

# Sparse Coding and Automatic Relevance Determination for Multi-way models

Morten Mørup and Lars Kai Hansen  
 DTU Informatics  
 Technical University of Denmark  
 2800 Kgs. Lyngby  
 Email: {mm.lkh}@imm.dtu.dk  
 Telephone: +45 4525 3900  
 Fax: +45 4587 2599

**Abstract**—Multi-way modeling has become an important tool in the analysis of large scale multi-modal data. An important class of multi-way models is given by the Tucker model which decomposes the data into components pertaining to each modality as well as a core array indicating how the components of the various modalities interact. Unfortunately, the Tucker model is not unique. Furthermore, establishing the adequate model order is difficult as the number of components are specified for each mode separately. Previously, rotation criteria such as VARIMAX has been used to resolve the non-uniqueness of the Tucker representation [7]. Furthermore, all potential models have been exhaustively evaluated to estimate the adequate number of components of each mode. We demonstrate how sparse coding can prune excess components and resolve the non-uniqueness of the Tucker model while Automatic Relevance Determination in Bayesian learning form a framework to learn the adequate degree of sparsity imposed. On a wide range of multi-way data sets the proposed method is demonstrated to successfully prune excess components thereby establishing the model order. Furthermore, the non-uniqueness of the Tucker model is resolved since among potential models the models giving the sparsest representation as measured by the sparse coding regularization is attained. The approach readily generalizes to regular sparse coding as well as the CandeComp/PARAFAC model as both models are special cases of the Tucker model.

## I. INTRODUCTION

Tensor decompositions are in frequent use today in a variety of fields including psychometric, image analysis, web data mining, bio-informatics, neuroimaging and signal processing [8]. While the analysis of tensors have been somewhat restricted due to memory and computational limitations growing attention has lately been given to the analysis of data with tensorial structure fueled by the increased memory capacity and computational power of modern computers. Tensors, i.e.,  $\mathcal{X} \in \mathfrak{R}^{I_1 \times I_2 \times \dots \times I_N}$ , also called multi-way arrays, multidimensional matrices or hypermatrices are generalizations of vectors (first order tensors) and matrices (second order tensors). The two most commonly used decompositions of tensors are the Tucker model [20] and the more restricted CandeComp/PARAFAC (CP) model [5].

The Tucker model reads

$$\mathcal{X}_{i_1, i_2, \dots, i_N} = \mathcal{R}_{i_1, i_2, \dots, i_N} + \mathcal{E}_{i_1, i_2, \dots, i_N} \\ = \sum_{j_1 j_2 \dots j_N} \mathcal{G}_{j_1, j_2, \dots, j_N} a_{i_1, j_1}^{(1)} a_{i_2, j_2}^{(2)} \dots a_{i_N, j_N}^{(N)} + \mathcal{E}_{i_1, i_2, \dots, i_N}.$$

where  $\mathcal{G} \in \mathfrak{R}^{J_1 \times J_2 \times \dots \times J_N}$  and  $\mathbf{A}^{(n)} \in \mathfrak{R}^{I_n \times J_n}$  while  $\mathcal{E}$  is the approximation error. To indicate how many vectors pertain to each modality it is customary also to denote the model a Tucker( $J_1, J_2, \dots, J_N$ ). Using the n-mode matricizing operation [8] the model can also be expressed as

$$\mathbf{X}_{(n)} = \mathbf{A}^{(n)} \mathbf{Z}_{(n)} + \mathbf{E}_{(n)},$$

where

$$\mathbf{Z}_{(n)} = \mathbf{G}_{(n)} (\mathbf{A}^{(1)} \otimes \dots \otimes \mathbf{A}^{(n-1)} \otimes \mathbf{A}^{(n+1)} \otimes \dots \otimes \mathbf{A}^{(N)})^\top.$$

Furthermore, using the n-mode tensor product  $\times_n$  [10] given by

$$(\mathcal{Q} \times_n \mathbf{P})_{i_1, i_2, \dots, j_n, \dots, i_N} = \sum_{i_n} \mathcal{Q}_{i_1, i_2, \dots, i_n, \dots, i_N} \mathbf{P}_{j_n, i_n},$$

the model is stated as

$$\mathcal{X} = \mathcal{G} \times_1 \mathbf{A}^{(1)} \times_2 \mathbf{A}^{(2)} \times_3 \dots \times_N \mathbf{A}^{(N)}.$$

Thus, the Tucker model represents the data spanning the  $n^{\text{th}}$  modality by the vectors (loadings) given by the  $J_n$  columns of  $\mathbf{A}^{(n)}$  such that the vectors of each modality interact with the vectors of all remaining modalities with strengths given by a so-called core tensor  $\mathcal{G}$ . As a result, the Tucker model encompasses all possible linear interactions between vectors pertaining to the various modalities of the data. The CP model is a special case of the Tucker model where the size of each modality of the core array  $\mathcal{G}$  is the same, i.e.,  $J_1 = J_2 = \dots = J_N$  while interaction is only between columns of same indices such that the only non-zero elements are found along the diagonal of the core, i.e.,  $\mathcal{G}_{j_1, j_2, \dots, j_N} \neq 0$  iff  $j_1 = j_2 = \dots = j_N$ . Thus, the CP model can by appropriate scaling of each component be expressed as a Tucker model with unit diagonal core, i.e.  $\mathcal{G}_{CP} = \mathcal{I}$ . As such, the regular sparse coding model based on a factor analysis type representation can be formulated as a 2-way Tucker model with diagonal core. Notice, in the Tucker model a rotation of a given loading matrix  $\mathbf{A}^{(n)}$  can be compensated by a counter rotation of the core  $\mathcal{G}$  [8]. For the CP model it is not possible in general to rotate the loadings and still keep the core diagonal. Thus, the CP model is unique up to scale and permutation [9].

As the CP model corresponds to the Tucker model with diagonal core – Tucker decompositions in which only some off diagonal elements are non-zero can be considered a representational interpolation between the Tucker and CP decomposition. Regularizing the core and loadings of the Tucker model by sparse priors using a sparse coding approach it becomes possible to simplify the core and turn of excess components. Furthermore, such restrictions on the core can potentially eliminate the rotation degeneracy of a Tucker decomposition. Thus, by adequately controlling the degree of pruning we can select the model order and simplify the core at the cost of estimating a conventional multi-way model.

We will use a standard approach in Bayesian inference referred to as Automatic Relevance Determination (ARD) [1], [17]. Traditionally, ARD has been based on Gaussian priors yielding a ridge regression type of selection. Here, we will derive an ARD approach based on the Laplace prior. Contrary to Gaussian priors, the Laplace prior favors sparse representations. Optimizing for sparse representation is related to the classic rotation criteria such as VARIMAX [6] and maximum Likelihood independent component analysis (ICA) based on sparse priors [13]. However, rather than rotating an estimated solution, the estimation process is directly posed as a tradeoff between simplicity of the representation and fitting the data. Thus, a sparse representation is strongly related to the principle of parsimony, i.e., among all possible accounts the simplest is considered the best [13]. If no formal prior information is given parsimony can be considered a reasonable guiding principle to avoid overfitting, see also [13] and references therein. In the present paper we describe the ARD approach based on sparse priors on the Tucker model. For a full analysis of the presented framework with comparison to existing model order heuristics as well as comparison between sparse and Gaussian priors see [?].

## II. SPARSE CODING AND AUTOMATIC RELEVANCE DETERMINATION FOR MULTI-WAY MODELS

Automatic Relevance Determination is a hierarchical Bayesian approach widely used for model selection [1], [17]. In ARD hyperparameters explicitly represents the relevance of different features by defining the range of variation for these features, usually by modeling the width of a zero-mean Gaussian prior imposed on the model parameters. If the width becomes zero, the corresponding feature cannot have any effect on the predictions. Hence, ARD optimizes these hyperparameters to discover which features are relevant. While ARD based on Gaussian priors can prune excess components Gaussian priors do not in general admit sparse representation within the active components hence does not necessarily favor simple parsimonious representations. However, the Laplace prior is known to admit sparse representation as it corresponds to a  $l_1$  regularization thus is the closest convex proxy to minimizing for the number of non-zero elements in the model [4]. Therefore, we consider Laplace priors on the model parameter  $\theta_d$ , i.e.  $P_{Laplace}(\theta_d|\alpha_d) = \prod_j \frac{\alpha_d}{2} \exp[-\alpha_d|\theta_{j,d}|]$ . In a Bayesian framework, the least squares objective  $SSE = \|\mathcal{X} - \mathcal{R}\|_F^2 =$

$\sum_{i_1, i_2, \dots, i_n} (\mathcal{X}_{i_1, i_2, \dots, i_n} - \mathcal{R}_{i_1, i_2, \dots, i_n})^2$ , corresponds to minimizing the negative log-likelihood assuming the entries in  $\mathcal{X}$  are independent, identically distributed (i.i.d.) with Gaussian noise, i.e.  $P(\mathcal{X}|\mathcal{R}, \sigma^2) = (2\pi\sigma^2)^{-\frac{I_1 I_2 \dots I_N}{2}} \exp[-\frac{\|\mathcal{X} - \mathcal{R}\|_F^2}{2\sigma^2}]$ .

Assigning Laplace priors for the loadings and core the posterior likelihood can be written as

$$L = P(\mathcal{G}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}|\mathcal{X}, \sigma^2, \alpha^{\mathcal{G}}, \alpha^{(1)}, \dots, \alpha^{(N)}) \\ \propto P(\mathcal{X}|\mathcal{R}, \sigma^2)P(\mathcal{G}|\alpha^{\mathcal{G}})P(\mathbf{A}^{(1)}|\alpha^{(1)}) \dots P(\mathbf{A}^{(N)}|\alpha^{(N)}).$$

Thus the negative log likelihood based on Laplace priors is proportional to

$$-\log L \propto c + \frac{1}{2\sigma^2} \|\mathcal{X} - \mathcal{R}\|_F^2 + \sum_n \sum_{j_n} \alpha_{j_n}^{(n)} |\alpha_{j_n}^{(n)}|_1 + \alpha^{\mathcal{G}} |\mathcal{G}|_1 \\ + \frac{1}{2} I_1 I_2 \dots I_N \log \sigma^2 - \sum_n \sum_{j_n} I_n \log \alpha_{j_n}^{(n)} - J_1 J_2 \dots J_n \log \alpha^{\mathcal{G}}.$$

Where  $c$  is a constant. Notice, how first line corresponds to the regular  $l_1$ -regularized least squares (sparse coding) problem when alternatingly solving for the loadings of each mode keeping the remaining loadings fixed. Thus, each alternating subproblem has the form

$$\arg \min_{\mathbf{A}^{(n)}} \frac{1}{2} \|\mathbf{X}_{(n)} - \mathbf{A}^{(n)} \mathbf{Z}_{(n)}\|_F^2 + \lambda_d \sum_j |\alpha_{j,d}|.$$

The normalization constants in the likelihood terms are given in the second line. It is due to these normalization terms that it is possible to learn the values of  $\sigma^2$ ,  $\alpha^{(n)}$  and  $\alpha^{\mathcal{G}}$ .

To solve the sparse coding problem of each alternating step we used the simple gradient based procedure given in Algorithm 1. In table IV it is demonstrated that this approach is very efficient for undercomplete representations, i.e.  $I_n \geq J_n$ , which is normally the case for the Tucker decomposition. The approach readily generalizes to non-negativity constrained optimization by truncating negative values to zero (i.e. using projected gradient) [12]. In Algorithm 2 the algorithm for sparse ARD Tucker is given. For further details consult the Matlab implementation available for download at [www.mortenmorup.dk](http://www.mortenmorup.dk). In general the crux of the ARD approach is that it estimates an optimal tradeoff between optimizing the likelihood of the data and the likelihood of the model parameters. Since  $\sigma^2, \alpha^{(n)}$  and  $\alpha^{\mathcal{G}}$  weights the importance of the likelihood of the data and model parameters in the objective respectively - good estimates of these parameters are the crux for the ARD approach to work well. Finally, the better the noise model as well as component priors fit the true structure of the data the better the ARD framework will work. Since estimating  $\sigma^2$  from the data has a tendency of underestimating the value of  $\sigma^2$  due to over-fitting, i.e. the models ability to fit noise we used the following more viable approach described in [?] to set  $\sigma^2$  based on the assumption that the modelled signal ( $\mathcal{R}$ ) and noise ( $\mathcal{E}$ ) are uncorrelated,

$$\sigma^2 = \|\mathcal{X}\|_F^2 / (I_1 I_2 \dots I_n (1 + 10^{\text{SNR}/10})).$$

where SNR is a user defined signal to noise ratio. In all the experiments we used a fixed value of SNR = 0dB assuming

**Algorithm 1** Gradient Based Sparse Coding (GBSC):  $\mathbf{A} = \text{GBSC}(\mathbf{X}, \mathbf{Z}, \lambda)$ , solves  $\arg \min_{\mathbf{A}} \frac{1}{2} \|\mathbf{X} - \mathbf{AZ}\|_F^2 + \sum_j \lambda_j |\mathbf{a}_j|_1$

---

```

1: repeat
2:   Take gradient step according to LS-objective
3:    $\mathbf{A}^{new} \leftarrow \mathbf{A}^{old} - \mu(\mathbf{AZ} - \mathbf{X})\mathbf{Z}^\top$ 
4:   Take gradient step according to  $l_1$ -regularization
5:   if  $|a_{i,j}^{new}| < \mu\lambda_j$  then
6:      $a_{i,j}^{new} = 0$ 
7:   else
8:      $a_{i,j}^{new} = a_{i,j}^{old} - \mu\lambda_j \text{sign}(a_{i,j}^{old})$ 
9:   end if
10:  Estimate  $\mu$  by line-search
11: until Convergence

```

---

the same degree of signal as noise in the data. In [?] the sensitivity of this parameter to the obtained decomposition was investigated and it was found that the parameter had little impact for conservative choices of SNR.

### III. ARD TUCKER ANALYSIS OF MULTI-WAY DATA

In figure 2 is given the estimated cores obtained by the Sparse ARD Tucker algorithm on the following five datasets: **Synthetic Data:** A data set with Tucker(3,4,5) structure was randomly generated with size  $30 \times 40 \times 50$ . All the factors as well as the core array were drawn from a normal  $N(0,1)$ -distribution, i.e. with zero mean and variance 1. Gaussian i.i.d. noise was added to the data such that  $\text{SNR} = 0\text{dB}$ .

**Flow Injection Analysis:** This data set is described in [14], [19] and is given by the absorption spectra over time for three different chemical analytes measured in 12 samples with different concentrations, i.e.  $12(\text{samples}) \times 100(\text{wavelengths}) \times 89(\text{times})$ , ideally this dataset form a Tucker(3,6,4) model.

**Amino Acid Fluorescence:** This data set is described in [3] and contains the excitation and emission spectra of five samples of different amounts of tyrosine, tryptophane and phenylalanine forming a  $5(\text{samples}) \times 51(\text{excitation}) \times 201(\text{emission})$  array. Hence the data can be described by a three component CP model.

**Sugar process data:** This data set contain emission and excitation spectra measurements in 265 samples forming a  $265(\text{samples}) \times 571(\text{emissions}) \times 7(\text{excitations})$  array [2]. The data was in [2] modeled by a four component CP model where the number of components were estimated based on an extensive split half analysis.

**Dorrit fluorescence data:** This data set contains the emission and excitation spectra of 27 synthetic samples containing different concentrations of four chemical analytes forming a  $27(\text{samples}) \times 551(\text{emissions}) \times 24(\text{excitations})$  array [18]. The data is adequately modeled by a four component CP model.

Since the components of the four chemometrics data sets are non-negative the estimated models for these data were constrained to be non-negative.

Clearly, regularization has both removed excess components and reduced the non-zero elements in the core, for a detailed

comparison to existing methods for estimating the Tucker and CP model order see [?].

### IV. TUNING THE PRUNING IN REGULAR SPARSE CODING

The proposed ARD approach readily generalize to the regular sparse coding model proposed in [16]. Which corresponds to the Tucker model given by  $\mathbf{X} = \mathbf{A}^{(1)}(\mathbf{GA}^{(2)})^\top = \mathbf{A}^{(1)}\mathbf{IS}^\top = \mathbf{A}^{(1)}\mathbf{S}^\top$ . As in sparse coding we will impose sparsity on  $\mathbf{A}^{(1)}$  henceforth denoted  $\mathbf{A}$  while requiring that the energy of each component  $\|\mathbf{s}_j\|_F = 1$ . Hence, we minimize the following objective

$$-\log L \propto c + \frac{1}{2\sigma^2} \|\mathbf{X} - \mathbf{AS}^\top\|_F^2 + \sum_j \alpha_j |\mathbf{a}_j|_1 + \frac{1}{2} I_1 I_2 \log \sigma^2 - \sum_j I_1 \log \alpha_j. \quad \text{s.t.} \quad \|\mathbf{s}_j\|_F = 1,$$

by alternately solving for  $\mathbf{A}$ ,  $\mathbf{S}$  and  $\alpha$ , i.e. such that  $\alpha_j = \frac{I_1}{|\mathbf{a}_j|_1}$  with  $\sigma^2 = \|\mathbf{X}\|_F^2 / (I_1 I_2 (1 + 10^{0/10}))$ . For details on the implementation see the ARDSC.M Matlab implementation available for download at [www.mortenmorup.dk](http://www.mortenmorup.dk).

The result obtained by analyzing the natural images described in [16] is given in figure 3. 20 components have been extracted that well correspond to the Gabor like simple cell receptive field properties reported in [16].

### V. CONCLUSION:

Model selection is perhaps one of the most challenging problems in unsupervised learning. We demonstrated how sparse coding and a simple Bayesian framework based on Automatic Relevance Determination could be adapted to the Tucker model. Sparsity enabled to prune excess components while the ARD framework enabled to learn the adequate degree of sparsity. Since the CP model and the regular sparse coding model can be considered the n-way and 2-way case of the Tucker model with diagonal core the proposed framework readily generalizes to these model.

### ACKNOWLEDGEMENT

This research was supported by the European Commission through the EU FP6 NEST Pathfinder grant PERCEPT (043261).

### REFERENCES

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, August 2006.
- [2] R. Bro. Exploratory study of sugar production using fluorescence spectroscopy and multi-way analysis. *Chemom. Intell. Lab. Syst.*, 46:133–147, 1999.
- [3] Rasmus Bro. Parafac: Tutorial and applications. *Chemometrics and Intelligent Laboratory Systems*, 38:149–171, 1997.
- [4] David Donoho. For most large underdetermined systems of linear equations the minimal  $l^1$ -norm solution is also the sparsest solution. *Communications on Pure and Applied Mathematics*, 59(6):797–829, 2006.
- [5] R. A. Harshman. Foundations of the PARAFAC procedure: Models and conditions for an "explanatory" multi-modal factor analysis. *UCLA Working Papers in Phonetics*, 16:1–84, 1970.
- [6] H. F. Kaiser. The varimax criterion for analytic rotation in factor analysis. *Psychometrika*, 23:187–200, 1958.

---

**Algorithm 2** Sparse Tucker estimation based on Automatic Relevance Determination (ARD)

---

- 1: set  $J_1, J_2, \dots, J_n$  large enough to encompass all potential models,  $\sigma^2 = \|\mathcal{X}\|_F^2 / (I_1 I_2 \cdots I_n (1 + 10^{\text{SNR}/10}))$ , set  $\alpha_G = 0$ ,  $\alpha^{(n)} = \mathbf{0}$  and initialize by random  $\mathbf{A}^{(n)}$  for  $n = 1, 2, \dots, N$
  - 2: **repeat**
  - 3:  $\mathbf{Q} = \mathbf{A}^{(1)} \otimes \mathbf{A}^{(2)} \otimes \dots \otimes \mathbf{A}^{(N)}$ ,  $\text{vec}(\mathcal{G}) \leftarrow \text{gbsc}(\text{vec}(\mathcal{X}), \mathbf{Q}, \sigma^2 \alpha_G)$ ,  $\alpha_G = \min\{\frac{J_1 J_2 \cdots J_N}{|\mathcal{G}|_1}, \frac{1}{\epsilon}\}$
  - 4:  $\mathcal{R} = \mathcal{G} \times_1 \mathbf{A}^{(1)} \times_2 \mathbf{A}^{(2)} \times_3 \dots \times_N \mathbf{A}^{(N)}$
  - 5: **for**  $n=1:N$  **do**
  - 6:  $\mathbf{Z}_{(n)} = (\mathcal{R} \times_n \mathbf{A}^{(n)\dagger})_{(n)}$ ,  $\mathbf{A}^{(n)} \leftarrow \text{gbsc}(\mathbf{X}_{(n)}, \mathbf{Z}_{(n)}, \sigma^2 \alpha^{(n)})$ ,  $\alpha_d^{(n)} = \min\{\frac{J_n}{|\mathbf{A}_d^{(n)}|_1}, \frac{1}{\epsilon}\}$
  - 7: **If**  $\alpha_{j_n}^{(n)} = \frac{1}{\epsilon}$  **then**  $J_n = J_n - 1$ ,  $\mathbf{A}^{(n)} = \mathbf{A}_{\setminus j_n}^{(n)}$ ,  $\mathcal{G} = \mathcal{G}_{\setminus j_n}$ ,  $\alpha^{(n)} = \alpha_{\setminus j_n}^{(n)}$  **end**
  - 8:  $\mathbf{R}_{(n)} = \mathbf{A}^{(n)} \mathbf{Z}_{(n)}$
  - 9: **end for**
  - 10: **until** convergence
- 

	100 × 256	256 × 256	1000 × 256	2500 × 256
SIGNSEARCH	0.0750 ± 0.0359	0.1984 ± 0.1342	<b>0.3734 ± 0.1759</b>	<b>1.6969 ± 0.6441</b>
CONJUGATE GRADIENT	0.4172 ± 0.0651	1.1219 ± 0.2560	9.0297 ± 1.8055	45.6297 ± 12.0142
LARS	0.0453 ± 0.0226	<b>0.1313 ± 0.0787</b>	0.4313 ± 0.1477	1.9813 ± 0.6342
NNQP	0.5703 ± 0.0696	0.9313 ± 0.0748	2.8719 ± 0.1389	15.5047 ± 0.7882
GBSC	<b>0.0125 ± 0.0066</b>	0.3172 ± 0.2121	2.0688 ± 1.0760	22.8828 ± 12.2846

TABLE I

COMPARISON OF THE CPU TIME USAGE FOR VARIOUS SPARSE CODING ALGORITHMS ON DIFFERENT PROBLEM SIZES. THE CONJUGATE GRADIENT ALGORITHM WAS OBTAINED FROM [www.ll-magic.org](http://www.ll-magic.org), WHEREAS THE NON-NEGATIVE QUADRATIC PROGRAMMING METHOD (NNQP) AND LARS METHOD WERE OBTAINED FROM [www.sparselab.stanford.edu](http://www.sparselab.stanford.edu). THE SIGNSEARCH ALGORITHM WAS KINDLY PROVIDED BY H. LEE [11]. THE PROBLEM SOLVED IS  $\arg \min_{\mathbf{a}} \frac{1}{2} \|\mathbf{x}^\top - \mathbf{a}\mathbf{Z}\|_F^2 + \lambda \|\mathbf{a}\|_1$ , FOR  $\lambda = 0.05$ .  $J \times I$  DENOTES THE SIZE OF  $\mathbf{Z}$ . THE MEAN AND STANDARD DEVIATION IS GIVEN FOR 10 RANDOMLY GENERATED PROBLEMS, EACH GIVEN BY SETTING  $\mathbf{Z}$  TO  $J$  RANDOMLY CHOSEN IMAGE PATCHES FROM THE NATURAL IMAGES DATA SET GIVEN IN [15] AND  $\mathbf{x}$  TO A RANDOMLY SELECTED IMAGE PATCH, NOT ALREADY USED IN THE DICTIONARY,  $\mathbf{Z}$ . NOTICE, SIGNSEARCH AND LARS FIND THE GLOBAL OPTIMUM. THE REMAINING ALGORITHMS WERE STOPPED, WHEN THEIR DEVIATION FROM THE TRUE MINIMUM WAS LESS THAN  $10^{-4}$ . FOR  $J \leq I$ , THE GBSC IS THE FASTEST OF ALL THE ALGORITHMS, BUT FOR OVER-COMPLETE PROBLEMS, I.E.,  $J \gg I$ , THE GBSC ALGORITHM IS NOT IN GENERAL AS EFFECTIVE AS THE ALGORITHMS WHICH USE HESSIAN INFORMATION. THEREFORE GBSC IS ATTRACTIVE FOR THE ARD TUCKER MODEL ESTIMATION GIVEN IN ALGORITHM 2 SINCE IN GENERAL  $J \leq I$ .

- [7] Henk A.L. Kiers. Joint orthomax rotation of the core and component matrices resulting from three-mode principal component analysis. *Journal of Classification*, 15:245–263, 1998.
- [8] T.G. Kolda and B.W. Bader. Tensor decompositions and applications. *SIAM Review*, to appear, 2008.
- [9] J.B. Kruskal. Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics. *Linear Algebra Appl.*, 18:95–138, 1977.
- [10] Lieven De Lathauwer, Bart De Moor, and Joos Vandewalle. Multilinear singular value decomposition. *SIAM J. MATRIX ANAL. APPL.*, 21(4):1253–1278, 2000.
- [11] H. Lee, A. Battle, R. Raina, and A. Y. Ng. Efficient sparse coding algorithms. In *Proceedings of the Neural Information Processing Systems (NIPS)*, 19, 2007.
- [12] C.-J. Lin. Projected gradient methods for non-negative matrix factorization. *Neural Computation*, 19:2756–2779, 2007.
- [13] M. Mørup. *Decomposition Methods for Unsupervised Learning*. PhD thesis, Technical University of Denmark, 2008.
- [14] L Nørgaard and C. Ridder. Rank annihilation factor analysis applied to flow injection analysis with photodiode-array detection. *Chemometrics and Intelligent Laboratory Systems*, 23(1):107–114, 1994.
- [15] B. A. Olshausen and D. J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381:607–609, 1996.
- [16] Bruno A. Olshausen and David J. Field. Sparse coding of sensory inputs. *Current Opinion in Neurobiology*, 14:481–487, 2004.
- [17] Yuan (Alan) Qi, Thomas P. Minka, Rosalind W. Picard, and Zoubin Ghahramani. Predictive automatic relevance determination by expectation propagation. In *ICML '04: Proceedings of the twenty-first international conference on Machine learning*, page 85, New York, NY, USA, 2004. ACM.
- [18] Jordi Riu and R Bro. Jack-knife estimation of standard errors and outlier detection in parafac models. *Chemometrics and Intelligent Laboratory Systems*, 65(1):35–49, 2003.
- [19] Age K. Smilde, Roma Tauler, Javier Saurina, and Rasmus Bro. Calibration methods for complex second-order data. *Analytica Chimica Acta*, 398:237–251, 1999.
- [20] L. R. Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31:279–311, 1966.

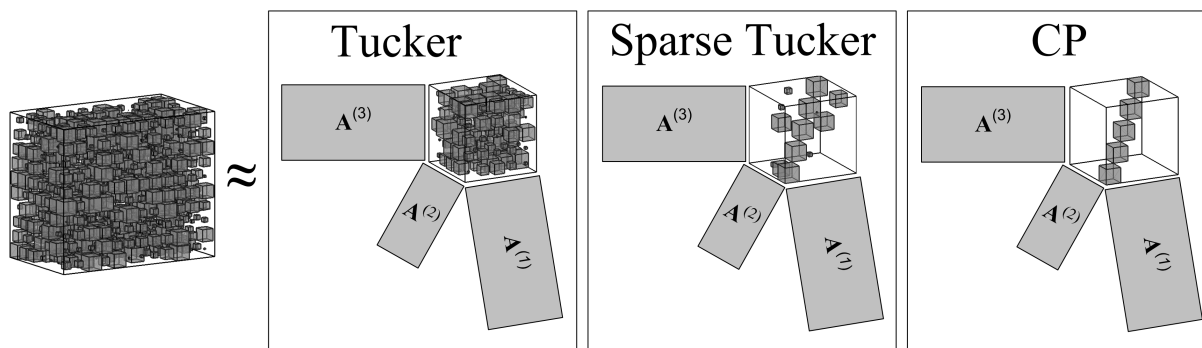


Fig. 1. Illustration of Tucker decomposition of a 3-way tensor given to the left. The Tucker model is given to the left, sparse Tucker model in the middle and CP model to the right. Whereas the Tucker model encompass all potential interaction between the components of each modality through the core array  $\mathcal{G}$ , the CP model only allow for interactions between columns of  $\mathbf{A}^{(n)}$  with same indices. The sparse Tucker model can be considered a model between the Tucker and CP model where interactions are present within a few of the components across the various modalities by imposing sparsity on the core. We will impose sparsity on the core and loadings to prune excess components while estimate the adequate degree of sparsity using a Bayesian approach named Automatic Relevance Determination (ARD).

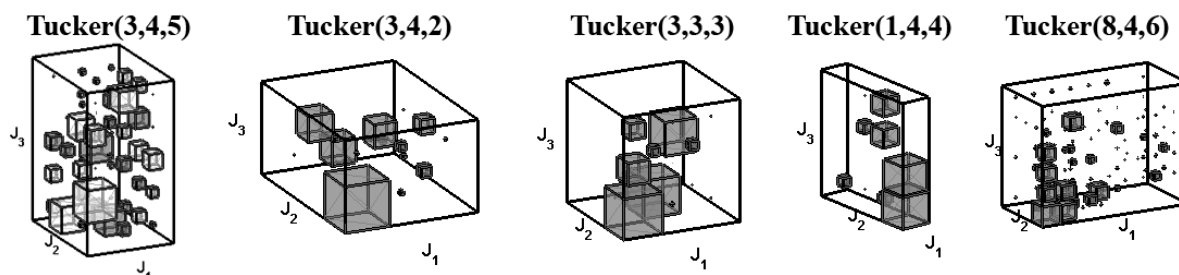


Fig. 2. The estimated cores for five different multi-way datasets – a Tucker(3,4,5) synthetic dataset as well as four 3-way chemometrics data sets obtained from [www.models.kvl.dk/research/data/](http://www.models.kvl.dk/research/data/). A Tucker(10,10,10) model was fitted to all datasets but using sparse coding to prune excess components and ARD to update the  $l_1$ -regularization strengths on the model parameters the models were reduced to form simpler models. Given are the estimated model orders and cores (size of boxes indicate interaction strengths such that gray boxes denote positive interactions and white boxes negative interactions). The estimated model orders relate well to the expected model order based on the number of true chemical compounds in the data.

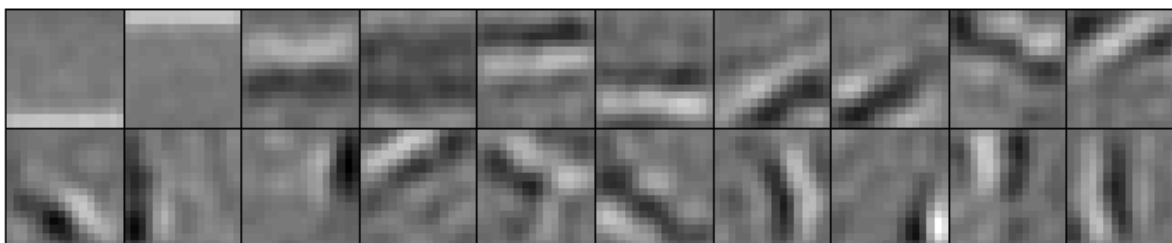


Fig. 3. Analysis of the regular sparse coding problem using the proposed ARD approach to tune the pruning parameters. 250 components were fitted to the data however all but 20 components were turned off by the proposed ARD framework. The components correspond well to the Gabor like simple cell receptive field properties reported in [16].