

On tensor rank of conditional probability tables in Bayesian networks



Jirka Vomlel and Petr Tichavský

ÚTIA, Academy of Sciences of the Czech Republic, Prague

Tensor (an object of multilinear algebra)

Tensor is a mapping $\mathcal{A} : \mathbb{I} \rightarrow \mathbb{R}$, where

- ▶ $\mathbb{I} = I_1 \times \dots \times I_k$,
- ▶ k is a natural number called the order of tensor \mathcal{A} ,
- ▶ $I_j, j = 1, \dots, k$ are index sets (typically, sets of integers),
- ▶ cardinalities of $I_j, j = 1, \dots, k$ are called dimensions of tensor \mathcal{A} .

Example

Example of a tensor of order $k = 4$ and dimensions $n_1 = n_2 = n_3 = n_4 = 2$:

$$\mathcal{A} = \begin{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} & \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \\ \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} & \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \end{pmatrix}$$

Tensor rank

Tensor \mathcal{A} has **rank one** if it can be written as an outer product of vectors $\mathbf{a}_j = (a_{j,i})_{i \in I_j}, j = 1, \dots, k$:

$$\mathcal{A} = \mathbf{a}_1 \otimes \dots \otimes \mathbf{a}_k,$$

with the outer product \otimes being defined as

$$\mathcal{A}_{i_1, \dots, i_k} = a_{1, i_1} \cdot \dots \cdot a_{k, i_k},$$

for all $(i_1, \dots, i_k) \in I_1 \times \dots \times I_k$.

Each tensor can be decomposed as a linear combination of rank-one tensors:

$$\mathcal{A} = \sum_{i=1}^r b_i \cdot \mathbf{a}_{i,1} \otimes \dots \otimes \mathbf{a}_{i,k},$$

The **rank of a tensor** \mathcal{A} , denoted $\text{rank}(\mathcal{A})$, is the minimal r over all such decompositions.

The decomposition of a tensor \mathcal{A} to tensors of rank one that sum up to \mathcal{A} is called **CP tensor decomposition**.

Example

The tensor \mathcal{A} from the above example can be written as:

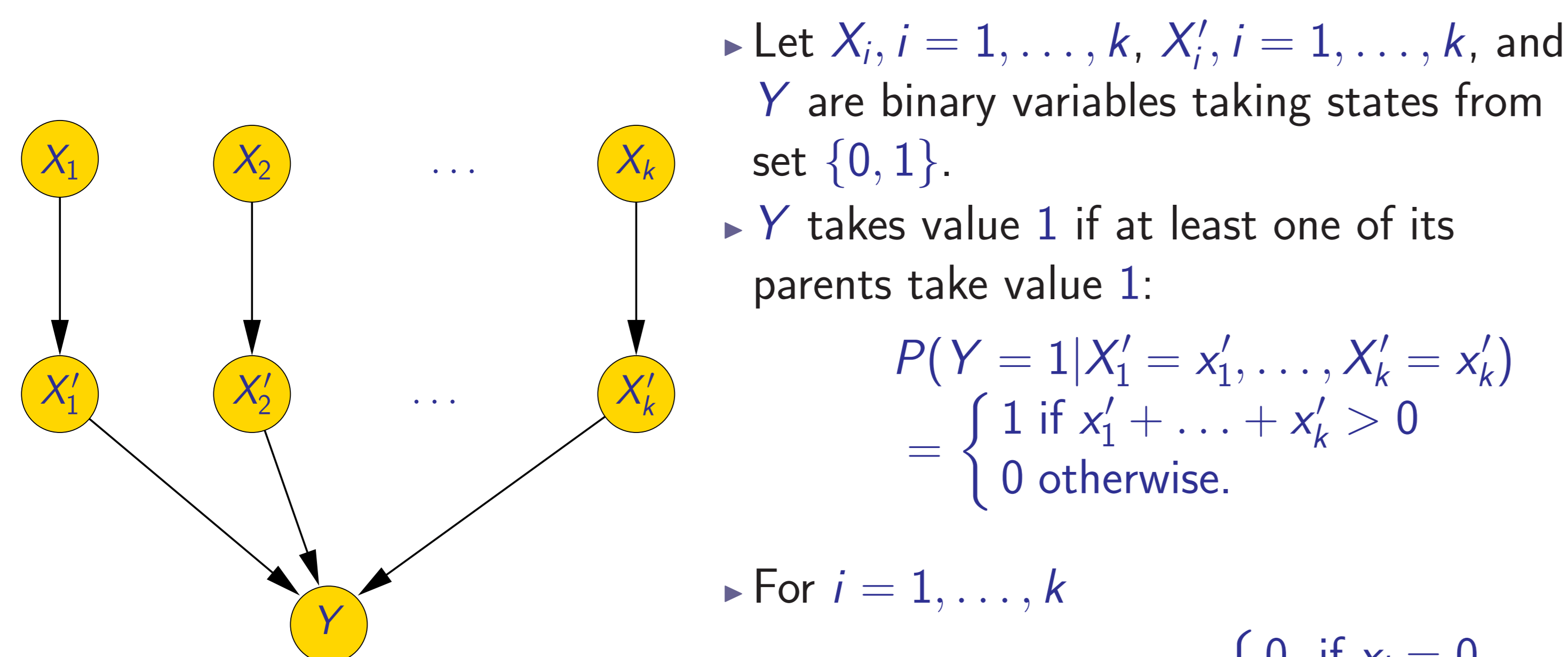
$$\mathcal{A} = (1, 1) \otimes (1, 1) \otimes (1, 1) \otimes (1, 1) - (1, 0) \otimes (1, 0) \otimes (1, 0) \otimes (1, 0).$$

This implies that its rank is at most 2.

Conditional probability tables (CPTs) as tensors

- ▶ Let $X_i, i = 1, \dots, k$ be discrete random variables taking states from $I_j, j = 1, \dots, k$ and
 - ▶ let Y be a discrete random variable with an observed state y .
- Then each CPT $P(Y = y | X_1, \dots, X_k)$ can be viewed as a **tensor**.

Noisy-or: an example of a CPT of a low rank



Example

The CPT of $P(Y = 1 | X'_1 = x'_1, \dots, X'_k = x'_k)$ corresponds to the tensor \mathcal{A} from the above example.

Probabilistic inference

Compute $P(X_j | e)$ where e is the collected evidence.

Example

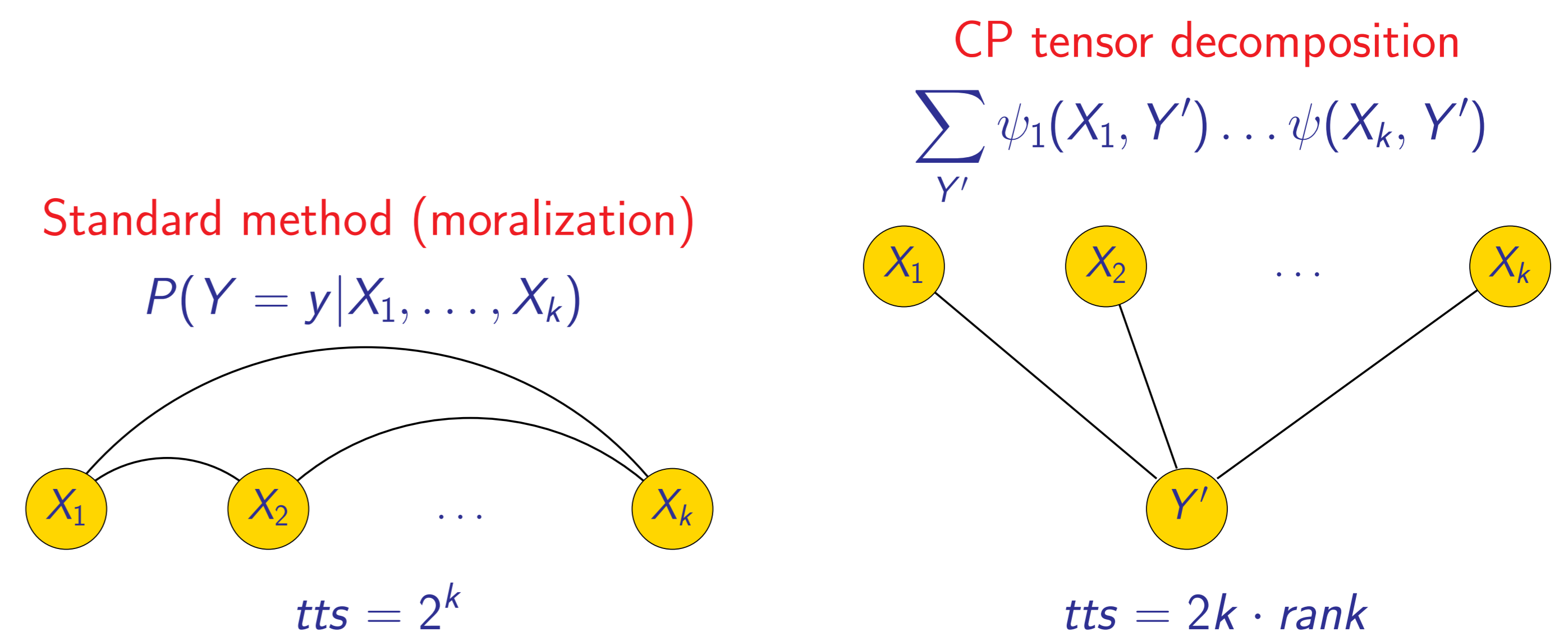
Compute $P(X_j | Y = y)$ for $j = 1, \dots, 4$.

Low rank means efficient inference

A measure of the complexity of probabilistic inference is the **total table size** (tts). It is the total size of the tables used in a probabilistic inference method applied to a Bayesian network.

Example

Bayesian network with one noisy-or:

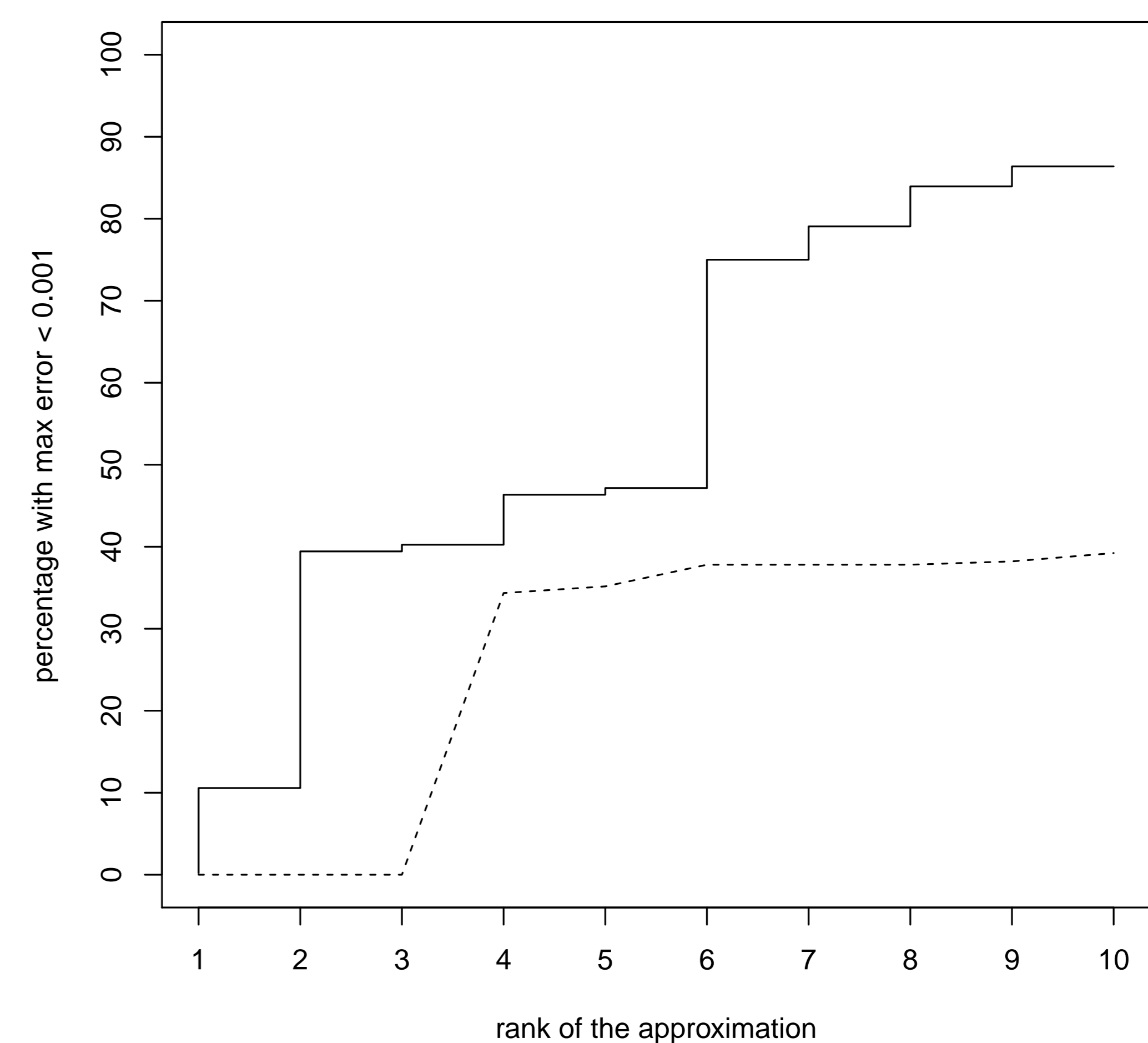


Real data and algorithms used

- ▶ 15 Bayesian networks from a Bayesian network repository: *alarm.net, hailfinder.net, barley2.net, pathfinder.net, munin1.net, munin2.net, munin3.net, munin4.net, mildew.net, hepar2.net, andes.net, win95pts.net, water.net, link.net, and insurance.net.*
- ▶ In the experiments we considered all CPTs for variables having at least three parents. It was 492 CPTs in total.
- ▶ For the CP tensor decomposition we used the Fast Damped Gauss-Newton (Levenberg-Marquard) algorithm implemented in the Matlab package TENSORBOX.

How common are low rank CPTs in real applications?

Percentage of CPTs that can be approximated with a maximum error smaller than 0.001 as a function of the rank of the approximation.



The full line corresponds to CPTs of real Bayesian networks, dashed line to CPTs of the same dimensions with their entries replaced by random numbers.

Conclusions

- ▶ Most of the conditional probability tables can be **very well approximated by tables that have lower rank** than one would expect from a general table of the same dimensions.
- ▶ The low rank approximation should be exploited (1) in **model elicitation** by using a compact parametrization of the internal structure of a CPT but also (2) during **probabilistic inference**.
- ▶ We conjecture that some low rank approximations of CPTs may actually correspond better to what a domain expert intended to model in the constructed CPT.