

# R. Belohlavek, M. Krmelova: Beyond Boolean Matrix Decompositions: Toward Factor Analysis and Dimensionality Reduction of Ordinal Data

Department of Computer Science, Palacky University, Olomouc (17. listopadu 12, CZ-77146 Olomouc, Czech Republic)

## General Matrix Decomposition

- input:  $n \times m$  object-attribute matrix  $I$  with entries  $I_{ij}$  expressing grades to which object  $i$  has attribute  $j$
- output: an  $n \times k$  object-factor matrix  $A$  and a  $k \times m$  factor-attribute matrix  $B$
- grades are taken from a bounded scale  $L$
- goal: find  $A$  and  $B$  with  $k$  (# factors) as small as possible

$$I = A \circ B$$

## Essential Parts of Matrices over Scales

- different role of matrix entries for decompositions
- essential part of  $I$ , a minimal set of entries whose coverage guarantees an exact decomposition of  $I$
- the number of such entries is significantly smaller than the number of all entries

**Definition 1**  $J \leq I$  is called an *essential part* of  $I$  if  $J$  is minimal w.r.t.  $\leq$  having the property that for every  $\mathcal{F} \subseteq \mathcal{B}(I)$  we have: if  $J \leq A_{\mathcal{F}} \circ B_{\mathcal{F}}$  then  $I = A_{\mathcal{F}} \circ B_{\mathcal{F}}$ .

- intervals in  $\mathcal{B}(I)$  play a crucial role for our considerations
- for  $C \in L^{1 \times n}$ ,  $D \in L^{1 \times m}$ , put  $\gamma(C) = \langle C^{\downarrow}, C^{\uparrow} \rangle$  and  $\mu(D) = \langle D^{\downarrow}, D^{\uparrow} \rangle$
- $\mathcal{I}_{C,D}$  the interval

$$\mathcal{I}_{C,D} = [\gamma(C), \mu(D)]$$

in  $\mathcal{B}(I)$ , i.e. the set

$$[\gamma(C), \mu(D)] = \{ \langle E, F \rangle \in \mathcal{B}(I) \mid \gamma(C) \leq \langle E, F \rangle \leq \mu(D) \}.$$

**Lemma 1** If  $\langle E, F \rangle \in \mathcal{I}_{C,D}$  then  $C^{\uparrow} \circ D \leq E^{\uparrow} \circ F$ .

**Lemma 2** Let  $\langle E, F \rangle \in \mathcal{B}(X, Y, I)$ ,  $a, b \in L$ . Then  $a \otimes b \leq E(i) \otimes F(j)$  if and only if for some  $c, d$  with  $a \otimes b \leq c \otimes d$  we have  $\langle E, F \rangle \in \mathcal{I}_{\{c\}^i, \{d\}^j}$ .

Now, for a given matrix  $I \in L^{n \times m}$ , let  $\mathcal{I}_{ij} = \{ \mathcal{I}_{\{a\}^i, \{b\}^j} \mid a \otimes b = I_{ij} \}$  and put

$$\mathcal{I}_{ij} = \bigcup \mathcal{I}_{ij}.$$

**Theorem 1** A rectangle corresponding to  $\langle E, F \rangle \in \mathcal{B}(X, Y, I)$  covers  $\langle i, j \rangle$  in  $I$  iff  $\langle E, F \rangle \in \mathcal{I}_{ij}$ .

Denote by  $\mathcal{E}(I) \in L^{n \times m}$  the matrix over  $L$  defined by

$$(\mathcal{E}(I))_{ij} = \begin{cases} I_{ij} & \text{if } \mathcal{I}_{ij} \text{ is non-empty and minimal w.r.t. } \subseteq, \\ 0 & \text{otherwise.} \end{cases}$$

**Theorem 2**  $\mathcal{E}(I)$  is the unique essential part of  $I$ .

**Theorem 3** Let  $\mathcal{G} \subseteq \mathcal{B}(\mathcal{E}(I))$  be a set of factor concepts of  $\mathcal{E}(I)$ , i.e.  $\mathcal{E}(I) = A_{\mathcal{G}} \circ B_{\mathcal{G}}$ . Then every set  $\mathcal{F} \subseteq \mathcal{B}(I)$  containing for each  $\langle C, D \rangle \in \mathcal{G}$  at least one concept from  $\mathcal{I}_{C,D}$  is a set of factor concepts of  $I$ , i.e.  $I = A_{\mathcal{F}} \circ B_{\mathcal{F}}$ .

## New Algorithms

The algorithms we present are inspired by GRESS [1] and ASSO [4], currently perhaps the best algorithms for the AFP and DBP, respectively.

### GreEss<sub>L</sub>

```

Input: matrix I with entries in scale L
Output: set F of factors for which I = A_F o B_F
1 G ← COMPUTEINTERVALS(I)
2 U ← {(i, j) | I_ij > 0}; F ← ∅
3 while U is non-empty do
4   foreach (C, D) ∈ G do
5     J ← D^↓ ⊗ C^↑; F ← ∅; s_{(C,D)} ← 0
6     while exists {a/j} ∈ C^↑ \ F s.t.
       cov(U, F ∪ {a/j}, J) > s_{(C,D)} do
7       select {a/j} maximizing cov(U, F ∪ {a/j}, J)
8       F ← (F ∪ {a/j})^↓ ⊗ J; E ← (F ∪ {a/j})^↓
9       s_{(C,D)} ← cov(U, F, J)
10    end
11    if s_{(C,D)} > s then
12      (E', F') ← (E, F)
13      (C', D') ← (C, D)
14      s ← s_{(C,D)}
15    end
16  end
17  add (E', F') to F
18  remove (C', D') from G
19  remove from U entries (i, j) covered by E' ⊗ F' in I
20 end
21 return F
    
```

### COMPUTEINTERVALS

```

Input: matrix I with entries in scale L
Output: set G ⊆ B(E(I))
1 E ← E(I)
2 U ← {(i, j) | E_ij > 0}
3 while U is non-empty do
4   D ← ∅; s ← 0
5   while exists {a/j} ∈ D s.t. cov_I(U, D ∪ {a/j}, E) > s do
6     select {a/j} maximizing cov_I(U, D ∪ {a/j}, E)
7     D ← (D ∪ {a/j})^↓ ⊗ E; C ← (D ∪ {a/j})^↓
8     s ← cov_I(U, D, E)
9   end
10  add (C, D) to G
11  remove from U entries (i, j) covered by C^↑ ⊗ D^↑ in I
12 end
13 return G
    
```

### ASSO<sub>L</sub>

```

Input: matrix I with entries in scale L, k ≥ 1, w^+, w^-, τ
Output: set F of factors
1 compute association matrix A
2 F ← ∅
3 for l = 1 ... k do
4   select (C, A_l) maximizing cover(F ∪ {(C, A_l)}, I, w^+, w^-)
5   add (C, A_l) to F
6 end
7 return F
    
```

The association matrix  $A$  is then defined by

$$A_{ij} = \text{round}_{\tau}(c(i \Rightarrow j, I)),$$

where  $\text{round}_{\tau}$  is defined for  $r \in [0, 1]$  by

$$\text{round}_{\tau}(r) = \begin{cases} r_+ = \min\{a \in L \mid a \geq r\} & \text{if } r_+ \leftrightarrow r \geq \tau, \\ r_- = \max\{a \in L \mid a < r\} & \text{otherwise.} \end{cases}$$

Here,  $r_+ \leftrightarrow r = \min(r_+ \rightarrow r, r \rightarrow r_+)$  is the biresiduum (logical equivalence). Note that  $\text{round}_{\tau}$  is used to obtain a matrix  $A$  with entries in  $L$  which is needed because the rows of  $A$  are the candidate basis vectors.

## Experimental Evaluation

- experimental evaluation of the presented algorithms on real and synthetic data
- the ability of the extracted factors to explain (i.e. reconstruct) the input data

## Real Data

### Characteristics of real data

dataset	$\ I\ $	$\ \mathcal{E}(I)\ $	$\ \mathcal{E}(I)\ /\ I\ $	size	$ L $
Breeds	1963	362	0.184	151 × 11	6
Decathlon	266	59	0.221	28 × 10	5
IPAQ	41624	1281	0.031	4510 × 16	3
Music	20377	5952	0.292	900 × 26	7
Music reduced	771	213	0.276	30 × 26	7

### Coverage of Data by Factors

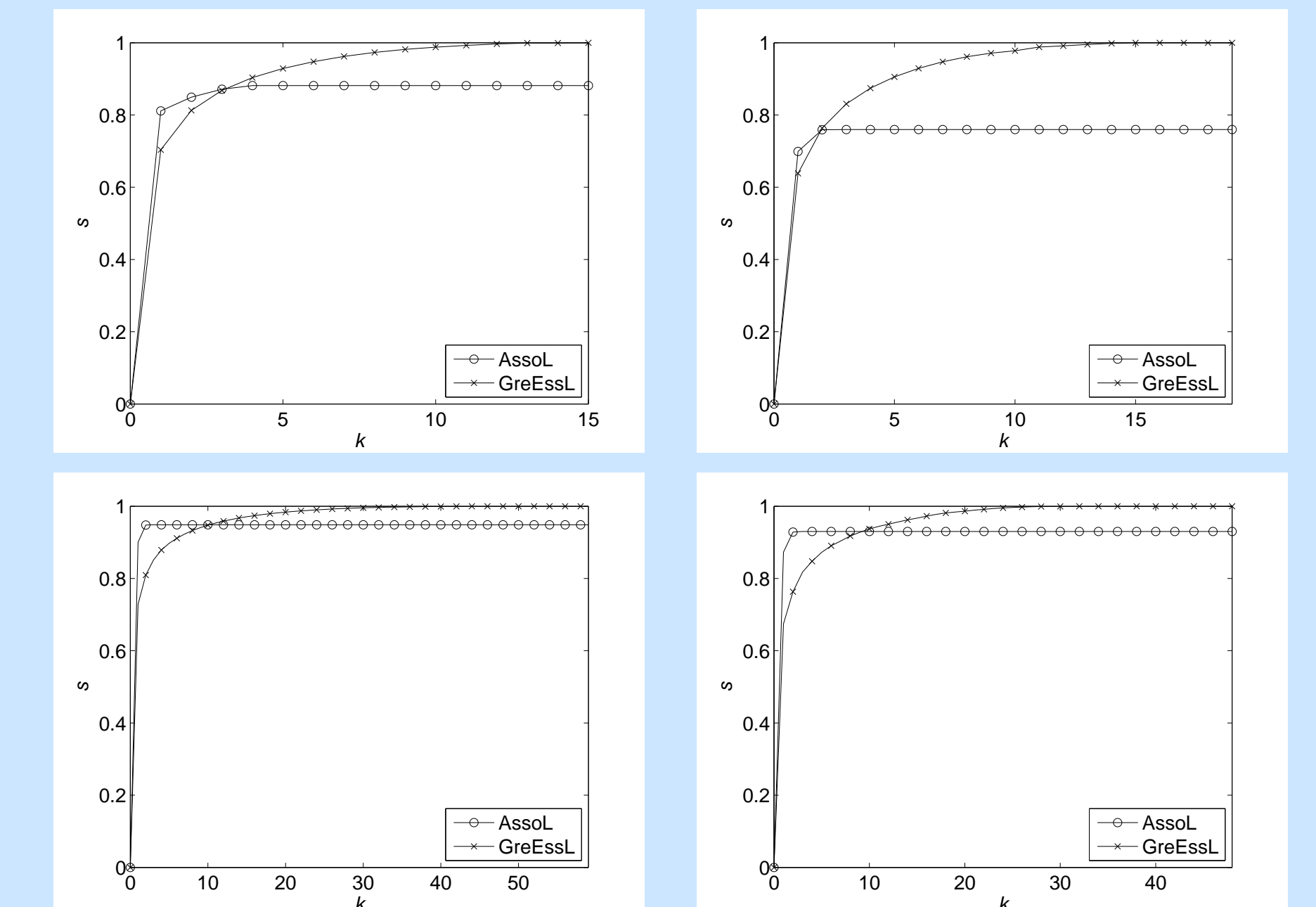
- $\|A\|$  denotes the number of non-zero entries in matrix  $A$
- numbers of factors needed to achieve a coverage  $s = \{0.75, 0.85, 0.95, 1\}$
- Breeds - ASSO<sub>L</sub> 2, 3, NA, NA; GRESS<sub>L</sub> 3, 7, 11, 15
- Decathlon - ASSO<sub>L</sub> 2, 4, NA, NA; GRESS<sub>L</sub> 3, 5, 8, 10
- IPAQ - ASSO<sub>L</sub> 1, 1, NA, NA; GRESS<sub>L</sub> 10, 12, 15, 17
- Music - ASSO<sub>L</sub> 2, NA, NA, NA; GRESS<sub>L</sub> 7, 14, 25, 29
- Music red. - ASSO<sub>L</sub> 1, 2, NA, NA; GRESS<sub>L</sub> 1, 3, 10, 30
- "NA" = prescribed coverage is not achievable

## Synthetic Data

### Characteristics of synthetic data

dataset	size	$ L $	$k$	distribution on $L$	avg $\ I\ $	avg $\ \mathcal{E}(I)\ $	avg $\ \mathcal{E}(I)\ /\ I\ $
Set 1	50×50	3	10	$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	2452	193	0.079
Set 2	50×50	5	10	$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}$	2499	358	0.143
Set 3	100×50	5	25	$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}$	4998	614	0.123
Set 4	100×100	5	20	$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}$	10000	2130	0.213
Set 5	150×150	10	25	$\frac{1}{10}$ for all	22498	5759	0.256

### Coverage $s$ by the first $k$ factors



## Discussion

- the first couple of factors produced by ASSO<sub>L</sub> has a better coverage compared to the same number of factors produced by GRESS<sub>L</sub>
- beyond certain coverage, ASSO<sub>L</sub> stops producing factors and is not able to compute an (exact) decomposition of  $I$ , while GRESS<sub>L</sub> always computes an exact decomposition
- GRESS<sub>L</sub> produces easier interpretable factors compared to ASSO<sub>L</sub>
- $|L| > 2$  (non-Boolean case), rectangles with values "around the middle" in  $L$ , such as 0.5, which may be produced as factors by ASSO<sub>L</sub>
- on average, GRESS<sub>L</sub> requires 30% less factors to achieve a prescribed coverage comparing with the fast greedy algorithm described in [2]

## Previous Work

- R. Belohlavek, M. Trnecka, *From-below approximations in Boolean matrix factorization: Geometry and new algorithm*. (submitted, available at arXiv).
- Belohlavek R., Vychodil V.: *Factor analysis of incidence data via novel decomposition of matrices*, LNAI 5548(2009), 83-97.
- R. Belohlavek, M. Krmelova, *Factor analysis of sports data via decomposition of matrices with grades*, Proceedings of the 9th International Conference on CLA (2012), 305-316.
- P. Miettinen, T. Mielikäinen, A. Gionis, G. Das, H. Mannila, *The discrete basis problem*, IEEE TKDE 20(2008), 1348-62.