# Visualizing Overlapping Biclusterings and Boolean Matrix Factorizations

THIBAULT MARETTE · PAULI MIETTINEN · STEFAN NEUMANN





# **MOTIVATION**

- Finding biclusters (sets of items closely related to each other) is a well-studied problem
- Visualization helps *assessing* outputs of biclustering algorithms

#### **Research Questions:**

- How can we visualize a bipartite graph given a set of biclusters?
- Can we define metrics to compare visualizations?



Figure 1: Visualization of the same biclustering using permutation from ADVISER [2] (top) and from our TSP-based heuristic (bottom).

# **OBJECTIVE FUNCTIONS**

The three criteria that we study are as follows:

Proximity: Rows and columns of each bicluster should be close to each other.



The code of our visualization tool is open-source and available online.

**EXPERIMENTS** 



github.com/tmarette/biclusterVisualization

We picked our **color scheme** using colorbrewer [1] to ensure readability and accessibility.



Ordering

Color scheme

**Qualitative Evaluation.** We notice an increase in the quality of the visualization by comparing our method to the baseline algorithm ADVISER.



Size of the consecutive cluster areas: Rows and columns of each bicluster should form large consecutive areas.

$$S_{\text{clArea}}(R_i, C_i) = \sum_{(X,Y) \in cons(R_i) \times cons(C_i)}$$

$$\sum_{(B_{\cdot})\times cons(C_{\cdot})} |X \times Y|^2.$$

**Related work.** Jin et al. [3] and Colantonio et al. [2] studied the visualization of a given set of overlapping biclusters as a biadjacency matrix.

- Colantonio et al. [2] proposed a greedy heuristic called ADVISER to minimize gaps in the visualization of each bicluster.
- However, similar but non-overlaping biclusters are not visualized close to one another.

### **PRE-PROCESSING**

**Problem setup.** Let  $G = (R \cup C, E)$  a bipartite graph, and  $((R_1, C_1), \ldots, (R_k, C_k))$  a biclustering of *G*, where  $R_i \subseteq R$  and  $C_i \subseteq C$  for all *i*. We Compute a row permutation  $\sigma_R$  and a column permutation  $\sigma_C$  of the biadjacency matrix A to *optimally* visualize the biclustering. In the visualization, we set  $A_{\sigma_R(r),\sigma_C(c)} = 1$  iff  $(r,c) \in E$ .

**Observation.** The pre-processing stems from:

- Similar items should be close to each other in the final permutation
- $\Rightarrow$  Rows and columns from the same clusters should be contiguous in the final permutation.

We partition the rows and columns into *row blocks* and *column blocks*. We ensure that each cluster can be expressed as the union of a set of blocks.

Formally, for  $r \in [m]$ , let  $clusters_R(r) = \{i : r \in R_i\}$ denote the set of indices of all row clusters that

Size of uninterrupted areas: Areas that belong to biclusters should form large uninterrupted areas (not limited to individual biclusters).



$$S_{\text{uninter}}^{R}(b_{i}^{R}) = \sum_{Y \in cons(nonzero(b_{\cdot}^{R}))} |b_{i}^{R} \times Y|^{2}.$$

- All three objective functions capture a different aspect of a good visualization
- They represent a *global* state of the visualization
- They are costly to compute

## DEMERIT

We introduce *demerit* as a fast-to-compute local distance function between two blocks. If two blocks are dissimilar, their demerit is high.

Formally, the *demerit* for row block  $b^R$  and column blocks  $b_i^C$  and  $b_i^C$  is given by:

$$lemerit(b^{R}; b_{i}^{C}, b_{j}^{C}) = \begin{cases} |b^{R}| \cdot (|c_{1} \cup c_{2}| + 1) \\ \text{if } c_{1} = \emptyset \text{ or } c_{2} = \emptyset \\ |b^{R}| \cdot (|c_{1} \cup c_{2}| - |c_{1} \cap c_{2}|) \\ \text{otherwise} \end{cases}$$

where  $c_1 = clusters_R(b^R) \cap clusters_C(b_i^C)$ , and  $c_2 =$  $clusters_R(b^R) \cap clusters_C(b_i^C).$ 





#### **ADVISER** Our method Increased size of uninterrupted areas.



**ADVISER** 



Our method Increased coherence.

Quantitative Evaluation. We gathered objective functions values of all methods across multiple datasets, using varying numbers of clusters.



contain row r.



Now, the *row block of* ris given by  $block_R(r) =$  $\{r' : clusters_R(r) =$  $clusters_R(r')$ }, i.e., it is the set of all rows r' that are contained in exactly the same row clusters as r.

 $\Rightarrow$  We only have to find permutations of the *row and* column blocks.

#### **R**EFERENCES

- [1] C. A. Brewer, M. Harrower, B. Sheesley, A. Woodruff, and D. Heyman. Colorbrewer 2.0: color advice for cartography. *The Pennsylvania State* University. http://colorbrewer2. org/.
- [2] A. Colantonio, R. Di Pietro, A. Ocello, and N. V. Verde. Visual role mining: A picture is worth a thousand roles. IEEE Trans. Knowl. Data *Eng.*, 24(6):1120–1133, 2011.
- [3] R. Jin, Y. Xiang, D. Fuhry, and F. F. Dragan. Overlapping matrix pattern visualization: A hypergraph approach. In IEEE Int. Conf. Data Min., pages 313-322, 2008.

To measure the demerit of the column block permuta*tion*  $\sigma_C$ , we set

$$demerit(\sigma_C) = \sum_{b^R \in \mathcal{B}^R} \sum_{i=1}^{t-1} demerit(b^R; b^C_{\sigma_C(i)}, b^C_{\sigma_C(i+1)}).$$

- The *demerit* is easier to optimize than the other three objective functions
- Can be computed locally
- Only have to consider consecutive pairs of column blocks  $b_{\sigma_C(i)}^C$  and  $b_{\sigma_C(i+1)}^C$
- Defined for all pairs of row and column blocks

**TSPheuristic.** We use a TSP solver in order to find a permutation minimizing *demerit* between blocks. The permutation returned by the TSP solver is the permutation used in the visualization.



 $\Rightarrow$  Minimizing the demerit achieves very high quality across all objective functions.

#### ACKNOWLEDGEMENT

This work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation.