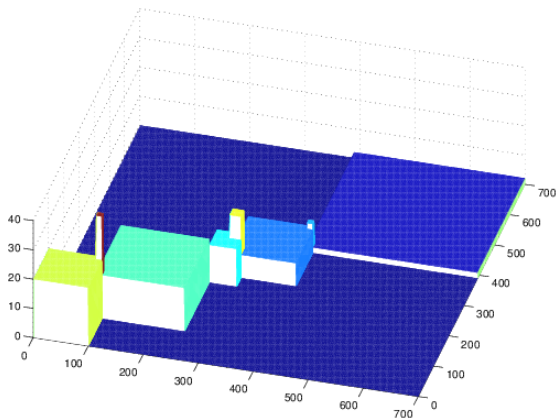
The data

An example data



The data after permuting rows and columns

An example data



The data in a 3D view

Can we find the bumps in the picture automatically (from unpermuted data)?

What is a bump?

$$\mathbf{A} = \begin{pmatrix} 3 & 1 & 3 \\ 2 & 3 & 1 \\ 3 & 2 & 3 \end{pmatrix}$$

$$I = \{1, 3\}$$

$$J = \{1, 3\}$$

$$\mathbf{x} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$$

$$\mathbf{A} \circ \mathbf{xy}^T = \begin{pmatrix} 3 & 0 & 3 \\ 0 & 0 & 0 \\ 3 & 0 & 3 \end{pmatrix}$$

- A **submatrix** of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ contains some rows of \mathbf{A} and some columns of those rows
 - ▶ Let $I \subseteq \{1, 2, \dots, m\}$ have the row indices and $J \subseteq \{1, 2, \dots, n\}$ have the column indices of the submatrix
 - ▶ If $\mathbf{x} \in \{0, 1\}^m$ has $x_i = 1$ iff $i \in I$ and $\mathbf{y} \in \{0, 1\}^n$ has $y_j = 1$ iff $j \in J$, then $\mathbf{xy}^T \in \{0, 1\}^{m \times n}$ has $(\mathbf{xy}^T)_{ij} = 1$ iff a_{ij} is in the submatrix
 - ▶ $\mathbf{A} \circ \mathbf{xy}^T$ has the values of the submatrix and zeros elsewhere
 - ★ $(\mathbf{A} \circ \mathbf{B})_{ij} = a_{ij}b_{ij}$ is the **Hadamard matrix product**
- The submatrix is uniform if all (or most) of its values are (approximately) the same
 - ▶ Exactly uniform submatrices with value δ can be written as $\delta \mathbf{xy}^T$ — a bump

The next bump and negative values

- Assume we know how to find the largest bump of a matrix
- To find another bump, we can subtract the found bump from the matrix and find the largest bump of the residual matrix
 - ▶ But after subtraction we might have negative values in the matrix
- We can generalize the uniform submatrices to require uniformity only in magnitude
 - ▶ Allow characteristic vectors \mathbf{x} and \mathbf{y} to take values from $\{-1, 0, 1\}$
 - ▶ If $\mathbf{x} = (-1, 0, -1)^T$ and $\mathbf{y} = (1, 0, -1)^T$, then

$$\delta \mathbf{xy}^T = \begin{pmatrix} -\delta & 0 & \delta \\ 0 & 0 & 0 \\ \delta & 0 & -\delta \end{pmatrix}$$

- This allows us to define bumps in matrices with negative values

Outline

- 1 Hunting the Bump
- 2 Semi-Discrete Decomposition**
- 3 The Algorithm
- 4 Applications
 - SDD alone
 - SVD + SDD
- 5 Wrap-Up

The definition

Semi-Discrete Decomposition

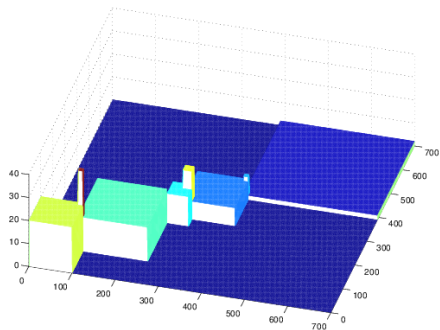
Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, the *semi-discrete decomposition* (SDD) of \mathbf{A} of dimension k is

$$\mathbf{A} \approx \mathbf{X}_k \mathbf{D}_k \mathbf{Y}_k^T,$$

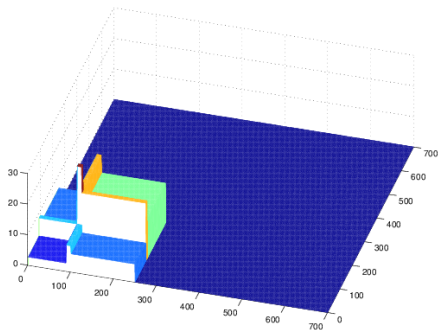
where

- $\mathbf{X}_k \in \{-1, 0, 1\}^{m \times k}$
- $\mathbf{Y}_k \in \{-1, 0, 1\}^{n \times k}$
- $\mathbf{D}_k \in \mathbb{R}_+^{k \times k}$ is a diagonal matrix

Example

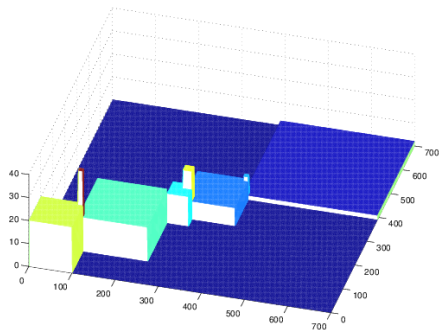


The data

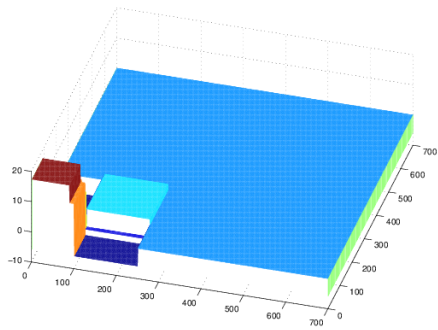


The first component $\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T$ using
SVD

Example



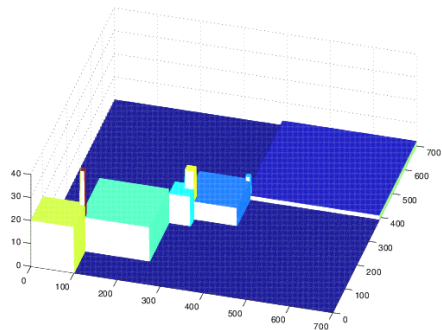
The data



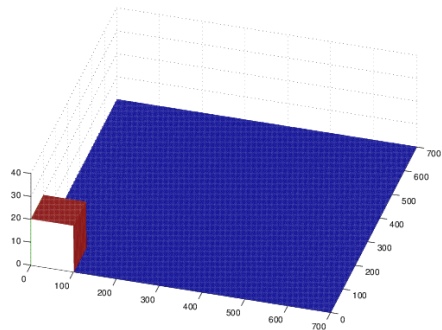
The second component $\sigma_2 \mathbf{u}_2 \mathbf{v}_2^T$
using SVD

The SVD cannot find the bumps

Example

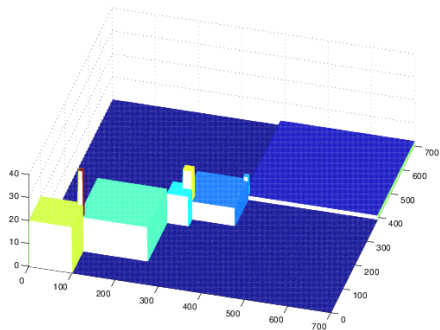


The data

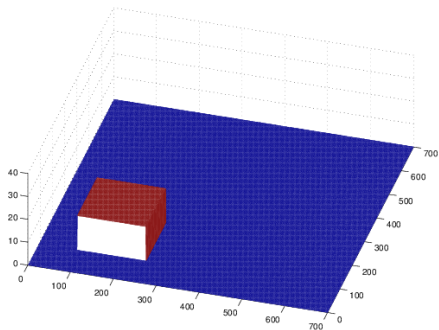


The first bump $d_1 \mathbf{x}_1 \mathbf{y}_1^T$ using SDD

Example

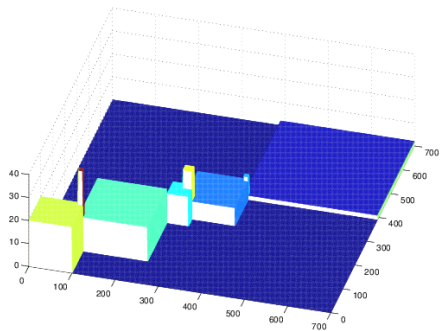


The data

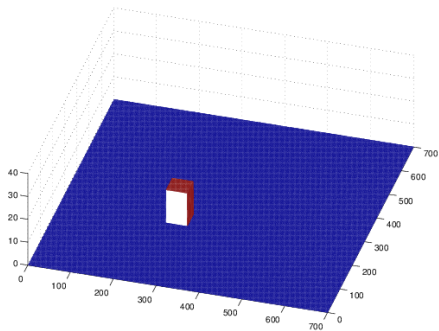


The second bump $d_2 \mathbf{x}_2 \mathbf{y}_2^T$ using SDD

Example

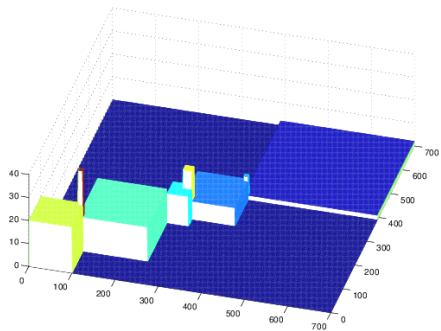


The data

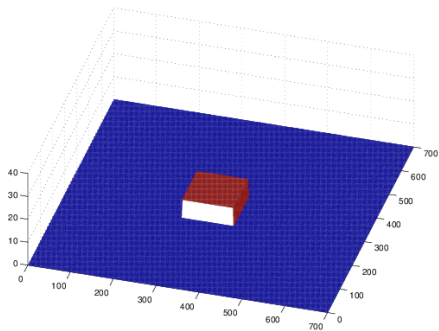


The third bump $d_3 \mathbf{x}_3 \mathbf{y}_3^T$ using SDD

Example

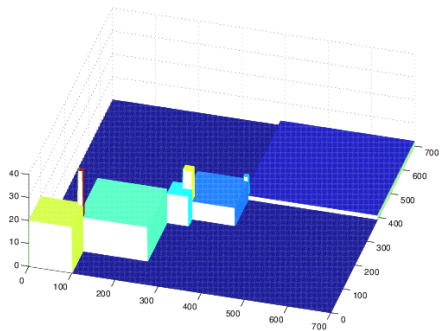


The data

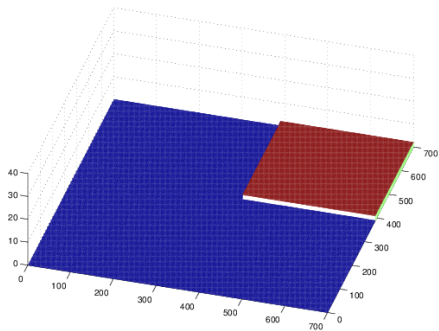


The fourth bump $d_4 \mathbf{x}_4 \mathbf{y}_4^T$ using SDD

Example

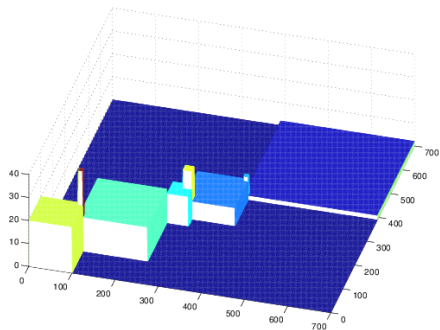


The data

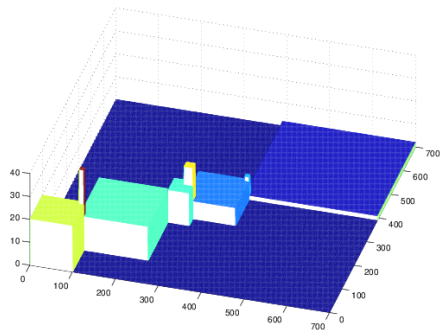


The fifth bump $d_5 \mathbf{x}_5 \mathbf{y}_5^T$ using SDD

Example



The data



The 5-dimensional SDD
approximation $\mathbf{X}_5 \mathbf{D}_5 \mathbf{Y}_5^T$

Properties of SDD

- The columns of \mathbf{X}_k and \mathbf{Y}_k do not need to be linearly independent
 - ▶ The same column can be even repeated multiple times
- The dimension k might need to be large for accurate approximation (compared to SVD)
 - ▶ $k = \min\{n, m\}$ is not necessarily enough for exact SDD
 - ★ $k = nm$ is always enough
 - ▶ First factors don't necessarily explain much about the matrix
- SDD factors are local
 - ▶ Only affect a certain submatrix, typically not every element
 - ▶ SVD factors typically change every value
- Storing an k -dimensional SDD takes less space than storing rank- k truncated SVD
 - ▶ \mathbf{X}_k and \mathbf{Y}_k are ternary and often sparse
- For every rank-1 layer of an SDD, all non-zero values in the layer have the same magnitude (d_{ij} for layer i)

Interpretation

- The factor interpretation is not very useful as the factors are not independent
 - ▶ A later factor can change just a subset of values already changed by an earlier factor
- The SDD can be interpreted as a form of bi-clustering
 - ▶ Every layer (bump) defines a group of rows and columns with homogeneous values in the residual matrix
- The component interpretation is natural to SDD
 - ▶ The SDD is a sum of local bumps
 - ▶ SDD doesn't model global phenomena (e.g. noise) well

Outline

- 1 Hunting the Bump
- 2 Semi-Discrete Decomposition
- 3 The Algorithm**
- 4 Applications
 - SDD alone
 - SVD + SDD
- 5 Wrap-Up

The outline of the algorithm

- 1 **Input:** Matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, non-negative integer k
- 2 **Output:** k -dimensional SDD of \mathbf{A} , i.e. matrices $\mathbf{X}_k \in \{-1, 0, 1\}^{m \times k}$, $\mathbf{Y}_k \in \{-1, 0, 1\}^{n \times k}$, and diagonal $\mathbf{D}_k \in \mathbb{R}_+^{k \times k}$
- 3 $\mathbf{R}_1 \leftarrow \mathbf{A}$
- 4 **for** $i = 1, \dots, k$
 - 1 Select $\mathbf{y}_i \in \{-1, 0, 1\}^n$
 - 2 **while** not converged
 - 1 Compute $\mathbf{x}_i \in \{-1, 0, 1\}^m$ given \mathbf{y}_i and \mathbf{R}_i
 - 2 Compute \mathbf{y}_i given \mathbf{x} and \mathbf{R}_i
 - 3 **end while**
 - 4 Set d_i to the average of $\mathbf{R}_i \circ \mathbf{x}_i \mathbf{y}_i^T$ over the non-zero locations of $\mathbf{x} \mathbf{y}^T$
 - 5 Set \mathbf{x}_i as the i th column of \mathbf{X}_i , \mathbf{y}_i the i th column of \mathbf{Y}_i , and d_i the i th value of \mathbf{D}_i
 - 6 $\mathbf{R}_{i+1} \leftarrow \mathbf{R}_i - d_i \mathbf{x}_i \mathbf{y}_i^T$
- 5 **end for**
- 6 **return** \mathbf{X}_k , \mathbf{Y}_k , and \mathbf{D}_k

Finding the bump

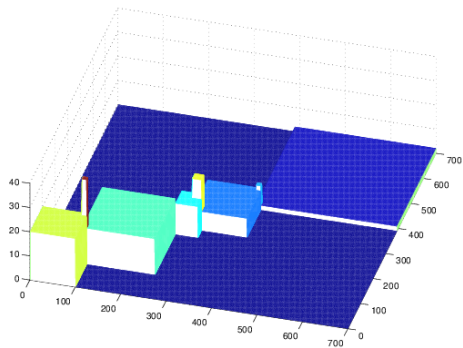
- **Problem:** Given $\mathbf{R} \in \mathbb{R}^{m \times n}$ and $\mathbf{y} \in \{-1, 0, 1\}^n$, find $\mathbf{x} \in \{-1, 0, 1\}^m$ such that $\|\mathbf{R} - d\mathbf{x}\mathbf{y}^T\|_F^2$ is minimized
 - ▶ We set $d \leftarrow \mathbf{x}^T \mathbf{R} \mathbf{y} / \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$ (the average of $\mathbf{R} \circ \mathbf{x}\mathbf{y}^T$ over the non-zero locations of $\mathbf{x}\mathbf{y}^T$)
 - ▶ We want to minimize the residual norm
- Set $\mathbf{s} \leftarrow \mathbf{R}\mathbf{y}$
- **Task:** Find \mathbf{x} that maximizes $F(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{s})^2 / \|\mathbf{x}\|_2^2$
 - ▶ Maximizing F equals minimizing the residual norm after d is set as above
 - ▶ Can be solved optimally by trying 2^m different binary vectors and setting the sign appropriately
- **Solution:** Order values s_i so that $|s_{i_1}| \geq |s_{i_2}| \geq \dots \geq |s_{i_m}|$ and set $x_{i_j} \leftarrow \text{sign}(s_{i_j})$ for the *first* J values s_i and 0 elsewhere
 - ▶ J is the number of nonzeros in \mathbf{x}
 - ★ Because we don't know J , we have to try every possibility and select the best
 - ▶ Values s_i contain the row sums of \mathbf{R} from those columns that are selected by \mathbf{y} and with sign set accordingly

Selecting the initial vector \mathbf{y}

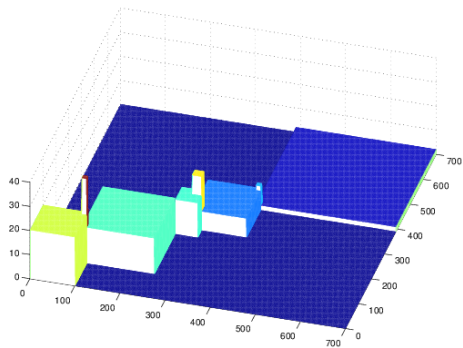
There are many ways to select the initial vector:

- MAX: set $y_j = 1$ for the column j that has the largest squared value of \mathbf{R} and rest to zero
 - ▶ Intuition: the very largest squared value is probably in the best bump
- CYC: set $y_j = 1$ for $j = (k \bmod n) + 1$
 - ▶ Cycle thru the columns
- THR: select a unit vector \mathbf{y} that satisfies $\|\mathbf{R}\mathbf{y}\|_F^2 \geq \|\mathbf{R}\|_F^2/n$
 - ▶ The selected column must have a squared sum that's above the average squared sum
 - ▶ The selection can be random or columns can be tried one-by-one
 - ★ The CYC and THR can be mixed

Example result

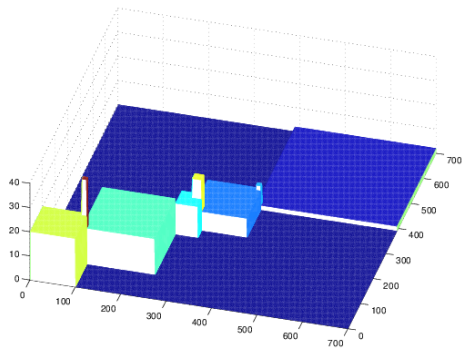


The data

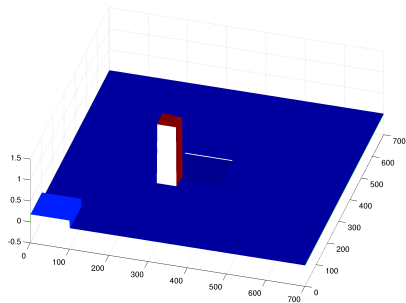


5-dimensional SDD

Example result



The data



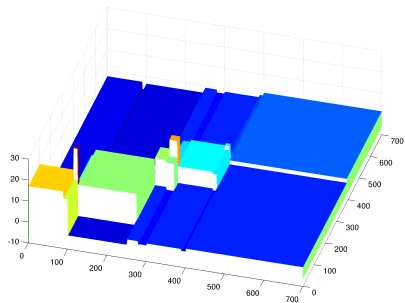
The matrix $\mathbf{X}_5 \mathbf{D}_5 \mathbf{Y}_5^T - \mathbf{A}$

Normalization

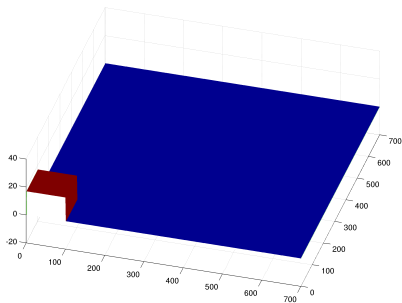
Normalization can have a profound effect on SDD

- Zero centering the columns will change the type of bumps found
 - ▶ The bumps in the original data have the largest-magnitude values
 - ▶ The bumps in the zero-centered data have the most extreme values
- Normalizing the variance will make the matrix to have more uniform values and thus changes the bumps
- Squaring the values will promote smaller bumps of exceptionally high values
- Square-rooting the values will promote larger bumps of smaller magnitude

Normalization example: zero-centered data

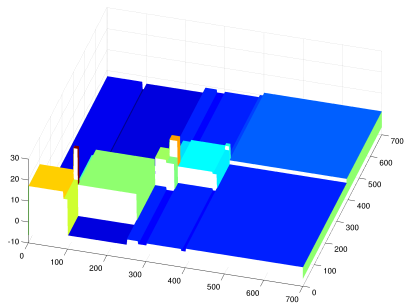


Zero-centered data

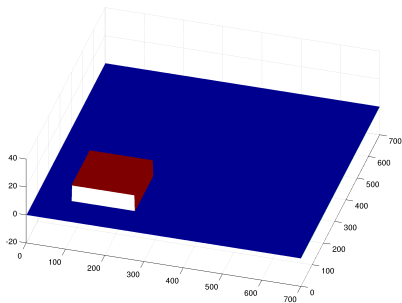


The first bump

Normalization example: zero-centered data

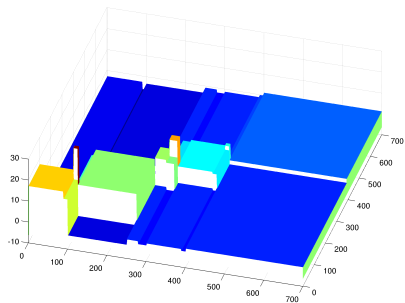


Zero-centered data

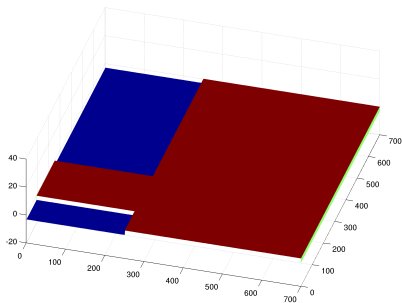


The second bump

Normalization example: zero-centered data



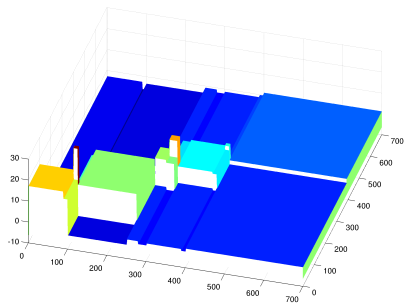
Zero-centered data



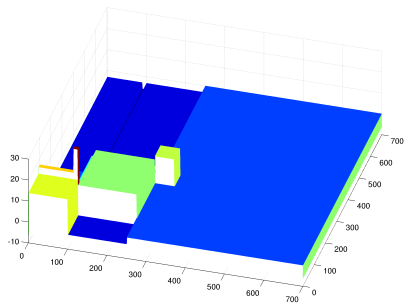
The third bump

Note that here red means 0

Normalization example: zero-centered data

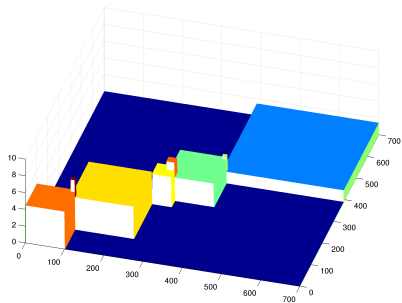


Zero-centered data

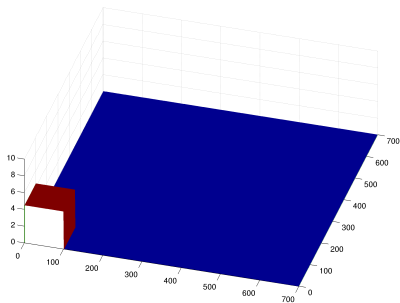


5-dimensional SDD

Normalization example: square-root of data

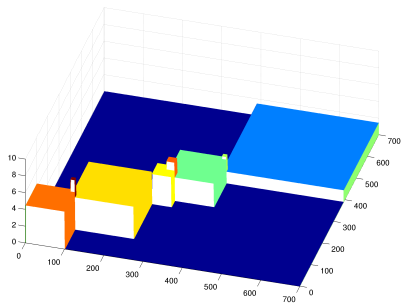


Data after taking element-wise square-root

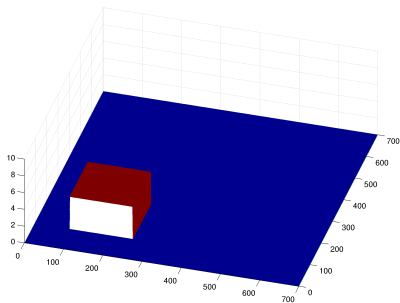


The first bump

Normalization example: square-root of data

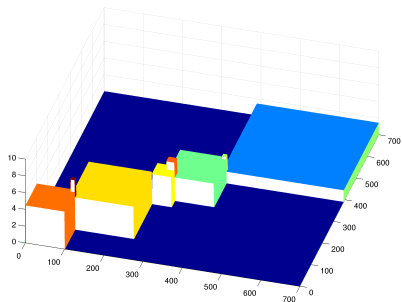


Data after taking element-wise square-root

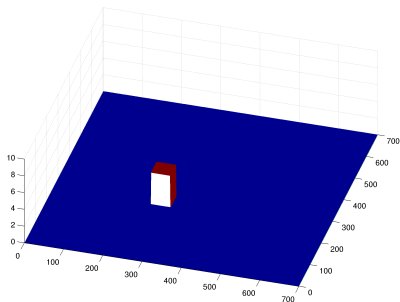


The second bump

Normalization example: square-root of data

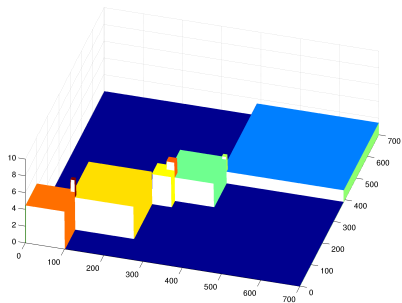


Data after taking element-wise square-root

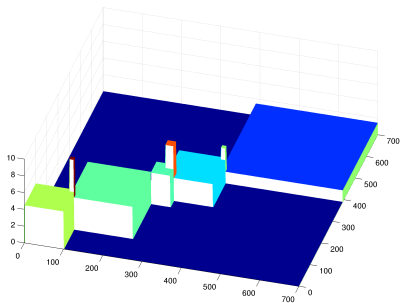


The third bump

Normalization example: square-root of data

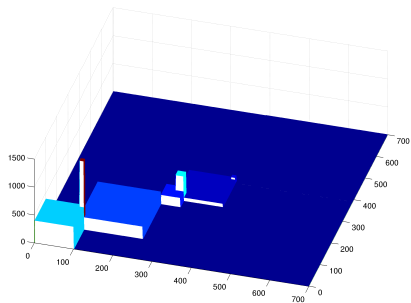


Data after taking element-wise square-root

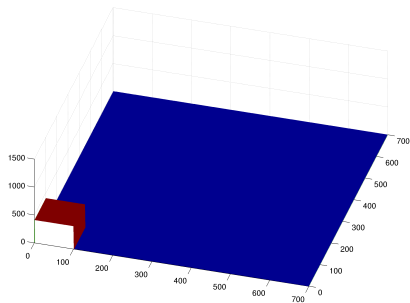


5-dimensional SDD

Normalization example: squared data

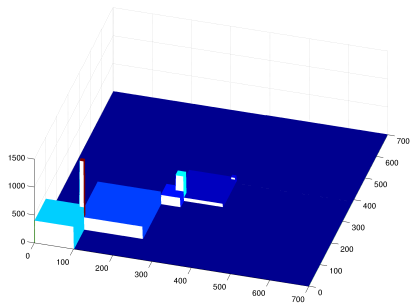


Squared data

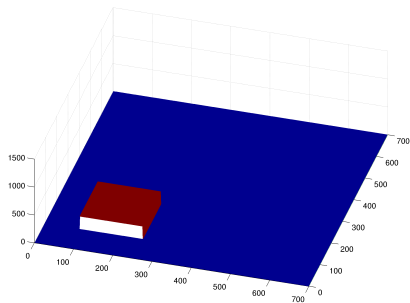


The first bump

Normalization example: squared data

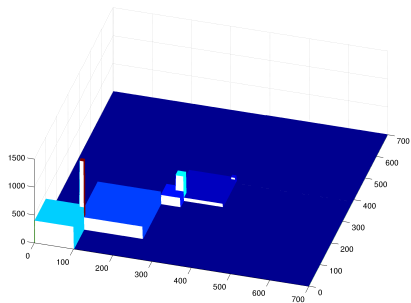


Squared data

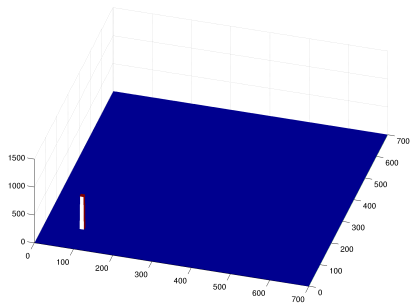


The second bump

Normalization example: squared data

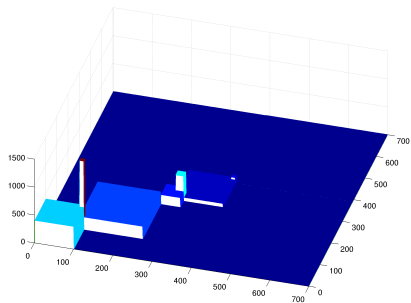


Squared data

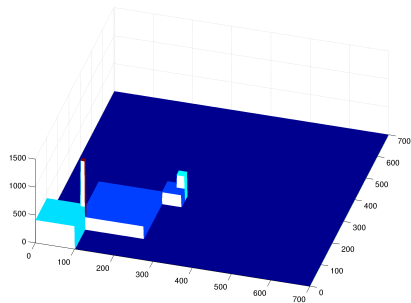


The third bump

Normalization example: squared data



Squared data



5-dimensional SDD

Outline

- 1 Hunting the Bump
- 2 Semi-Discrete Decomposition
- 3 The Algorithm
- 4 Applications**
 - SDD alone
 - SVD + SDD
- 5 Wrap-Up

Outline

- 1 Hunting the Bump
- 2 Semi-Discrete Decomposition
- 3 The Algorithm
- 4 Applications
 - SDD alone
 - SVD + SDD
- 5 Wrap-Up

Clustering

- SDD performs a type of bi-clustering of the matrix
 - ▶ Every bump $d\mathbf{x}\mathbf{y}^T$ gives a cluster of rows $x_i \neq 0$ and a cluster of columns $y_j \neq 0$, together with 'centroid' d
 - ▶ This is **not** a partition clustering: the clusters can overlap and not every row or column has to belong to some cluster
- We can impose an ordering of the bumps based on the values of d_i
 - ▶ The algorithm usually returns the bumps in that order
- This ordering can be used to obtain a hierarchical clustering
 - ▶ First column of \mathbf{X} clusters the rows of \mathbf{A} into three sets $(-1, 0, 1)$, and same for the first column of \mathbf{Y} and columns of \mathbf{A}
 - ▶ The second column of \mathbf{X} splits the previous clusters again into three sets
 - ★ Some of these sets can be empty
 - ▶ And so on and so forth
- Distance between two objects in the hierarchical clustering is not the usual dendrogram depth, but depends on whether we use the 0 or ± 1 branches

Other applications

- Image compression (O'Leary & Peleg, 1983)
 - ▶ A grayscale image is compressed using its SDD
 - ▶ The original application, modern image compression techniques are better
- Latent topic models (Kolda & O'Leary, 1998)
 - ▶ Used similarly as SVD is used to compute LSA
 - ▶ Compute the SDD of the term-document matrix
- SDD to correlation matrices
 - ▶ A bump in the correlation matrix \mathbf{AA}^T corresponds to rows of \mathbf{A} with similar values
 - ▶ A bump in $\mathbf{A}^T\mathbf{A}$ corresponds to columns with similar values

Outline

- 1 Hunting the Bump
- 2 Semi-Discrete Decomposition
- 3 The Algorithm
- 4 Applications
 - SDD alone
 - SVD + SDD
- 5 Wrap-Up

General approaches

- Most common way to combine SVD and SDD is to first use SVD to denoise the data and then to compute the SDD on the clean data
 - ▶ clean data = truncated SVD (\mathbf{A}_k)
 - ▶ SVD is good at finding global structure, SDD at finding local structure
- Another option is to first compute the \mathbf{A}_k with SVD and then apply SDD to $\mathbf{A}_k \mathbf{A}_k^T$ (or $\mathbf{A}_k^T \mathbf{A}_k$)
 - ▶ SDD finds the objects with similar values
- The results can be visualized using the first 2–3 columns of \mathbf{U} from SVD and the first layers of the hierarchical clustering from SDD

Classifying galaxies

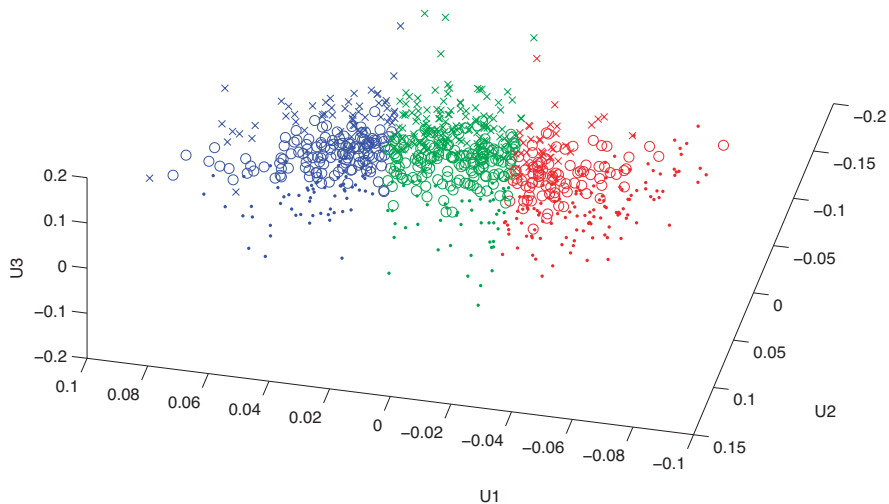


Figure 8. Plot of the SVD of galaxy data, overlaid with the SDD classification.

Finding minerals

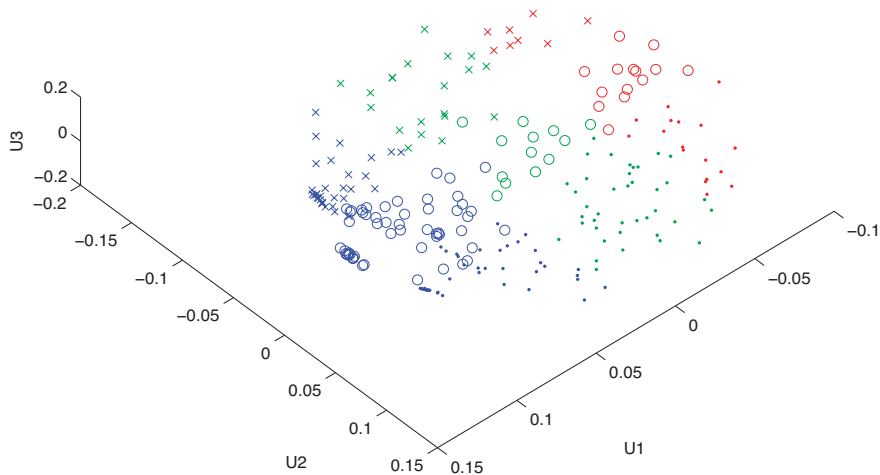


Figure 11. *Plot with position from the SVD, and color and shape labelling from the SDD.*

Clustering information from SDD added, first bump defines the colour, second the marker. Colour corresponds to the depth of the sample.

Outline

- 1 Hunting the Bump
- 2 Semi-Discrete Decomposition
- 3 The Algorithm
- 4 Applications
 - SDD alone
 - SVD + SDD
- 5 Wrap-Up

Lessons learned

- SDD finds the local areas of values with uniform magnitude
→ easier interpretation, 'orthogonal' view to SVD
- Finding SDD is hard and requires a heuristic
- Together SVD and SDD provide a strong analysis toolset

Suggested reading

- Skillicorn, Chapters 5 & 6
- Tamara G. Kolda & Diane P. O'Leary, 1998. A Semidiscrete Matrix Decomposition for Latent Semantic Indexing in Information Retrieval. *ACM Trans. Inf. Syst.* 16(4), pp. 322–346
DOI: 10.1145/291128.291131
- Tamara G. Kolda & Diane P. O'Leary, 2000. Algorithm 805: Computation and Uses of the Semidiscrete Matrix Decomposition. *ACM Trans. Math. Software* 26(3), pp. 415–435
DOI: 10.1145/358407.358424