

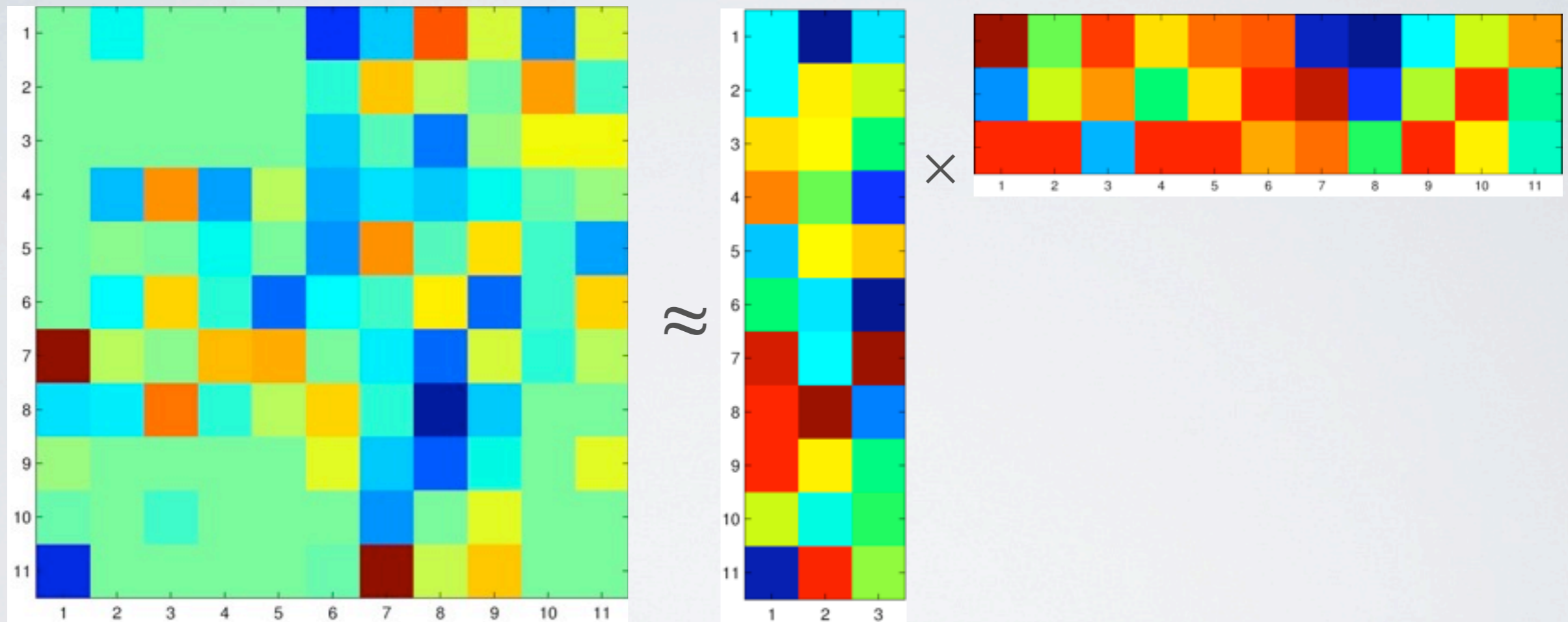
BOOLEAN MATRIX FACTORIZATIONS

Pauli Miettinen
Leap day, 2012



mp | max planck institut
informatik

MATRIX FACTORIZATIONS

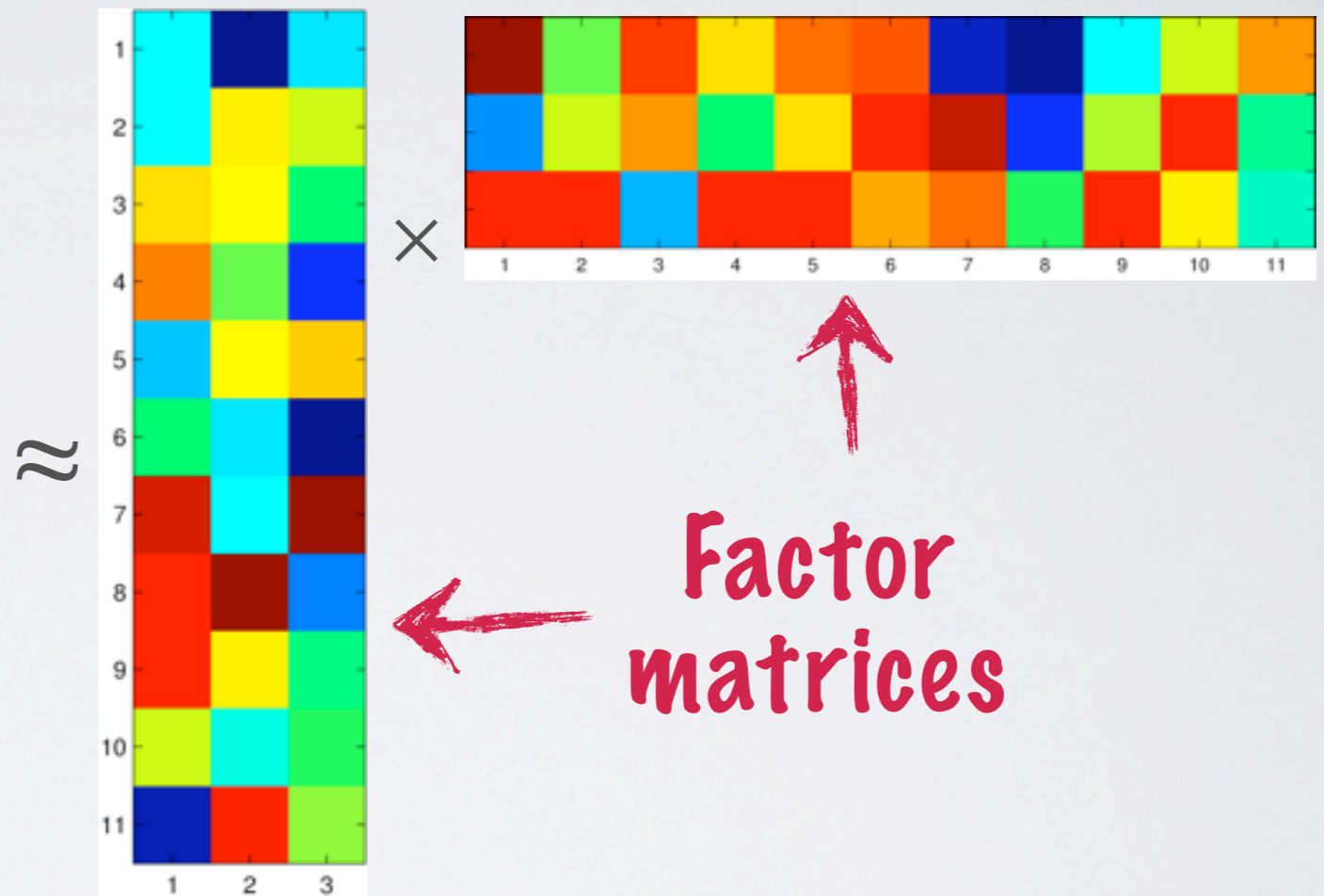
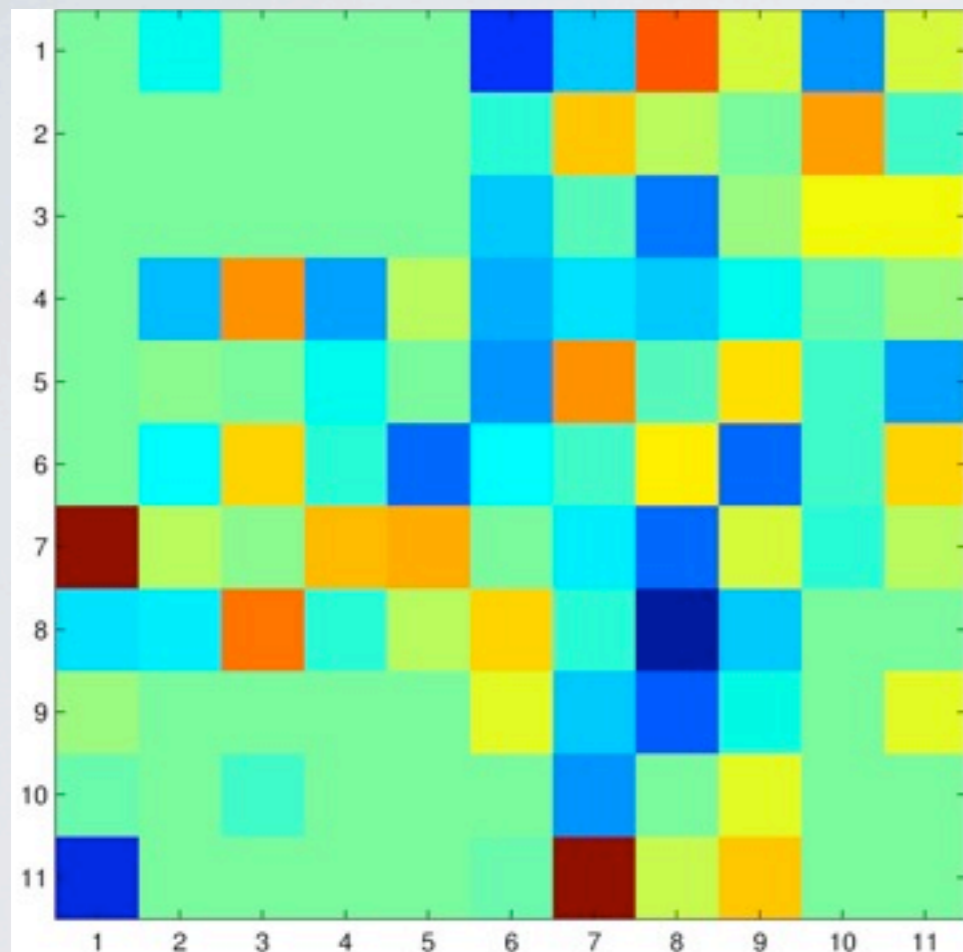


MATRIX FACTORIZATIONS

- A **factorization** of matrix \mathbf{X} represents it as a product of two (or more) **factor matrices**: $\mathbf{X} = \mathbf{AB}$
 - \mathbf{X} is n -by- m , \mathbf{A} is n -by- k , and \mathbf{B} is k -by- m
 - k is the **size** (or **rank**) of the factorization
- Factorization can be **exact** ($\mathbf{X} = \mathbf{AB}$) or **approximate** ($\mathbf{X} \approx \mathbf{AB}$)



MATRIX FACTORIZATIONS



Factor matrices

Rank = 3



SOME LINEAR ALGEBRA

- A set of vectors is *linearly independent* if no vector in the set can be expressed as a linear combination of the others
- A matrix \mathbf{X} is orthogonal if and only if $\mathbf{X}\mathbf{X}^T = \mathbf{X}^T\mathbf{X} = \mathbf{I}$
- The **column rank** of a matrix is the number of linearly independent columns it has
 - Equals the row rank of the matrix
 - ⇒ **the rank** of a matrix is its column rank = row rank



ON MATRIX RANK

- Matrix \mathbf{X} has $\text{rank}(\mathbf{X}) = 1$ iff $\mathbf{X} = \mathbf{a}\mathbf{b}^T$
 - Outer product of column vectors \mathbf{a} and \mathbf{b}
- Matrix \mathbf{X} has $\text{rank}(\mathbf{X}) \leq k$ if it can be represented as a sum of k rank-1 matrices
 - Smallest such k is the rank of \mathbf{X}
- Equivalently, $\text{rank}(\mathbf{X}) \leq k$ iff there is a rank- k factorization of \mathbf{X}

$$\mathbf{X} = \sum_{i=1}^k \mathbf{a}_i \mathbf{b}_i^T = \mathbf{A}\mathbf{B}$$



MATRIX DISTANCES

- The Frobenius norm: $\|\mathbf{X}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^m x_{ij}^2}$
 - We drop the F in Frobenius for now...
- The sum of absolute values: $|\mathbf{X}| = \sum_{i=1}^n \sum_{j=1}^m |x_{ij}|$
 - If \mathbf{X} is binary, $|\mathbf{X}| = \|\mathbf{X}\|^2$



FAMOUS MATRIX FACTORIZATIONS

- Eigendecomposition: $\mathbf{X} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$
 - \mathbf{X} is square; \mathbf{Q} is orthogonal with the eigenvectors of \mathbf{X} ; $\mathbf{\Lambda}$ is diagonal and has the eigenvalues
- Singular value decomposition: $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
 - \mathbf{U} and \mathbf{V} are orthogonal, $\mathbf{\Sigma}$ is diagonal with the singular values
- Non-negative matrix factorization: $\mathbf{X} = \mathbf{WH}$
 - All matrices are non-negative



OTHER FAMOUS MATRIX FACTORIZATIONS

- k -means clustering
- tiling databases



K-MEANS AS MATRIX FACTORIZATION

- Given m data points (in \mathbf{R}^n), partition them in k clusters such that

$$\sum_{i=1}^k \sum_{\mathbf{x}_j \in C_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|_2^2$$

is minimized **Distance of data points to cluster centroid**

- Equivalently, minimize $\|\mathbf{X} - \mathbf{MC}\|^2$, where
 - \mathbf{X} is the data (n -by- m), \mathbf{M} (n -by- k) has the centroids as its columns, and \mathbf{C} (k -by- m) is a **cluster assignment matrix**
 - Each column of \mathbf{C} has exactly one 1, and rest is 0s



TILING AS MATRIX FACTORIZATION

- Maximum k -tiling: find at most k **tiles** such that the tiling has maximum area
 - Data is binary matrix, tiles are submatrices full of 1s
 - Area of a tiling is the number of 1s in the data that belong to at least one tile
- We turn this to *minimum-error tiling*
 - Minimize the number of 1s in the data that do not belong to any tile



TILING AS MATRIX FACTORIZATION

- We want to find factor matrices **A** and **B** such that $(\mathbf{AB})_{ij} = 1$ iff element (i, j) belongs to at least one tile
 - Minimize $|\mathbf{X} - \mathbf{AB}|$
- Single tile is an outer product of two binary vectors: \mathbf{ab}^T
 - $b_j = 1$ if an item j belongs to the tile; $a_i = 1$ if a transaction i belongs to the tile
- But how to combine the tiles?



COMBINING THE TILES

- The problem: $\sum_{i=1}^k \mathbf{a}_i \mathbf{b}_i^T$ is not binary
 - $|\mathbf{X} - \mathbf{AB}|$ will add an error every time $x_{ij} = 1$ belongs to more than one tile
- Solution: don't count multiplicity
 - Define $1 + 1 = 1$



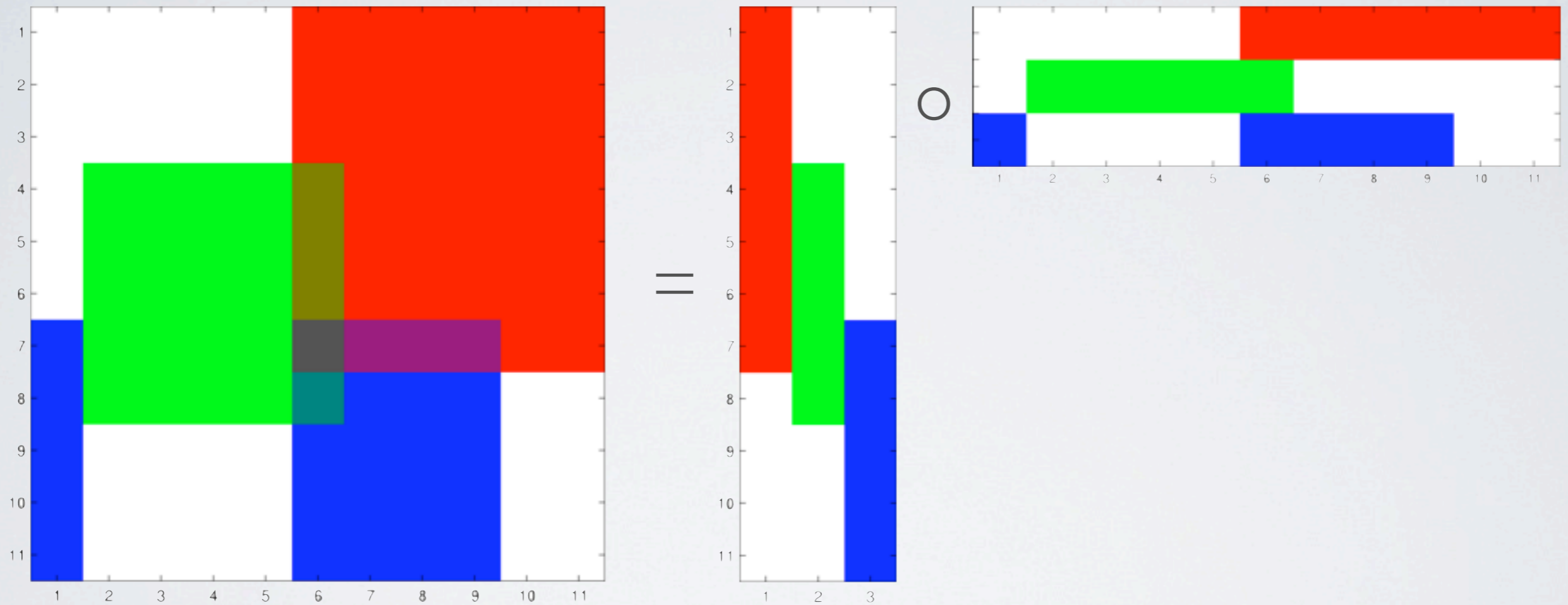
THE BOOLEAN MATRIX PRODUCT

- As normal matrix product, but with addition defined as $1 + 1 = 1$ (logical OR)
- Closed under binary matrices
- Corresponds to set union operation

$$(\mathbf{X} \circ \mathbf{Y})_{ij} = \bigvee_{l=1}^k x_{il} y_{lj}$$



THE BOOLEAN MATRIX PRODUCT

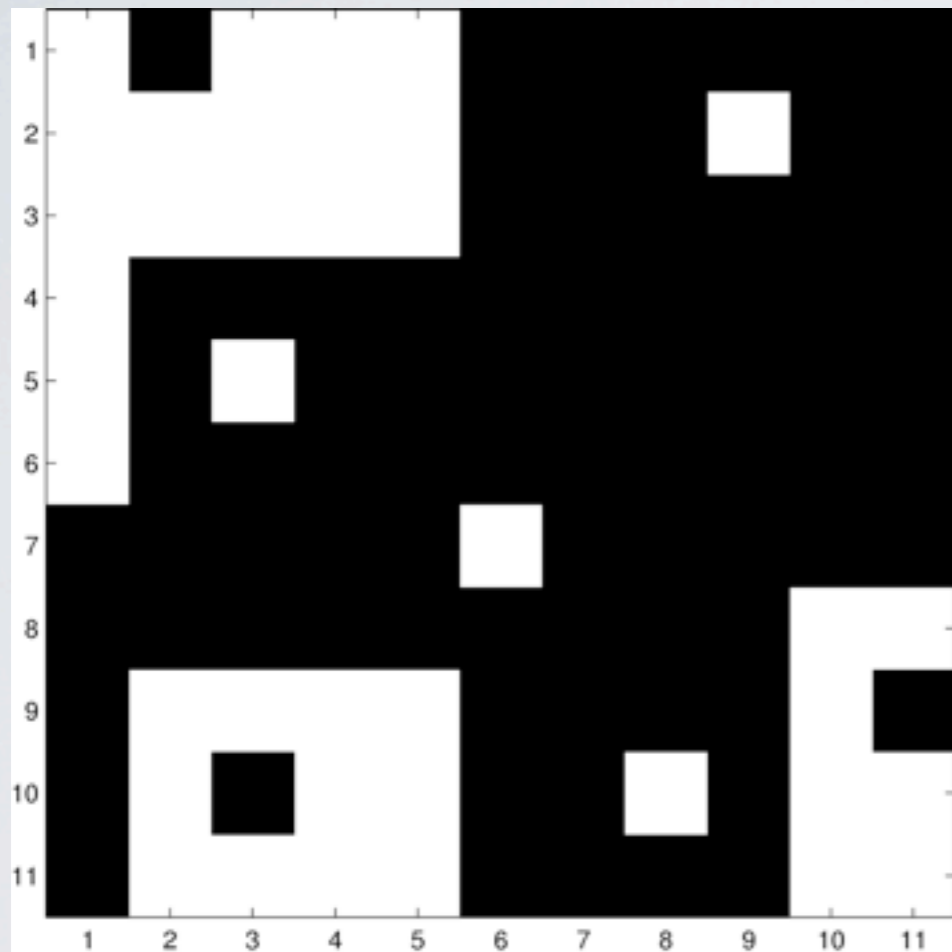


TILING REVISITED

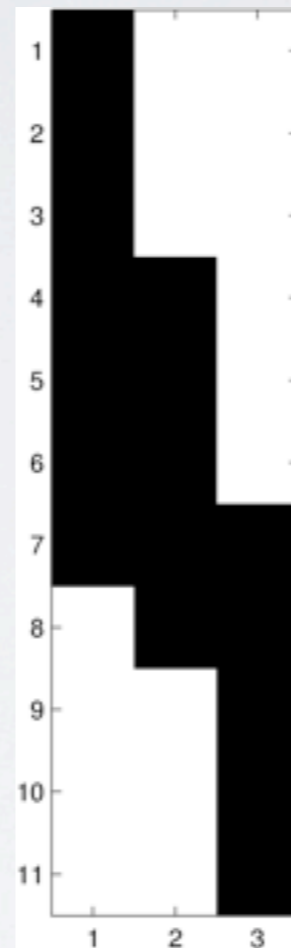
- Given transaction data as an n -by- m binary matrix \mathbf{X} and integer k , find binary matrices \mathbf{A} (n -by- k) and \mathbf{B} (k -by- m) such that if $(\mathbf{A} \circ \mathbf{B})_{ij} = 1$, then $\mathbf{X}_{ij} = 1$ and $|\mathbf{X} - \mathbf{A} \circ \mathbf{B}|$ is minimized
 - Requirement makes sure that tiles have only 1s that appear in the data
 - What happens if we remove this restriction?



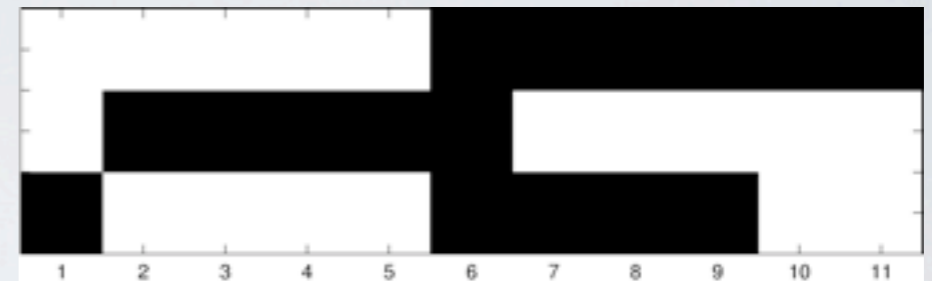
BOOLEAN MATRIX FACTORIZATIONS



\approx



\circ



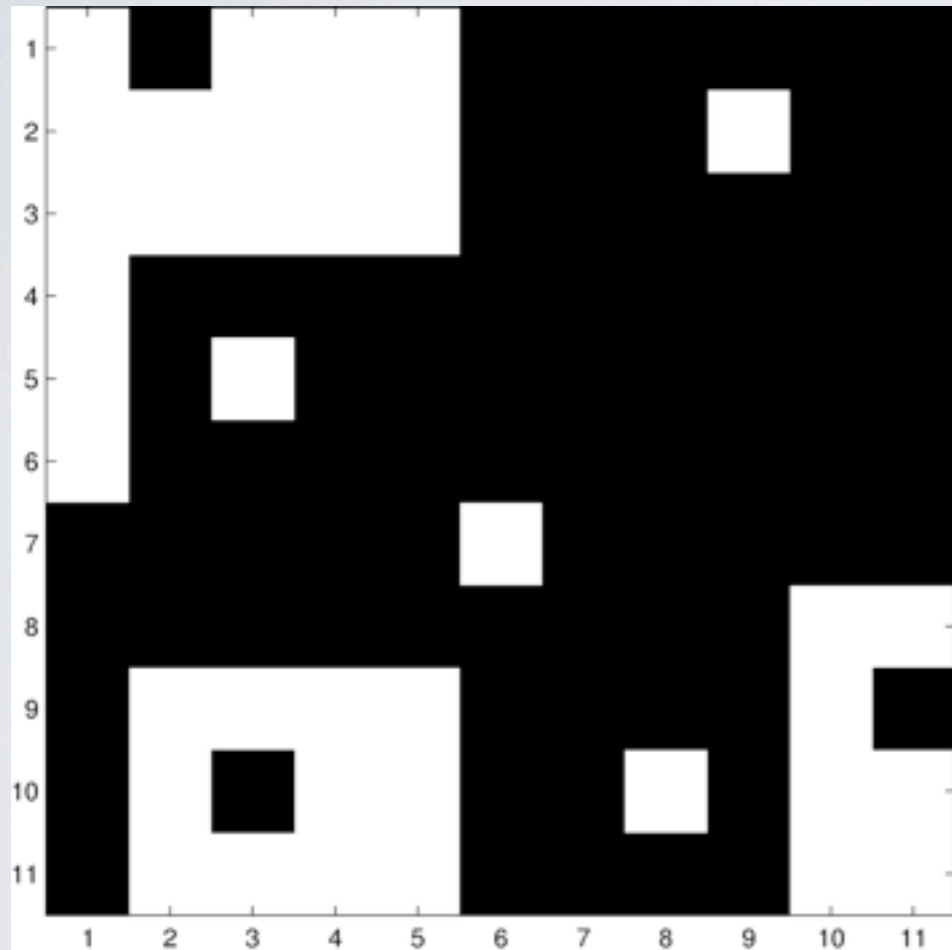
BOOLEAN MATRIX FACTORIZATIONS

Definition (BMF). Given an n -by- m binary matrix \mathbf{A} and non-negative integer k , find n -by- k binary matrix \mathbf{B} and k -by- m binary matrix \mathbf{C} such that they minimize

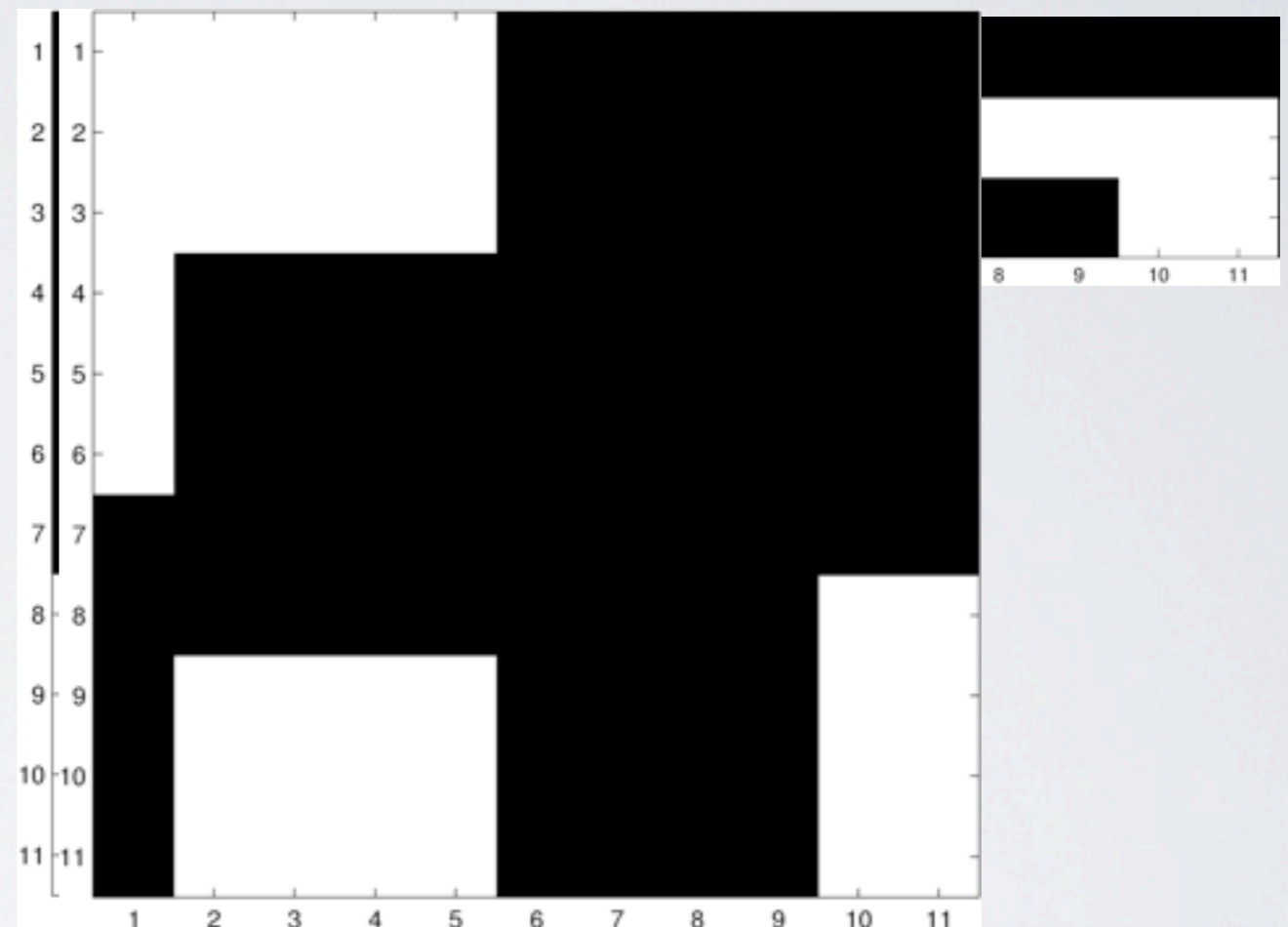
$$|\mathbf{A} \otimes (\mathbf{B} \circ \mathbf{C})| = \sum_{i,j} |a_{ij} - (\mathbf{B} \circ \mathbf{C})_{ij}|$$



BOOLEAN MATRIX FACTORIZATIONS



\approx



WHAT ABOUT DATA MINING?

- Factors provide groups of objects that ‘go together’
 - Everything is binary \Rightarrow factors are sets (unlike NMF or SVD)
- Factors can overlap (unlike clustering)
- Provides a global view (unlike frequent item sets)
- Allows missing ones and zeros (unlike tiling)



BMF: A DM EXAMPLE



long-haired
well-known
male

long-haired	✓	✓	✗
well-known	✓	✓	✓
male	✗	✓	✓



BMF: A DM EXAMPLE



long-haired
well-known
male



BMF: A DM EXAMPLE

Alice & Bob: long-haired and well-known
Bob & Charles: well-known males

$$\begin{matrix}
 \text{long-haired} \\
 \text{well-known} \\
 \text{male}
 \end{matrix}
 \begin{pmatrix}
 1 & 0 \\
 1 & 1 \\
 0 & 1
 \end{pmatrix}
 \circ
 \begin{matrix}
 \mathbf{A} & \mathbf{B} & \mathbf{C} \\
 \begin{pmatrix}
 1 & 1 & 0 \\
 0 & 1 & 1
 \end{pmatrix}
 \end{matrix}$$



SOME APPLICATIONS

- Explorative data mining
 - Factors tell something about the data
- Role mining
 - Naïve approach not very good
- Entity disambiguation / synonym finding
 - Allows synonymity and polysemy
 - Might need tensors



SOME THEORY



BOOLEAN RANK

Matrix rank. The **rank** of an n -by- m matrix **A** is the least integer k such that there exists n -by- k matrix **B** and k -by- m matrix **C** for which **A** = **BC**.

Boolean matrix rank. The **Boolean rank** of an n -by- m binary matrix **A** is the least integer k such that there exists n -by- k binary matrix **B** and k -by- m binary matrix **C** for which **A** = **B** \circ **C**.



SOME PROPERTIES OF BOOLEAN RANK

- For some matrices, Boolean rank is higher than normal rank
 - Twice the normal rank is the biggest known difference
- For some matrices, Boolean rank is much smaller
 - Can be a logarithm of the normal rank
 - **Boolean matrix factorization can have smaller reconstruction error than SVD of same size**



AN EXAMPLE

$$\begin{aligned} & \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{pmatrix} \leftarrow \text{Original matrix} \\ & = \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{pmatrix} \circ \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix} \leftarrow \text{Exact Boolean rank-2 decomposition} \\ & \approx \begin{pmatrix} 1/2 & 1/\sqrt{2} \\ 1/\sqrt{2} & 0 \\ 1/2 & -1/\sqrt{2} \end{pmatrix} \begin{pmatrix} \frac{\sqrt{2}+1}{2} & \frac{\sqrt{2}+2}{2} & \frac{\sqrt{2}+1}{2} \\ 1/\sqrt{2} & 0 & 1/\sqrt{2} \end{pmatrix} \leftarrow \text{The best approximate normal rank-2 decomposition} \end{aligned}$$



COMPUTATIONAL COMPLEXITY

- Approximating the Boolean rank is as hard as approximating the minimum chromatic number of a graph
 - Read: hard to even approximate
 - Except with some sparse matrices; more on that later



COMPUTATIONAL COMPLEXITY

- Finding minimum-error BMF is NP-hard
 - NP-hard to approximate within any poly computable factor
 - Because best answer = 0 is NP-hard to recognize
 - NP-hard to approximate within additive error of $n^{1/4}$



A SUBPROBLEM AND ITS COMPLEXITY

Basis Usage (BU). Given binary matrices \mathbf{A} and \mathbf{B} , find a binary matrix \mathbf{C} that minimizes $|\mathbf{A} - \mathbf{B} \circ \mathbf{C}|$.

- Corresponds to a problem where \mathbf{A} and \mathbf{C} are just column vectors
- Error NP-hard to approximate better than in superpolylogarithmic factor

$$\Omega \left(2^{\log^{1-\varepsilon} |\mathbf{a}|} \right)$$



AN ALGORITHM



THE ASSO ALGORITHM

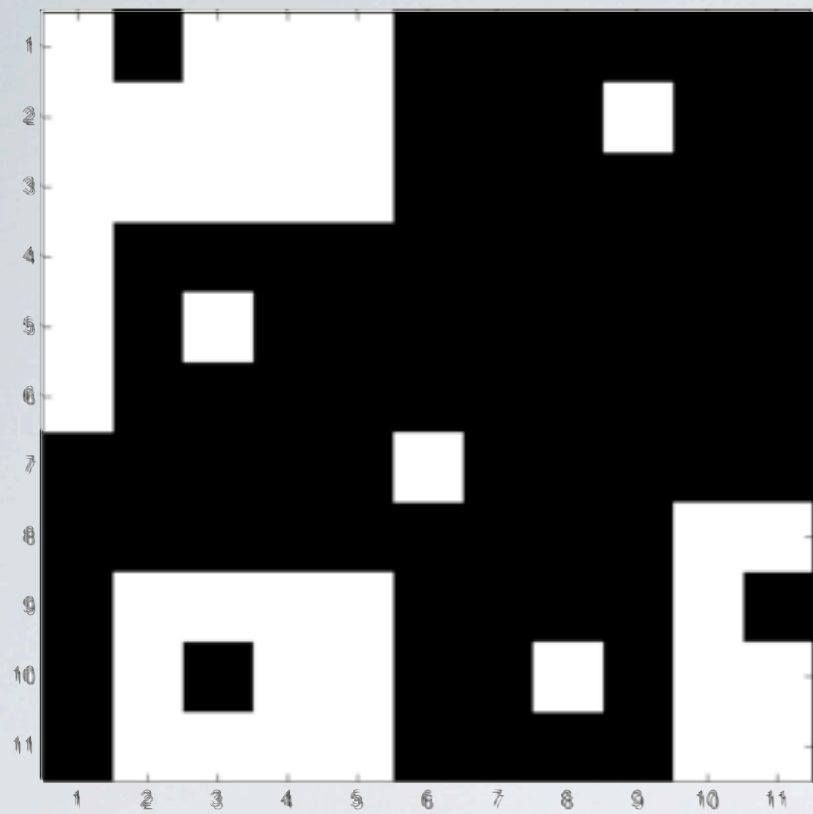
- Heuristic – too many hardness results to hope for good provable results in any case
- Intuition: If two columns share a factor, they have 1s in same rows
 - Noise makes detecting this harder
 - Pairwise row association rules reveal (some of) the factors



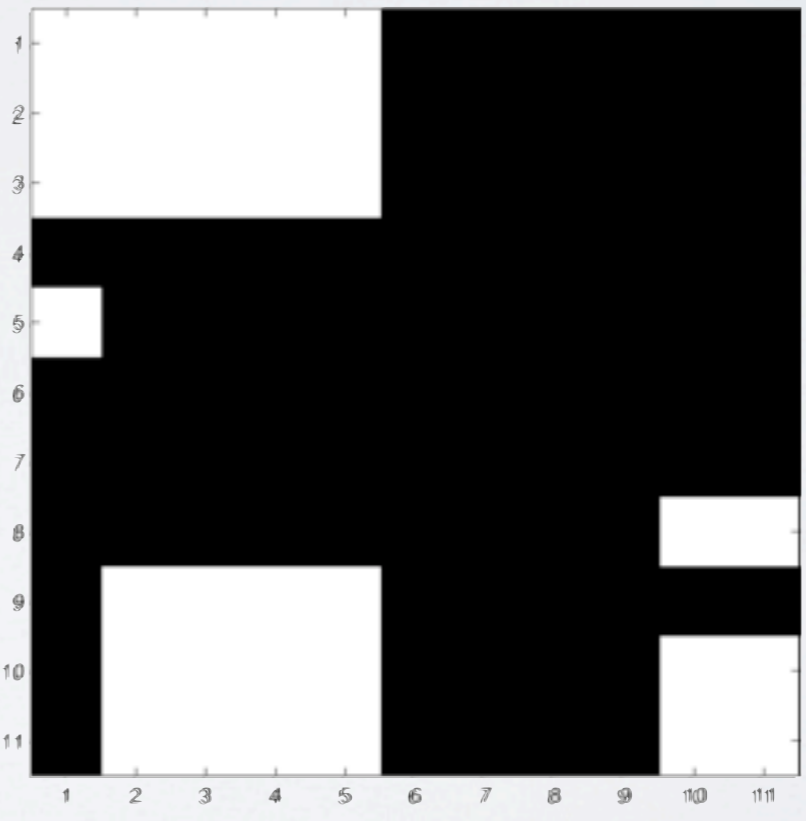
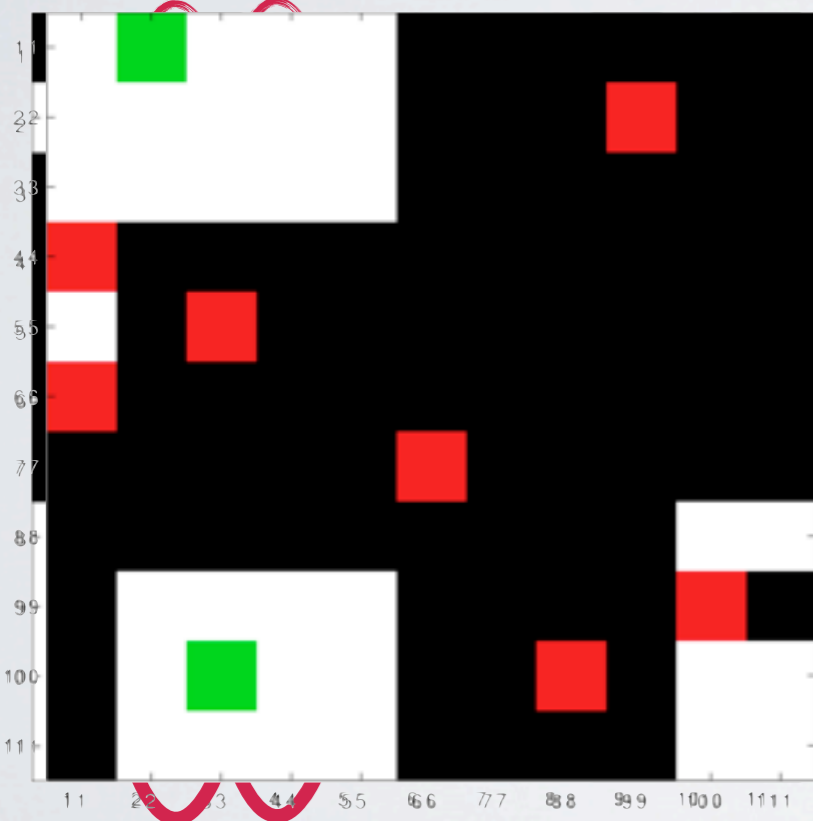
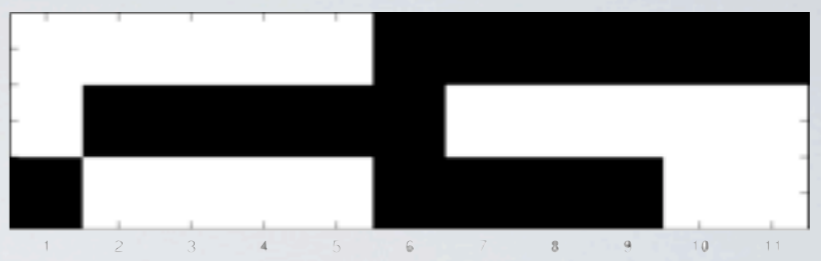
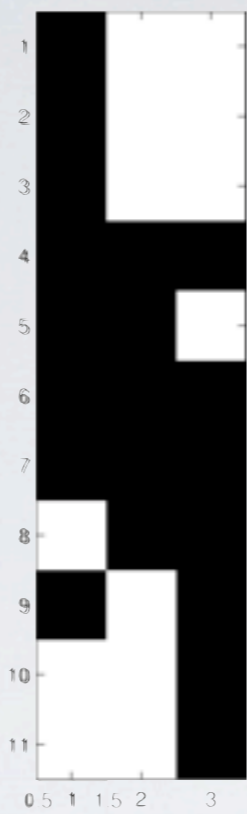
THE ASSO ALGORITHM

1. Compute pairwise association accuracies between rows of **A**
2. Round these (from a user-defined point t) to get a binary n -by- n matrix of candidate columns
3. Select greedily the candidate column that covers most of the not-yet covered 1s of **A**
4. Mark the 1s covered by the selected vector and return to 3 or quit if enough factors have been selected





\approx



SPARSE MATRICES



MOTIVATION

- Many real-world data are sparse
- With sparse input, we hope for sparse output (factors)
- Sparsity should also help with computational complexity
 - Less degrees of freedom



SPARSE FACTORIZATIONS

- Ideally, sparse matrices have sparse factors
 - Not true with many factorization methods
- Sparse Boolean matrices have sparse decompositions

Theorem 1. For any n -by- m 0/1 matrix \mathbf{A} of Boolean rank k , there exist n -by- k and k -by- m 0/1 matrices \mathbf{B} and \mathbf{C} such that $\mathbf{A} = \mathbf{B} \circ \mathbf{C}$ and

$$|\mathbf{B}| + |\mathbf{C}| \leq 2|\mathbf{A}|.$$


APPROXIMATING BOOLEAN RANK IN SPARSE MATRICES

- Intuition: Sparse matrices cannot have as complex structure as dense matrices – rank could be easier to approximate
- Recently, Belohlavek and Vychodil (2010) proposed a reduction to Set Cover, giving $O(\log n)$ approximation
 - Can yield exponential increase in instance size
- Sparsity helps!



APPROXIMATING THE BOOLEAN RANK

- Sparsity is not enough; we need some structure in it
- An n -by- m 0/1 matrix \mathbf{A} is $f(n)$ -uniformly sparse, if all of its columns have at most $f(n)$ 1s

Theorem 2. The Boolean rank of $\log(n)$ -uniformly sparse matrix can be approximated to within $O(\log(m))$ in time $\tilde{O}(m^2n)$.



NON-UNIFORMLY SPARSE MATRICES

- Uniform sparsity is very restricted; what can we do
 - Trade non-uniformity with approximation accuracy

Theorem 3. If there are at most $\log(m)$ columns with more than $\log(n)$ 1s, then we can approximate the Boolean rank in polynomial time to within $O(\log^2(m))$.

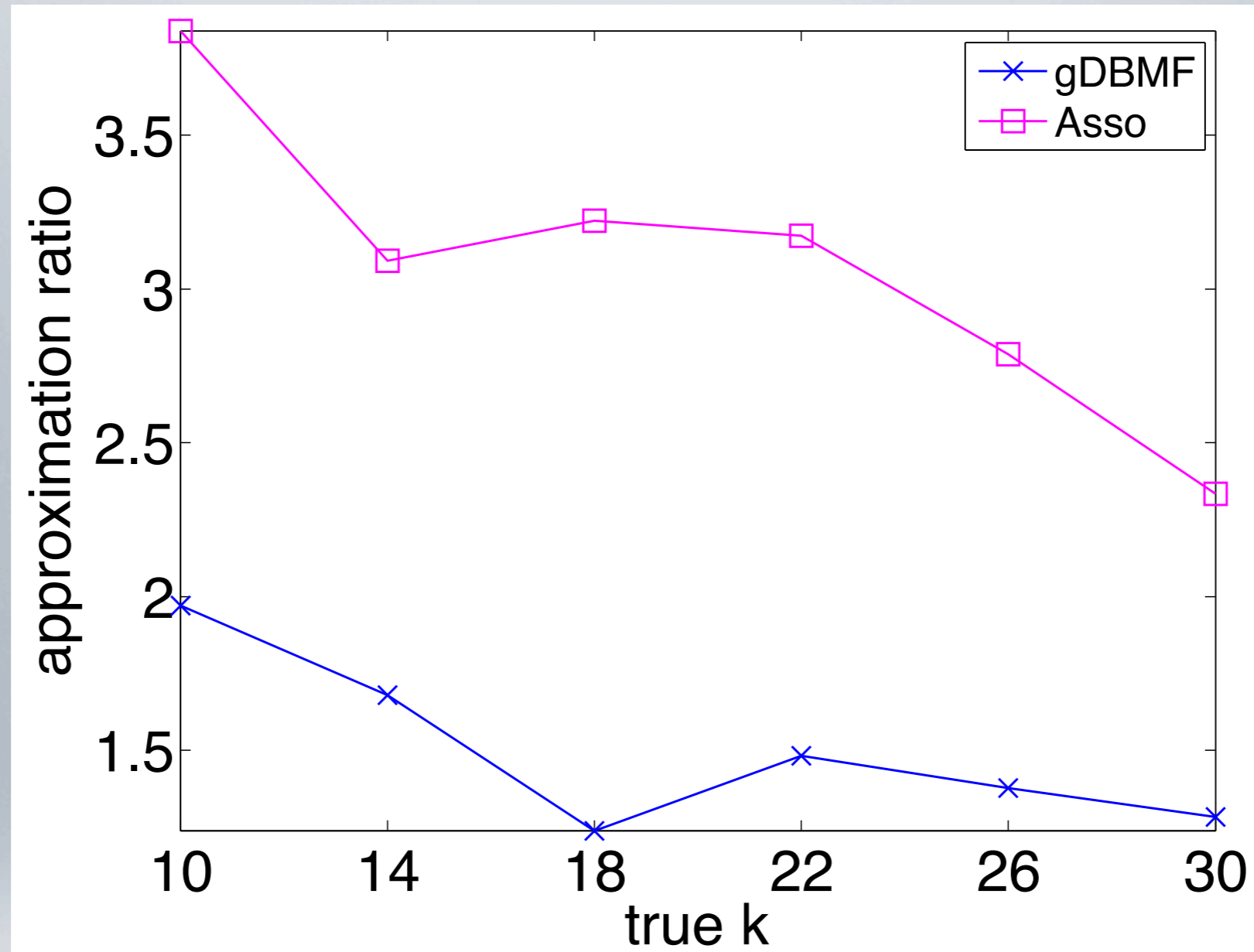


APPROXIMATING DOMINATED COVERS

Theorem 4. If n -by- m 0/1 matrix \mathbf{A} is $O(\log n)$ -uniformly sparse, we can approximate the best dominated k -cover of \mathbf{A} by $e/(e-1)$ in polynomial time.

- Dominated k -cover: The rank is k and if $(\mathbf{B} \circ \mathbf{C})_{ij} = 1$, then $\mathbf{A}_{ij} = 1$
 - This is tiling!





APPROXIMATING THE RANK



MODEL ORDER SELECTION



HOW DO I KNOW WHAT K TO USE?

Definition (BMF). Given an n -by- m binary matrix \mathbf{A} and non-negative integer k , find n -by- k binary matrix \mathbf{B} and k -by- m binary matrix \mathbf{C} such that they minimize

$$|\mathbf{A} \otimes (\mathbf{B} \circ \mathbf{C})| = \sum_{i,j} |a_{ij} - (\mathbf{B} \circ \mathbf{C})_{ij}|$$

N.B. This is nothing special to BMF!



PRINCIPLES OF GOOD K

- **Goal:** Separate noise from structure
- We assume data has BMF-type structure
 - There are k factors explaining the BMF structure
 - Rest of the data does not follow the BMF structure (noise)
- But how to decide where structure ends and noise starts?



ENTER MDL



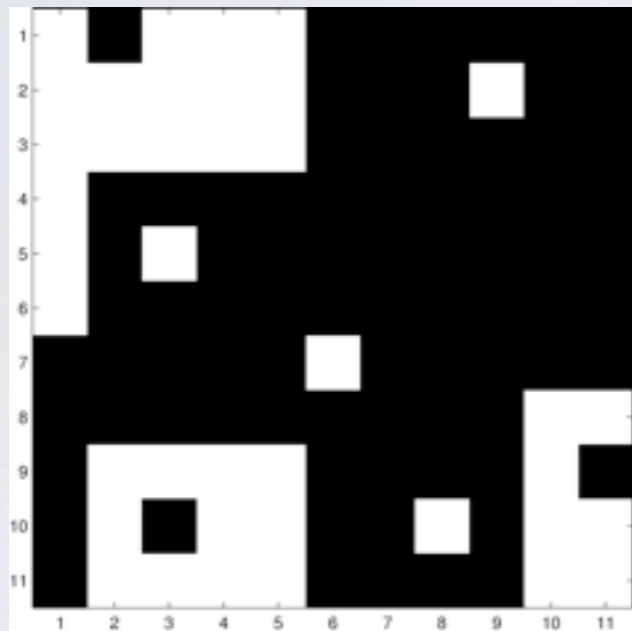
THE MINIMUM DESCRIPTION LENGTH PRINCIPLE

- Selecting k = model order selection problem
- The best model (order) is the one that allows us to represent the data with least number of bits
- **Intuition:** Using factor matrices to represent the BMF structure in the data saves space, but using them to represent noise wastes space

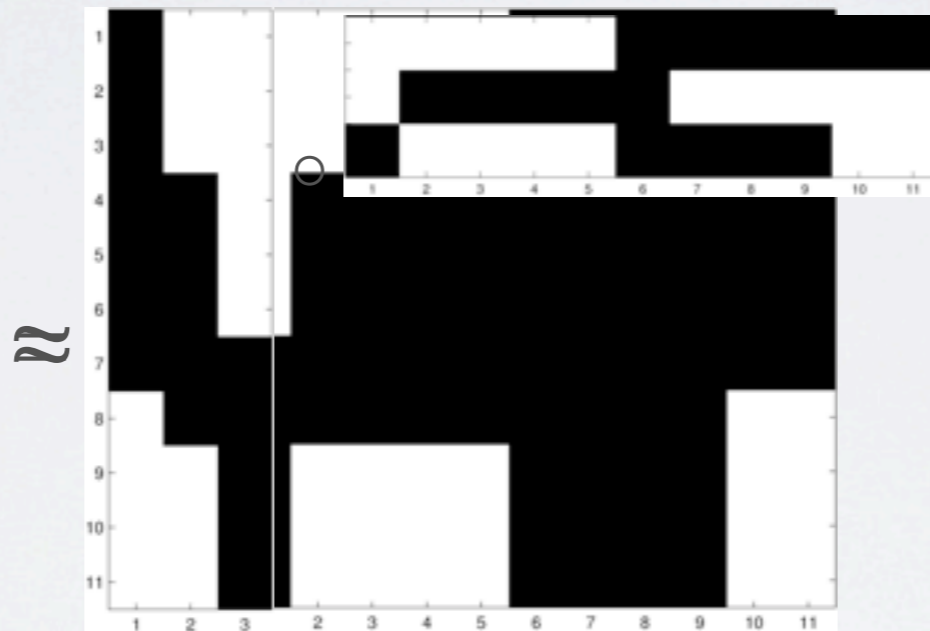


FITTING BMF TO MDL

- MDL requires exact representation

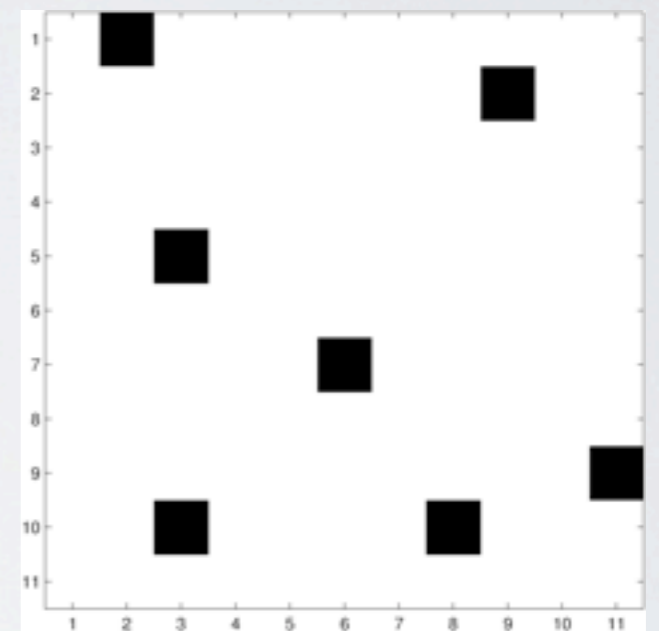


A



B \circ **C**

\otimes

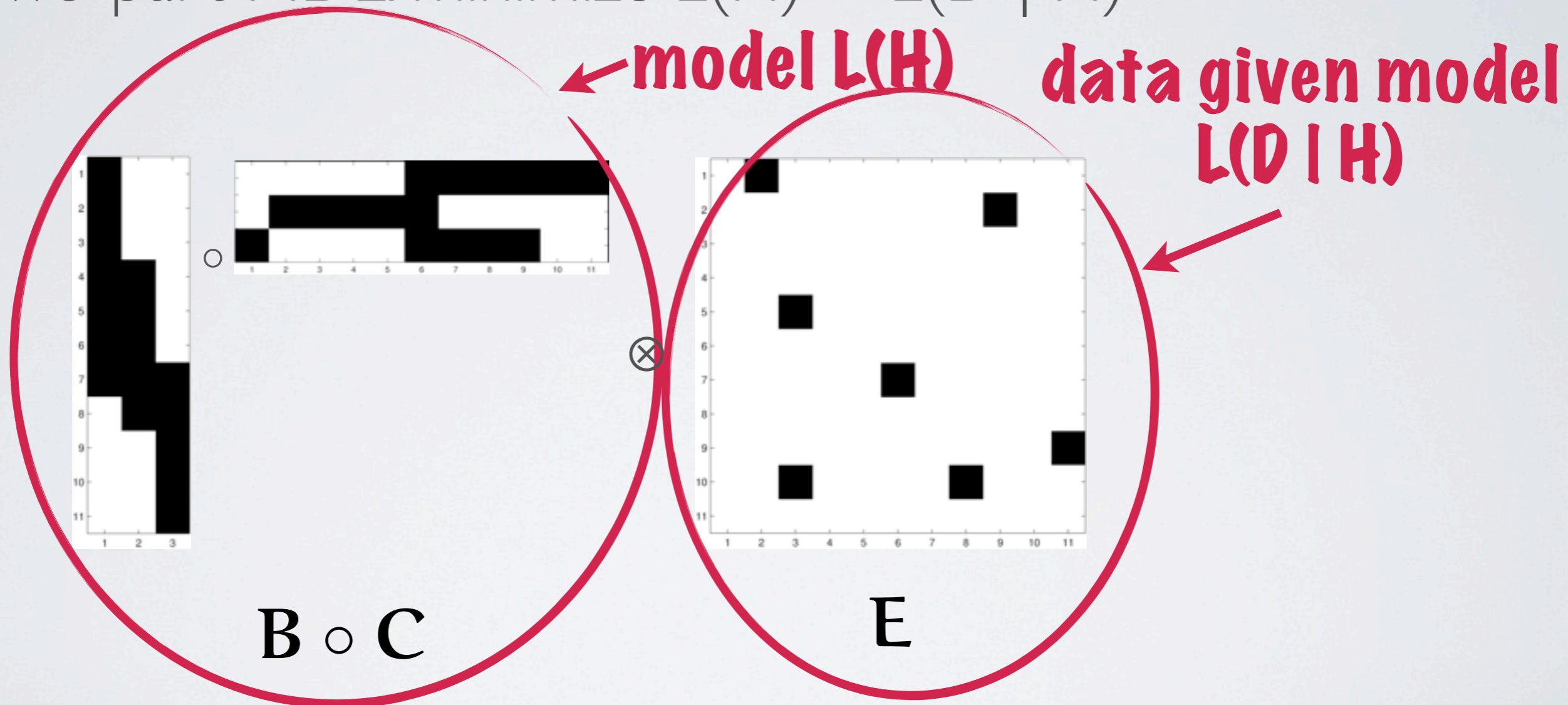


E



FITTING BMF TO MDL

- Two-part MDL: minimize $L(H) + L(D | H)$



ENCODING THE MODEL

- Model includes factor matrices **B** and **C** and their dimensions (n , m , and k)
- Each factor (row of **B** and column of **C**) is encoded using an optimal prefix code

$$H(\mathbf{B}) = k \log n - \sum_{i=1}^k \left(|\mathbf{b}_i| \log \frac{|\mathbf{b}_i|}{n} + (n - |\mathbf{b}_i|) \log \frac{n - |\mathbf{b}_i|}{n} \right)$$



HOW HARD CAN IT BE?

- MDL itself is an approximation of Kolmogorov complexity
- Finding minimum-error BMF is NP hard (even to approximate)
- But how hard it is to find the MDL-optimal decomposition?
 - Not necessarily minimum-error decomposition
 - Hardness depends on encoding
 - We know that there exists an encoding for which it is NP-hard to find the MDL-optimal decomposition



USING ASSO WITH MDL

- **The Good**

- Asso is hierarchical and deterministic
 - The k^{th} factor does not change the previous $k - 1$ factors

- **The Bad**

- Asso is heuristic

- **The Ugly**

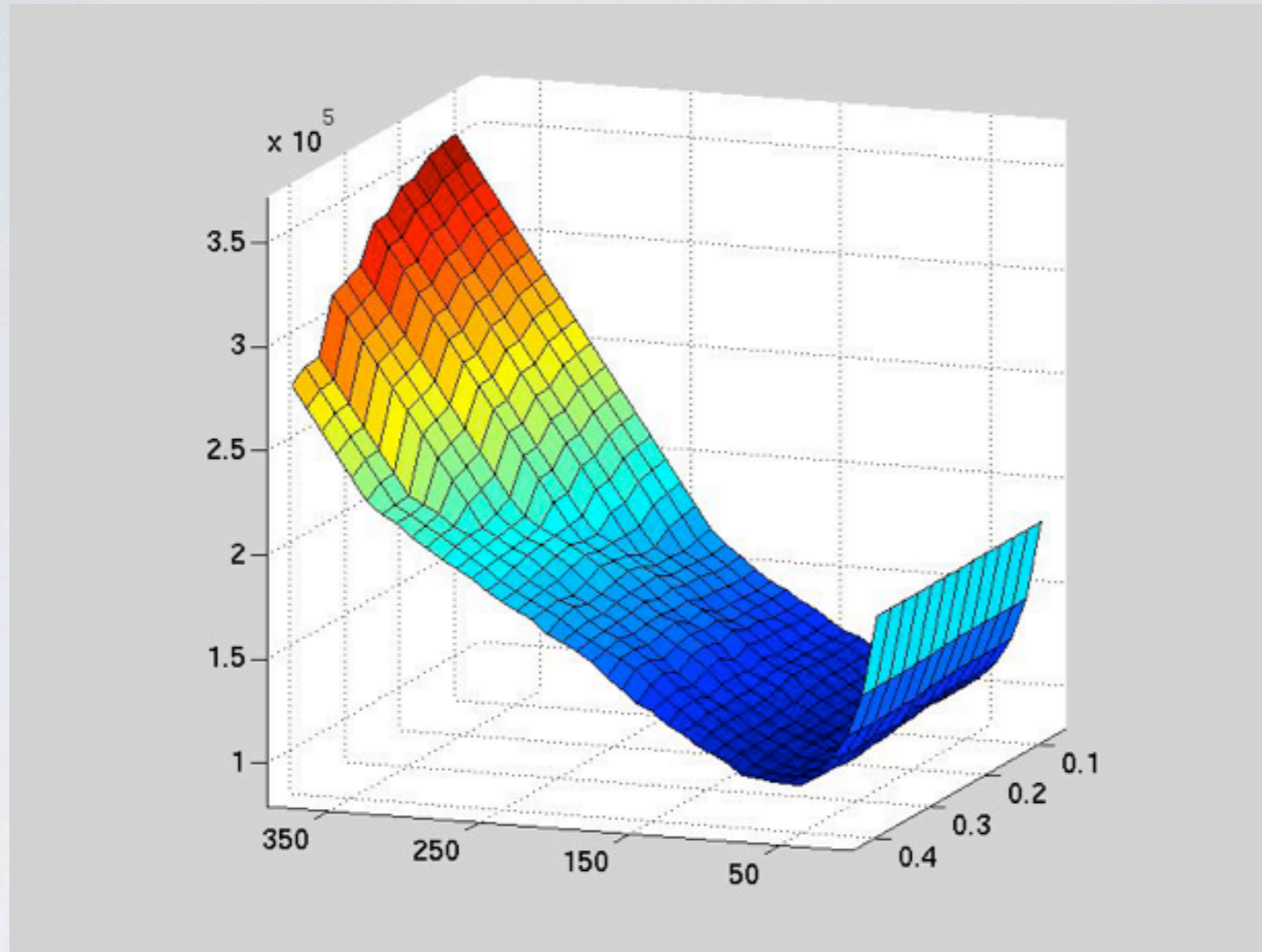
- Asso requires extra parameter t — but MDL can be used to find this, too



HASN'T THIS BEEN DONE BEFORE?

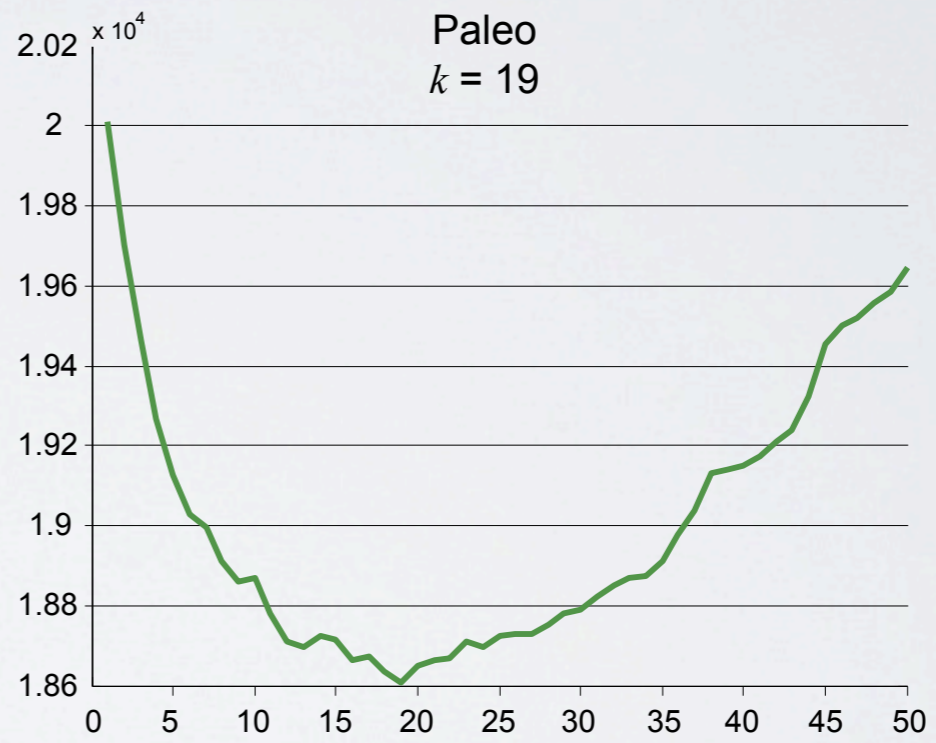
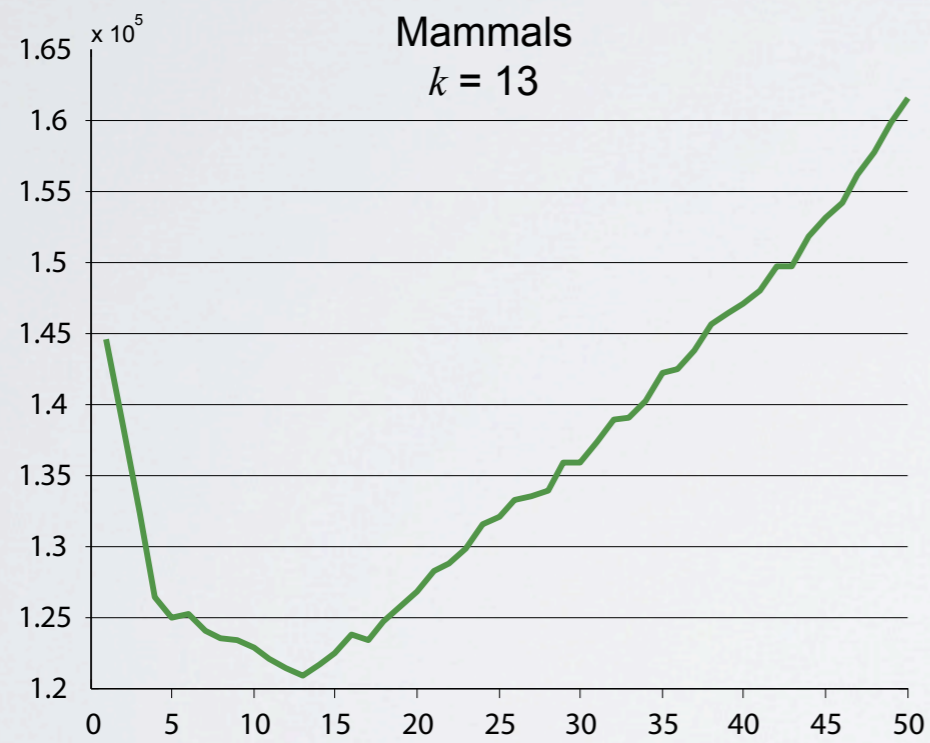
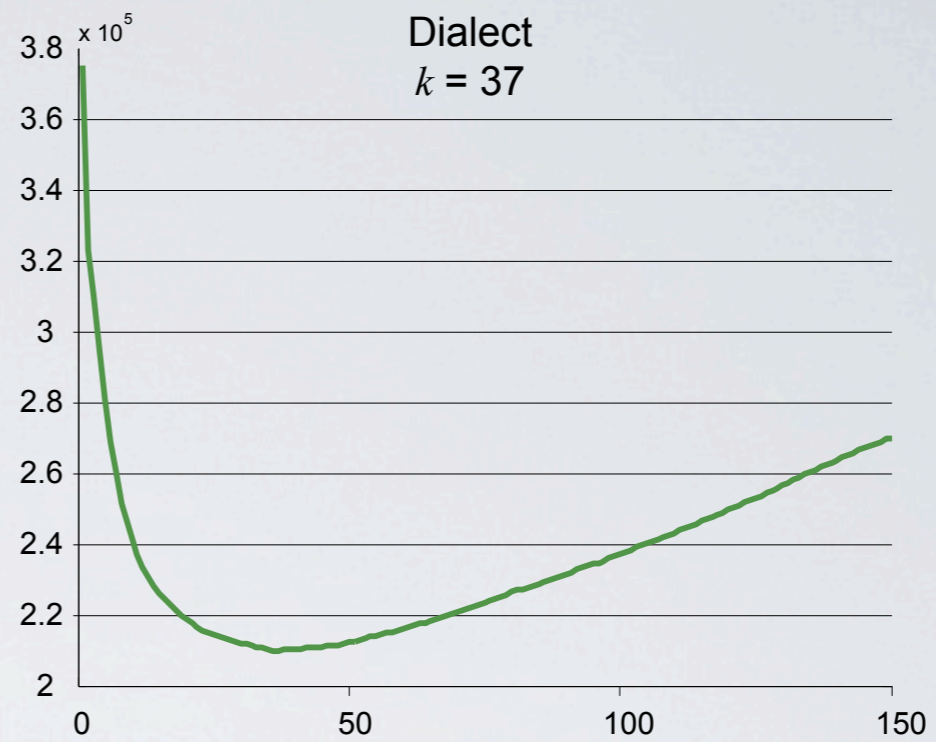
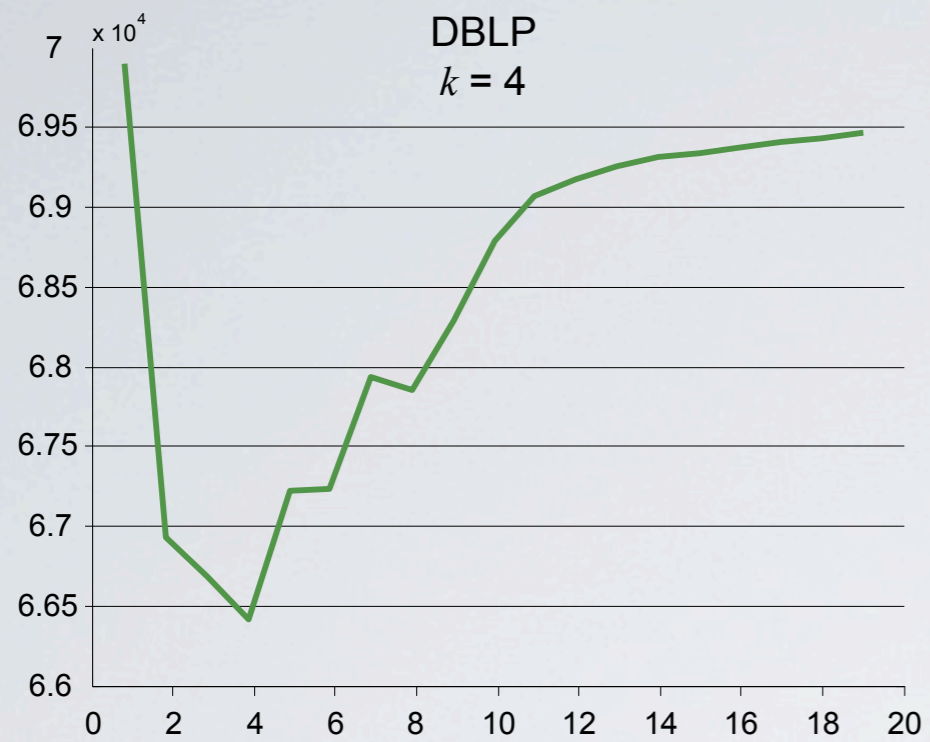
- Model order selection for matrix factorizations is studied before (mostly with SVD/PCA)
- Methods such as Guttman–Kaiser criterion (c. 1950) or Cattell's scree test (1966) are not suitable
 - Poor performance and need for subjective decisions
- Cross validation doesn't work, either
 - Well-known problem with matrix factorizations





THE DNA DATA





REAL-WORLD DATA



FUTURE WORK

- Binary tensors
- Maximize similarity vs. minimize dissimilarity
- Solve BMF via LP optimization
 - And better algorithms in general
- Joint subspaces
-



CONCLUSIONS

- BMF is a strong data mining technique
 - If your data are binary, consider BMF
- Computationally hard, but sparsity helps
- Model order selection can be solved with MDL
 - Irrespective of algorithm used
- Lots of things to do...

