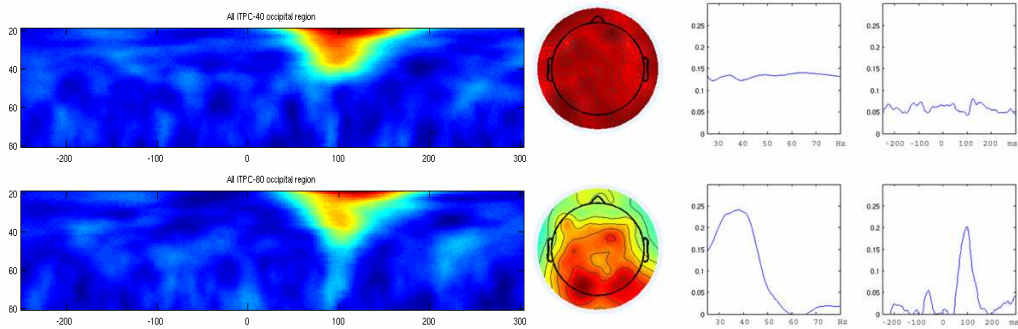
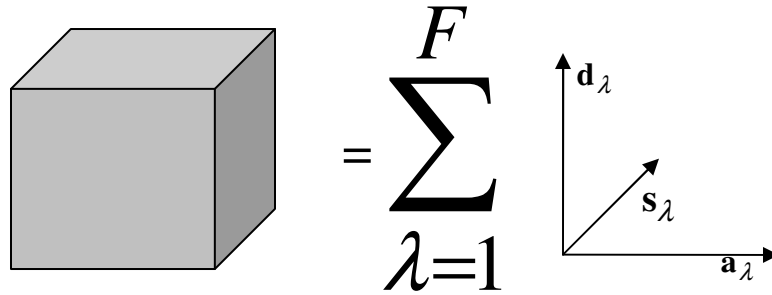
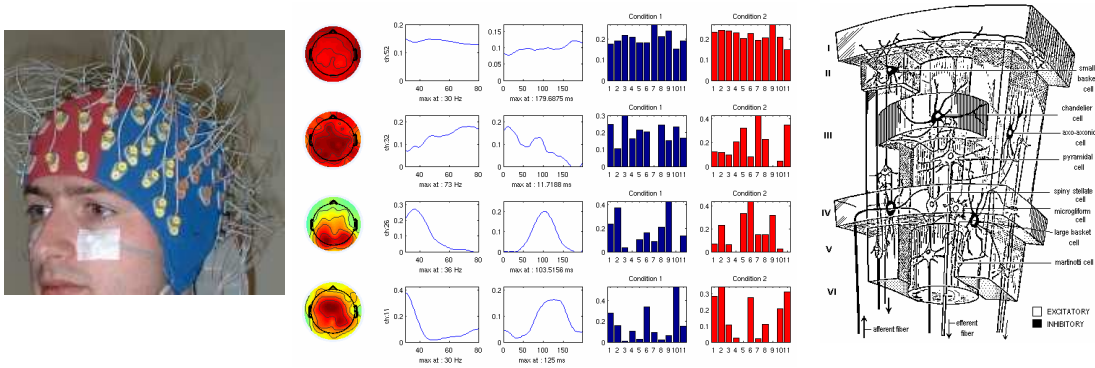


Analysis of Brain Data

Using Multi-Way Array Models on the EEG

Written by
Morten Mørup



Preface

This master thesis serves as documentation for the final assignment in the requirements to achieve the degree Master of Science in Engineering. The work has been carried out in the period from the 1st of September 2004 to the 14th of February 2005 at the Intelligent Signal Processing group at the Institute of Informatics and Mathematical Modeling, Technical University of Denmark. The work has been supervised by Professor Lars Kai Hansen and part of the work has been done in corporation with Sidse Arnfred at Cognitive Research Unit, Department of Psychiatry, Hvidovre Hospital. I wish to thank both for great discussions, enthusiasm and guidance during the project.

Kgs. Lyngby, February 14, 2005

Morten Mørup, s012198

Abstract

In this thesis the multi-way array model Parallel Factors (PARAFAC) also known as Canonical Decomposition (CANDECOMP) was applied to the event related potential (ERP) of electroencephalographic (EEG) recordings. Previous work done analyzing the ERP by PARAFAC had encountered great problems of degeneracy in the factors. However, in this thesis it is shown that the problem of degeneracy can be effectively circumvented by imposing non-negativity. Furthermore, the PARAFAC analysis was, to my knowledge, for the first time used to analyze the wavelet transformed data of the ERP. Through this analysis, it was shown that PARAFAC is able to access the correct components of the data. Finally, a novel PARAFAC algorithm based on independent component analysis on data having the concept of Combined Independence was proposed. This algorithm proved both fast and efficient in accessing the correct components of simulated as well as real data. In dealing with noise, the algorithm performed even better than the popular PARAFAC algorithm based on alternating least squares.

***Keywords:** PARAFAC, CANDECOMP, ERP, EEG, coherence, ITPC, multi-way arrays, tensors, Independent Component Analysis, gamma activity, wavelet analysis, Combined Independence, HOSVD, TUCKER, Core Consistency Diagnostic.*

Abstrakt

I denne afhandling blev multi-way array modellen Parallel Factors (PARAFAC) også kendt som Canonical Decomposition (CANDECOMP) brugt til at analysere event relaterede potentialer (ERP) fra elektroencefalografiske optagelser (EEG). Tidligere forsøg på at analysere ERP'et ved hjælp af PARAFAC havde stødt på problemer med degeneration i faktorerne. I denne afhandling blev det vist, at dette problem kan løses ved at indføre ikke-negativitets begrænsninger. Derudover blev PARAFAC, så vidt jeg ved, for første gang brugt til at analysere wavelet-transformeret ERP-data. Det viste sig, at PARAFAC også her var i stand til at finde de rigtige komponenter i data. Endelig blev en ny PARAFAC algoritme baseret på independent component analysis foreslået til brug på data havende begrebet Combined Independence. Denne algoritme viste sig både hurtig og effektiv til at finde de korrekte komponenter i simuleret såvel som rigtig data. Algoritmen var endda bedre til at håndtere støj end den populære PARAFAC algoritme baseret på alternating least squares.

***Nøgleord:** PARAFAC, CANDECOMP, ERP, EEG, coherence, ITPC, multi-way arrays, tensors, Independent Component Analysis, Gamma activity, wavelet analysis, Combined Independence, HOSVD, TUCKER, Core Consistency Diagnostic.*

Table of Contents

Notation.....	6
Abbreviations.....	8
Introduction.....	9
1 PARAFAC.....	10
1.1 Wavelet Analysis.....	10
1.2 Bayesian Learning.....	12
1.2.1 The Expectation Maximization (EM) Algorithm.....	13
1.2.2 Variational Bayesian EM (VBEM) algorithm.....	18
1.3 Multi-way arrays.....	20
1.3.1 Unfolding.....	20
1.4 Models.....	22
1.4.1 Parallel Factor Analysis (PARAFAC).....	22
1.4.2 TUCKER and Higher Order Singular Value Decomposition.....	25
1.4.3 Model Relations.....	30
1.5 PARAFAC Algorithms.....	31
1.5.1 Core Consistency Diagnostic.....	33
1.5.2 PARAFAC by Alternating Least Squares.....	34
1.5.3 PARAFAC by multi-way rank one decomposition.....	35
1.5.4 PARAFAC by EM and VBEM.....	36
1.5.5 PARAFAC combined with ICA.....	40
1.5.6 Algorithm relations.....	42
2 Electroencephalography recording (EEG).....	44
2.1 Dipoles.....	44
2.2 EEG and coherence.....	48
2.3 Synaptic potentials and action potentials.....	51
2.3.1 Synaptic potentials.....	51
2.3.2 Action potentials.....	52
2.3.3 Synaptic Potentials versus Action Potentials.....	54
2.4 Features of the EEG and ERP.....	55
2.4.1 Event Related Potentials, ERP.....	55
2.4.2 Noise.....	59
2.4.3 EEG/ERP and PARAFAC.....	60
3 Data analysis.....	62
3.1 Simulated data.....	62
3.2 Real data.....	69
4 Discussion.....	87
5 Conclusion.....	89
References.....	91
Appendix A: Theorems with proofs.....	94
Appendix B: Multi-way array algebra.....	111
Appendix C: MATLAB implementation of multi-way array manipulations.....	113
Appendix D: ICA- and ALSPARAFAC on Chemometric Data.....	116

Notation

$\mathbf{x}, \bar{\mathbf{X}}$	Vector
\mathbf{X}	Matrix
\mathcal{X}	Multi-way array in the literature also referred to as tensors, higher-order tensors or multidimensional matrices
$x_{i_1}, x_{i_1 i_2}, x_{i_1 i_2 \dots i_N}$	Denotes the element at the subscribed indices for the corresponding vector, matrix and multi-way array.
$\begin{matrix} a \\ \parallel \\ b \end{matrix}$	Same as $a=b$
$\underset{nr}{=} , \underset{nr}{\Leftrightarrow} , \underset{nr}{\Rightarrow}$	nr refers to the index of an explanation given below the equation.
$\mathbf{X}_{(:,i)}$	The MATLAB notation in this case for the vector consisting of the i^{th} column of \mathbf{X} .
\mathbf{x}_i	Vector containing the i^{th} column of \mathbf{X} , i.e. $\mathbf{x}_i = \mathbf{X}_{(:,i)}$
$\mathbf{X}^{(i)}$	Denotes a vector, but where size $\mathbf{X}^{(i)}$ not necessarily equals $\mathbf{X}^{(j)}$
\mathbf{X}^+	The pseudo inverse of \mathbf{X}
$\text{vec}(\mathbf{X})$	The vectorization of \mathbf{X} given by: $\text{vec}\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_N \end{bmatrix}$
$\text{diag}(\cdot)$	On a matrix: Sets off diagonal elements of a matrix to zero. On a vector: Creates the diagonal matrix where the vector elements are along the diagonal.
$\ \cdot\ $	The Frobenius-norm, see Definition 3, page 111
$\mathbf{A} \otimes \mathbf{B}$	The tensor or Kronecker product, i.e. $\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \dots & a_{1m}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{n1}\mathbf{B} & \dots & a_{nm}\mathbf{B} \end{bmatrix}$
$\mathbf{A} \otimes \mathbf{B}$	The Khatri-Rao product: $\mathbf{A} \otimes \mathbf{B} = [\mathbf{a}_1 \otimes \mathbf{b}_1 \quad \dots \quad \mathbf{a}_f \otimes \mathbf{b}_f]$ Requires that \mathbf{A} and \mathbf{B} have same number of columns.
$\mathbf{a} \circ \mathbf{b}$	The outer product, i.e. $\mathbf{C} = \mathbf{a} \circ \mathbf{b} \Leftrightarrow c_{ij} = a_i b_j$
$\mathbf{X}^{I \times K}$	Unfolding \mathcal{X} by the j^{th} -way onto the i^{th} way, i.e. for a 3-way array $\mathbf{X}^{I \times K} = \begin{bmatrix} \mathbf{X}_{(:,1,:)} \\ \vdots \\ \mathbf{X}_{(:,J,:)} \end{bmatrix}$
$\mathbf{X}^{I \times JK}$	Unfolding \mathcal{X} by the K^{th} -way onto the J^{th} way, i.e. for a 3-way array

	$\mathbf{X}^{I \times JK} = [\mathbf{X}_{(:, :, 1)} \quad \dots \quad \mathbf{X}_{(:, :, K)}]$
$\mathbf{X}_{(n)}$	The n-mode-matricizing of the multi-way array $\mathbf{X} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_N}$ to the matrix $\mathbf{X}_{(n)} \in \mathcal{R}^{I_n \times I_1 \dots I_{n-1} I_{n+1} \dots I_N}$, the inverse operation is denoted $\mathbf{X}_{(n)}^{-1}$
$\mathbf{A} \times_n \mathbf{B}$	The n-mode Multiplication, see also Definition 1, page 111.
$\langle \mathbf{A}, \mathbf{B} \rangle$	The scalar product of two tensors, see also Definition 2, page 111.
$rank(\mathbf{A})$	The rank of the multi-way array \mathbf{A} , see also Definition 4, page 111.
$k_{\mathbf{A}}$	The k-rank of the matrix \mathbf{A} , see also Definition 6, page 112.
$\mathcal{N}(\mathbf{s}_i \mathbf{0}, \mathbf{I})$	The vector \mathbf{s}_i is normal distributed with mean $\mathbf{0}$ and covariance \mathbf{I} .
$\dim(\mathbf{X})$	The number of dimensions of \mathbf{X} , i.e. if $\mathbf{X} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_N}$ then $\dim(\mathbf{X}) = N$
$\langle \cdot \rangle$	The expected value of x , i.e. $\langle x \rangle = \int xp(x)dx$
$tr(\cdot)$	The sum of the diagonal elements of a matrix
$cov(\mathbf{x}, \mathbf{y})$	The covariance of \mathbf{x} and \mathbf{y} .
$\mathbf{A} \bullet \mathbf{B}$	The Hadamard product (element wise product of two matrices)
$ \mathbf{A} $	The determinant of the matrix \mathbf{A} .
∇, ∇^2	$\nabla = \begin{bmatrix} \frac{\partial}{\partial x} & \frac{\partial}{\partial y} & \frac{\partial}{\partial z} \end{bmatrix}$ $\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} & \frac{\partial f}{\partial z} \end{bmatrix}$ <p>$\nabla \cdot \mathbf{f}$ - the divergence of \mathbf{f}: $\frac{\partial f}{\partial x} + \frac{\partial f}{\partial y} + \frac{\partial f}{\partial z}$</p> <p>$\nabla^T \times \mathbf{f}$ - the curl of \mathbf{f} where \times is the cross product.</p> <p>$\nabla^2 f$ - the Laplacian of f: $\frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} + \frac{\partial^2 f}{\partial z^2}$</p>
\exists	Denotes there exists, i.e. $\exists i; \ \mathbf{x}_i\ = 1$ means that there is at least one column of \mathbf{X} having the norm 1.
\forall	Denotes for all variable, i.e. $\ \mathbf{x}_i\ = 1 \forall i$ means that the norm of each column of \mathbf{X} is 1.
ε	Infinitesimal small value
$\mathbf{e}, \mathbf{E}, E$	The model error

Abbreviations

ALS	Alternating Least Squares
ARD	Automatic Relevance Determination
CANDECOMP	CANonical DECOMPosition (same as PARAFAC)
CCD	Core Consistency Diagnostic
CI _{m,n}	Combined Independent in the m th and n th dimension
BIC	Bayesian Information Criterion
ECG	Electro-Cardiogram
EEG	Electroencephalography
EM	Expectation Maximization
EOG	Electro-Oculogram
EP	Evoked Response Potential
EPSP	Excitatory Post Synaptic Potential
ERD	Evoked Response Desynchronization
ERP	Event Related Potentials
ERPCOH	Event Related Cross Coherence
ERSP	Event Related Spectral Perturbation
FFT	Fast Fourier Transform
fMRI	Functional Magnetic Resonance Imaging
HOSVD	Higher Order Singular Value Decomposition
ICA	Independent Component Analysis
i.i.d.	Independent and identically distributed
IPSP	Inhibitory Post Synaptic Potential
ITPC	Inter Trial Phase Coherence
KL	Kullback-Leibler divergence
LTM	Long Term Memory
MUM	Match and Utilization Model
NMF	Non-negative Matrix Factorization
PARAFAC -ALSPARAFAC -SR1PARAFAC -EMPARAFAC -VBPARAFAC -ICAPARAFAC	Parallel Factor Analysis Based on: -ALS -Sum of Rank One Components -Expectation Maximization -Variational Bayesian Expectation Maximization - PARAFAC model based on ICA
PCA	Principal Component Analysis
RE	Reticular thalamic nucleus
SNR	Signal to Noise Ratio
STFA	Short Time Fourier Analysis
SVD	Singular Value Decomposition
TCR	Thalamic Relay Cells
VBEM	Variational Bayesian Expectation Maximization

Introduction

Electroencephalography (EEG) refers to electrical activity measured at the scalp that arises from neural activity in the brain. EEG signals generated in response to sensory stimuli events are also referred to as event related potentials (ERP). Traditionally the EEG/ERP has been analyzed by looking at trial averages and spectrums. As much of the focus in the interpretation of the EEG/ERP is based on frequency changes in the data, wavelet analysis has become a popular tool. However, wavelet analysis increases the dimensionality as it adds a frequency dimension to the data giving a multi-way array of *channel* \times *time* \times *frequency*. Consequently, to be able to effectively interpret the wavelet analyzed data there is a need to decompose these multi-modal EEG/ERP data into easily interpretable components.

In this thesis multi-way array analysis of ERP will be explored. The thesis is inspired by the work of Miwakeichi and Martínez-Montes et al. [24], [25] who applied the multi-way decomposition method Parallel Factor (PARAFAC) to analyze the space-time-frequency components of the EEG. The PARAFAC model used by Miwakeichi and Martínez-Montes et al. will be compared to other PARAFAC models taken from the framework of higher order singular value decomposition (HOSVD) [19] and a more statistical framework using the expectation maximization algorithm (EM) and variational Bayesian expectation maximization (VBEM) described by Beal [3]. Finally, a PARAFAC model based on Independent Component Analysis will be proposed.

The PARAFAC algorithms will be evaluated on real as well as simulated ERP data. The real data was collected by Sidse Arnfred at Cognitive Research Unit, Department of Psychiatry, Hvidovre Hospital. The data reproduces a well known experiment described by Herrmann et al. [15] in which evoked gamma oscillations are found in the posterior regions of the brain. The PARAFAC model will be used both on the ERP as previously done by Field et al.[10], but also for the first time, to my knowledge, to analyze the wavelet transformed ERP-data in terms of the Inter Trial Phase Coherence (ITPC).

1 PARAFAC

“If the only tool you have is a hammer,
you tend to see every problem as a nail.”
Abraham H. Maslow

Before addressing the Parallel Factor (PARAFAC) model and algorithms an introduction to wavelet analysis, Bayesian learning and multi-way array algebra will first be given.

1.1 Wavelet Analysis

A wavelet analysis transforms an EEG signal of *channel × time* into a multi-way array of *channel × time × frequency*.

The spectrum of a signal $x(t)$ is given by its Fourier transform:

$$X(F) = \int_{-\infty}^{\infty} x(t)e^{-i2\pi Ft} dt \quad \text{eq. 1.1}$$

However, the Fourier transform can't reveal frequency changes through the signal. This has led to the development of the Short-Time Fourier Analysis (STFA). In STFA the signal is Fourier transformed within a finite time-window – giving a temporal resolution of the frequency components of the signal. Unfortunately, the time-window is fixed disabling good temporal resolution for high frequencies. The wavelet transform resolves this problem.

$$C(\text{scale}, \text{shift}) = \int_{-\infty}^{\infty} x(t)\varphi^*(\text{scale}, \text{shift}, t)dt \quad \text{eq. 1.2}$$

A wavelet is a waveform of effectively limited duration that has an average value of zero. Scaling a wavelet simply means stretching or compressing it, and shifting a wavelet delaying or hastening its onset. The wavelet analysis has grown to become a huge discipline in the analysis of EEG from noise reduction to feature extraction.

Wavelets are separated into continuous and discrete wavelets based on the characteristic of the wavelet rather than the signal's characteristic as is the case for the Fourier transform. A wavelet is called continuous if it can be scaled and shifted to any value. An example of a continuous wavelet is the popular complex Morlet wavelet used in [14],[15], [24]:

$$\tilde{\varphi}(t) = \frac{1}{\sqrt{\pi F_b}} \exp(i2\pi F_c t) \exp\left(-\frac{t^2}{F_b}\right) \quad \text{eq. 1.3}$$

F_c is the center frequency and F_b is a bandwidth parameter. The scaling factor a and shift factor p changes $\tilde{\varphi}$ by:

$$\varphi(a, p, t) = \frac{1}{\sqrt{a}} \tilde{\varphi}\left(\frac{t-p}{a}\right) = \frac{1}{\sqrt{\pi F_b a}} \exp\left(i2\pi F_c \frac{(t-p)}{a}\right) \exp\left(-\frac{(t-p)^2}{a^2 F_b}\right) \quad \text{eq. 1.4}$$

$\tilde{\varphi}$ is also called the mother wavelet as it is φ without scaling and shifting. The effect of scaling is illustrated in Figure 1.1.

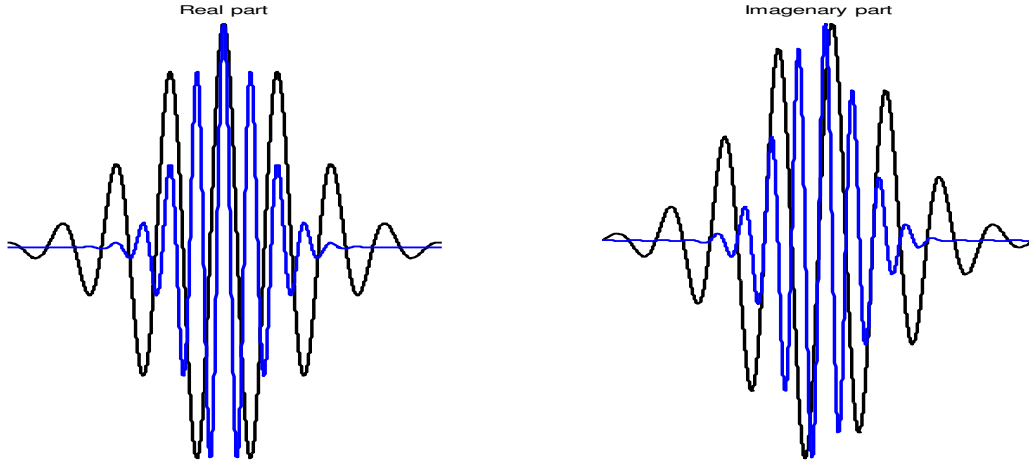


Figure 1.1: The effect of scaling the complex Morlet function. As seen scaling results in a temporal compression of the functions, black has twice the scaling factor of blue.

From the scale of the wavelet transform the frequency of the signal can be estimated as [32]:

$$F = \frac{F_c}{a} \quad \text{eq. 1.5}$$

There is an inherent tradeoff for wavelets between good frequency resolution and good time resolution. This is explained by the Heisenberg-Gabor inequality [17]. As seen from eq. 1.3 a relatively large bandwidth of the wavelet gives a good frequency resolution but the length of the wavelet makes the time point less accurate. Furthermore, there is no simple relation between center frequency and bandwidth as frequency changes with scale according to eq. 1.5 but the bandwidth changes according to eq. 1.4 by $a^2 F_b$.

Consequently, in some literature the bandwidth is denoted $\sigma_b^2 = a^2 F_b$ making $\sigma_b \propto \frac{1}{F}$.

Although the wavelet's estimate of the frequency at a given time isn't exact and the whole analysis is slightly influenced by the choice of wavelet, the wavelet analysis is considered a very powerful tool in the analysis of the temporal development of the frequency of the EEG. In the following analysis the complex Morlet wavelet having a bandwidth parameter $F_b = 2$ and a center frequency of $F_c = 1$ will be used, as it has been well accepted in the literature, see also [14],[15].

1.2 Bayesian Learning

Reverend Thomas Bayes (1702-1761) was the first recorded to notice [4]:

$$p(x, y) = p(x|y)p(y) \Rightarrow \left. \begin{array}{l} p(a, b) = p(a|b)p(b) \\ \parallel \\ p(b, a) = p(b|a)p(a) \end{array} \right\} \Rightarrow \underbrace{p(b|a) = \frac{p(a|b)p(b)}{p(a)}}_{\text{Bayes' Theorem}} \quad \text{eq. 1.6}$$

This theorem has become the cornerstone in a probabilistic modeling approach named Bayesian learning.

Given the data \mathbf{D} , the model m , and the model parameters θ the posterior probability distribution of the parameters can be expressed using Bayes' theorem as:

$$p(\theta|\mathbf{D}, m) = \frac{p(\mathbf{D}|\theta, m)p(\theta|m)}{p(\mathbf{D}|m)} \quad \text{eq. 1.7}$$

Where $p(\theta|m)$ is the prior probability of the parameters given the model. $p(\mathbf{D}|\theta, m)$ is the likelihood of the parameters also called the likelihood function. As $\int p(\theta|\mathbf{D}, m) \mathcal{D}\theta = 1$, $p(\mathbf{D}|m)$ is a normalization constant also denoted the *marginal likelihood*, given by:

$$p(\mathbf{D}|m) = \int p(\mathbf{D}|\theta, m)p(\theta|m) \mathcal{D}\theta \quad \text{eq. 1.8}$$

In probabilistic modeling the goal is to develop models that explain the given data but also generalize well on new data. In Bayesian learning this becomes the two main goals [3]:

1. Approximating the marginal Likelihood of the observed data $p(\mathbf{D}|m)$
2. Approximating the posterior distribution over the parameters of a model $p(\theta|\mathbf{D}, m)$

Consequently, in probabilistic modeling two main problems must be addressed; finding the right model and the optimal parameters. The solution will be based on the Expectation Maximization algorithm (EM) and the Variational Bayesian Expectation Maximization algorithm (VBEM) based on the analysis given by Beal [3].

The problem of finding the optimal parameters will first be addressed. In maximum likelihood learning equal priors of the parameters given the model is assumed. The maximum likelihood parameter is given as the parameter that is the most probable given the model and the data:

$$\begin{aligned}\boldsymbol{\theta}_{MAP} &= \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathbf{D}, m) \stackrel{1}{=} \arg \max_{\boldsymbol{\theta}} \frac{p(\mathbf{D}|\boldsymbol{\theta}, m)p(\boldsymbol{\theta}|m)}{p(\mathbf{D}|m)} \stackrel{2}{=} \arg \max_{\boldsymbol{\theta}} p(\mathbf{D}|\boldsymbol{\theta}, m)p(\boldsymbol{\theta}|m) \\ \boldsymbol{\theta}_{ML} &= \arg \max_{\boldsymbol{\theta}} p(\mathbf{D}|\boldsymbol{\theta}, m)\end{aligned}\tag{eq. 1.9}$$

- 1) Follows by eq. 1.7.
- 2) Result as the denominator is a constant independent of $\boldsymbol{\theta}$.
- 3) Equality holds as maximum likelihood assumes equal priors of $\boldsymbol{\theta}$ given the model. If equal priors can't be assumed, the maximum a posteriori (MAP) estimate is found instead.

Whereas the EM algorithm is based on the maximum likelihood parameter estimate $\boldsymbol{\theta}_{ML}$, the VBEM algorithm is based on the maximum a posteriori estimate $\boldsymbol{\theta}_{MAP}$.

1.2.1 The Expectation Maximization (EM) Algorithm

We consider a model having the hidden variables \mathbf{S} and the observed data \mathbf{D} . The parameters describing the (potentially) stochastic dependencies between the hidden and observed variables are given by $\boldsymbol{\theta}$. We assume further that the data $\mathbf{D} = \{\mathbf{d}_1, \dots, \mathbf{d}_n\}$ consist of n independent and identically distributed (i.i.d.) items, generated using a set of hidden variables $\mathbf{S} = \{\mathbf{s}_1, \dots, \mathbf{s}_n\}$ such that the likelihood can be written as a function of $\boldsymbol{\theta}$ in the following way [3]:

$$p(\mathbf{D}|\boldsymbol{\theta}) = \prod_{i=1}^n p(\mathbf{d}_i|\boldsymbol{\theta}) = \prod_{i=1}^n \int p(\mathbf{d}_i, \mathbf{s}_i|\boldsymbol{\theta}) \delta \mathbf{s}_i\tag{eq. 1.10}$$

The logarithm of the likelihood $\mathcal{L}(\boldsymbol{\theta})$ is defined as:

$$\mathcal{L}(\boldsymbol{\theta}) \equiv \ln p(\mathbf{D}|\boldsymbol{\theta}) = \sum_{i=1}^n \ln p(\mathbf{d}_i|\boldsymbol{\theta}) = \sum_{i=1}^n \ln \int p(\mathbf{s}_i, \mathbf{d}_i|\boldsymbol{\theta}) d\mathbf{s}_i\tag{eq. 1.11}$$

By introducing an auxiliary distribution of the hidden variables given by $q_s(\mathbf{s})$ we can find a lower bound of $\mathcal{L}(\boldsymbol{\theta})$:

$$\begin{aligned}
\mathcal{L}(\boldsymbol{\theta}) &= \sum_{i=1}^n \ln \int p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta}) d\mathbf{s}_i = \sum_{i=1}^n \ln \int q_{s_i}(\mathbf{s}_i) \frac{p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta})}{q_{s_i}(\mathbf{s}_i)} d\mathbf{s}_i \geq \int_{\mathbf{s}_i} q_{s_i}(\mathbf{s}_i) \ln \frac{p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta})}{q_{s_i}(\mathbf{s}_i)} d\mathbf{s}_i \\
&= \sum_{i=1}^n \int q_{s_i}(\mathbf{s}_i) \ln p(\mathbf{d}_i | \boldsymbol{\theta}) d\mathbf{s}_i - \int q_{s_i}(\mathbf{s}_i) \ln \frac{q_{s_i}(\mathbf{s}_i)}{p(\mathbf{s}_i | \mathbf{d}_i, \boldsymbol{\theta})} d\mathbf{s}_i \equiv \mathcal{F}(q_{s_1}(\mathbf{s}_1), \dots, q_{s_n}(\mathbf{s}_n), \boldsymbol{\theta})
\end{aligned} \tag{eq. 1.12}$$

1) Result of Jensen's inequality, see Theorem 1.

It's worth noticing:

$$\int q_{s_i}(\mathbf{s}_i) \ln \frac{q_{s_i}(\mathbf{s}_i)}{p(\mathbf{s}_i | \mathbf{d}_i, \boldsymbol{\theta})} d\mathbf{s}_i \equiv KL[q_{s_i}(\mathbf{s}_i) \| p(\mathbf{s}_i | \mathbf{d}_i, \boldsymbol{\theta})] \geq 0 \tag{eq. 1.13}$$

Is the Kullback-Leibler (KL) divergence - a measure of distance between two distributions.

The Expectation-Maximization (EM) algorithm alternates between an E step, which infers posterior distributions over hidden variables given a current parameter setting, and an M step, which maximizes $\mathcal{L}(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$ given the statistics gathered from the E step [3]:

$$\begin{aligned}
\text{Estep : } & q_{s_i}^{t+1}(\mathbf{s}_i) \leftarrow \arg \max_{q_{s_i}} \mathcal{F}(q_{s_i}(\mathbf{s}_i), \boldsymbol{\theta}^t), \forall i \in \{1, \dots, n\} \\
\text{Mstep : } & \boldsymbol{\theta}^{t+1} \leftarrow \arg \max_{\boldsymbol{\theta}} \mathcal{F}(q_{s_i}^{t+1}(\mathbf{s}_i), \boldsymbol{\theta})
\end{aligned}$$

The E step

To find the optimal distribution $q_{\mathbf{s}_i}(\mathbf{s}_i)$ we differentiate $\mathcal{F}(q_{\mathbf{s}_i}(\mathbf{s}_i), \boldsymbol{\theta}^t)$ with respect to $q_{\mathbf{s}_i}(\mathbf{s}_i)$ subject to the constraint $\int q_{\mathbf{s}_i}(\mathbf{s}_i) d_{\mathbf{s}_i} = 1$ implemented by the Lagrange multipliers $\{\lambda_i\}_{i=1}^n$, we find:

$$\frac{d\left(\mathcal{F}(q_{\mathbf{s}_i}(\mathbf{s}_i), \boldsymbol{\theta}^t) + \sum_{i=1}^n \lambda_i [\int q_{\mathbf{s}_i}(\mathbf{s}_i) d_{\mathbf{s}_i} - 1]\right)}{d\lambda_i} = \int q_{\mathbf{s}_i}(\mathbf{s}_i) d_{\mathbf{s}_i} - 1 = 0 \Leftrightarrow \int q_{\mathbf{s}_i}(\mathbf{s}_i) d_{\mathbf{s}_i} = 1$$

$$\left. \begin{aligned} \frac{d\left(\mathcal{F}(q_{\mathbf{s}_i}(\mathbf{s}_i), \boldsymbol{\theta}^t) + \sum_{i=1}^n \lambda_i [\int q_{\mathbf{s}_i}(\mathbf{s}_i) d_{\mathbf{s}_i} - 1]\right)}{dq_{\mathbf{s}_i}(\mathbf{s}_i)} &= \ln p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta}^t) - \ln q_{\mathbf{s}_i}(\mathbf{s}_i) + \lambda_i - 1 = 0 \\ &\Downarrow \\ q_{\mathbf{s}_i}^{t+1}(\mathbf{s}_i) &= \exp(\lambda_i - 1) p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta}^t) \\ \int q_{\mathbf{s}_i}(\mathbf{s}_i) d_{\mathbf{s}_i} = 1 &\Rightarrow \int \exp(\lambda_i - 1) p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta}^t) d_{\mathbf{s}_i} = 1 \\ &\Downarrow \\ \exp(\lambda_i - 1) &= \frac{1}{\int p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta}^t) d_{\mathbf{s}_i}} \end{aligned} \right\} \begin{aligned} & q_{\mathbf{s}_i}^{t+1}(\mathbf{s}_i) \\ & \parallel \\ & \exp(\lambda_i - 1) p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta}^t) \quad \text{eq. 1.14} \\ & \parallel \\ & \frac{p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta}^t)}{\int p(\mathbf{s}_i, \mathbf{d}_i | \boldsymbol{\theta}^t) d_{\mathbf{s}_i}} \\ & \parallel \\ & p(\mathbf{s}_i | \mathbf{d}_i, \boldsymbol{\theta}^t) \end{aligned}$$

The optimal choice of the distribution $q_{\mathbf{s}_i}^{t+1}(\mathbf{s}_i)$ is the posterior distribution of \mathbf{s}_i given by the data and model parameters $p(\mathbf{s}_i | \mathbf{d}_i, \boldsymbol{\theta}^t)$. This choice of $q_{\mathbf{s}_i}^{t+1}(\mathbf{s}_i)$ fulfills according to eq. 1.13 that the KL-divergence is zero.

The M step

To find the optimal parameters we make use of the result from the E step, i.e.

$$q_{s_i}^{t+1}(s_i) = p(s_i | \mathbf{d}_i, \boldsymbol{\theta}).$$

$$\begin{aligned} \mathcal{F}(q_{s_1}^{t+1}(s_1), \dots, q_{s_n}^{t+1}(s_n), \boldsymbol{\theta}) &= \sum_{i=1}^N \int q_{s_i}(s_i) \ln \frac{p(s_i, \mathbf{d}_i | \boldsymbol{\theta})}{q_{s_i}(s_i)} ds_i = \\ \sum_{i=1}^N \int p(s_i | \mathbf{d}_i, \boldsymbol{\theta}) \ln \frac{p(s_i, \mathbf{d}_i | \boldsymbol{\theta})}{p(s_i | \mathbf{d}_i, \boldsymbol{\theta})} ds_i &= \sum_{i=1}^N \int p(s_i | \mathbf{d}_i, \boldsymbol{\theta}) \ln p(\mathbf{d}_i | \boldsymbol{\theta}) ds_i = \sum_{i=1}^N \ln p(\mathbf{d}_i | \boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta}) \end{aligned} \quad \text{eq. 1.15}$$

The M step defined by $\boldsymbol{\theta}^{t+1} \leftarrow \arg \max_{\boldsymbol{\theta}} \mathcal{F}(q_{s_i}^{t+1}(s_i), \boldsymbol{\theta})$ therefore becomes a matter of maximizing the likelihood $\mathcal{L}(\boldsymbol{\theta})$.

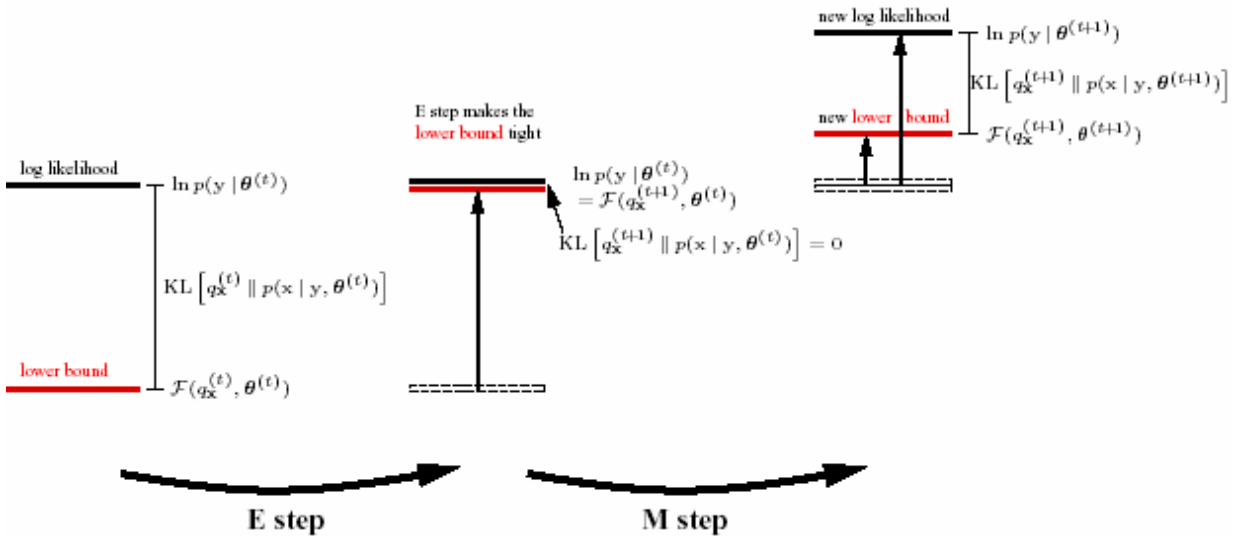


Figure 1.2: The EM algorithm for maximum likelihood learning. In the E step the hidden variable posterior is set to the exact model posterior, making the KL-divergence zero. In the M step, the lower bound of the likelihood of the parameters $\mathcal{F}(q_{s_1}^{t+1}(s_1), \dots, q_{s_n}^{t+1}(s_n), \boldsymbol{\theta})$ is maximized. (Taken from [3]).

Bayesian Information Criterion (BIC)

Finding the right model for the EM algorithm, i.e. choosing the number of source signals amounts to finding the optimal model M that according to the Bayesian Information Criterion satisfies:

$$M_{opt} = \arg \max_M p(M | \mathbf{D}) \quad \text{eq. 1.16}$$

Where $p(M | \mathbf{D})$ is given by Bayes' theorem:

$$p(M | \mathbf{D}) = \frac{p(\mathbf{D} | M)p(M)}{p(\mathbf{D})} \quad \text{eq. 1.17}$$

Where:

$$p(\mathbf{D}) = \sum_M p(\mathbf{D} | M)p(M) \quad \text{eq. 1.18}$$

$p(M)$ is the prior of the model and assumed to be uniform. For a particular choice of model, the probability of finding the observed data \mathbf{D} is given by the integral over all model parameters:

$$\begin{aligned} p(\mathbf{D} | M) &= \int p(\mathbf{D}, \boldsymbol{\theta} | M) d\boldsymbol{\theta} = \int p(\mathbf{D} | M, \boldsymbol{\theta}) p(\boldsymbol{\theta} | M) d\boldsymbol{\theta} \\ &= \int e^{\log p(\mathbf{D} | M, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta} | M)} d\boldsymbol{\theta} = \int e^{-f(\boldsymbol{\theta})} d\boldsymbol{\theta} \end{aligned} \quad \text{eq. 1.19}$$

Where

$$f(\boldsymbol{\theta}) = -\log p(\mathbf{D} | M, \boldsymbol{\theta}) - \log p(\boldsymbol{\theta} | M) \quad \text{eq. 1.20}$$

As equal priors are assumed in maximum likelihood learning the optimum of $f(\boldsymbol{\theta})$ is given by $\boldsymbol{\theta}_{ML}$. Making a second order Taylor expansion around the optimum given by $\boldsymbol{\theta}_{ML}$ yields:

$$f(\boldsymbol{\theta}) \approx f(\boldsymbol{\theta}_{ML}) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_{ML})^T \mathbf{H}(\boldsymbol{\theta} - \boldsymbol{\theta}_{ML}) \quad \text{eq. 1.21}$$

Which gives:

$$p(\mathbf{D} | M) = \int e^{-f(\boldsymbol{\theta})} d\boldsymbol{\theta} \approx e^{-f(\boldsymbol{\theta}_{ML})} \int e^{-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_{ML})^T \mathbf{H}(\boldsymbol{\theta} - \boldsymbol{\theta}_{ML})} d\boldsymbol{\theta} \quad \text{eq. 1.22}$$

As the integral has a Gaussian form it can be written by:

$$p(\mathbf{D} | M) \approx p(\mathbf{D} | \boldsymbol{\theta}_{ML}, M) p(\boldsymbol{\theta}_{ML} | M) (2\pi)^{\frac{D}{2}} |\mathbf{H}|^{-1/2} \quad \text{eq. 1.23}$$

Where D is the number of free parameters. Neither the prior $p(\boldsymbol{\theta}_{ML} | M)$ nor $(2\pi)^{\frac{D}{2}}$ depends on the number of samples N . The Hessian H holds a $D \times D$ product over samples that can be factored out as $|H| = N^D |\tilde{H}|$. Neglecting $|\tilde{H}|$ gives:

$$p(\mathbf{D} | M) \approx p(\mathbf{D} | \boldsymbol{\theta}_{ML}, M) N^{-\frac{D}{2}} \quad \text{eq. 1.24}$$

The probability of the observed data given the model is therefore the probability of the observed data given the optimal parameters of the model weighted by a function of the number of observations N and free parameters D .

1.2.2 Variational Bayesian EM (VBEM) algorithm

Once more we consider a model having the hidden variables \mathbf{S} and the observed data \mathbf{D} . The parameters describing the (potentially) stochastic dependencies between the hidden and observed variables are given by $\boldsymbol{\theta}$. We again assume that the data $\mathbf{D} = \{\mathbf{d}_1, \dots, \mathbf{d}_n\}$ consist of n independent and identically distributed (i.i.d.) items, generated using a set of hidden variables $\mathbf{S} = \{\mathbf{s}_1, \dots, \mathbf{s}_n\}$ such that the likelihood can be written as a function of \mathbf{S} and $\boldsymbol{\theta}$ in the following way [3]:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}, \mathbf{S}) &\equiv \ln p(\mathbf{D} | m) = \ln \int p(\mathbf{D}, \mathbf{S}, \boldsymbol{\theta} | m) d\boldsymbol{\theta} d\mathbf{S} = \ln \int q(\boldsymbol{\theta}, \mathbf{S}) \frac{p(\mathbf{D}, \mathbf{S}, \boldsymbol{\theta} | m)}{q(\boldsymbol{\theta}, \mathbf{S})} d\boldsymbol{\theta} d\mathbf{S} \geq_1 \\ &\int q(\boldsymbol{\theta}, \mathbf{S}) \ln \frac{p(\mathbf{D}, \mathbf{S}, \boldsymbol{\theta} | m)}{q(\boldsymbol{\theta}, \mathbf{S})} d\boldsymbol{\theta} d\mathbf{S} = \int_2 q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S}) \ln \frac{p(\mathbf{D}, \mathbf{S}, \boldsymbol{\theta} | m)}{q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S})} d\boldsymbol{\theta} d\mathbf{S} = \\ &= \mathcal{F}(q_{\mathbf{S}}(\mathbf{S}), q_{\boldsymbol{\theta}}(\boldsymbol{\theta})) \end{aligned} \quad \text{eq. 1.25}$$

1. Result of Jensens inequality, see Theorem 1.
2. Comes from the assumption that $\boldsymbol{\theta}$ and \mathbf{S} are mutually independent.

We now have:

$$\begin{aligned} \ln p(\mathbf{D} | m) \cdot \mathcal{F}(q_{\mathbf{S}}(\mathbf{S}), q_{\boldsymbol{\theta}}(\boldsymbol{\theta})) &= \ln p(\mathbf{D} | m) - \int q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S}) \ln \frac{p(\mathbf{D}, \mathbf{S}, \boldsymbol{\theta} | m)}{q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S})} d\boldsymbol{\theta} d\mathbf{S} = \\ \ln p(\mathbf{D} | m) - \int q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S}) \ln \frac{p(\mathbf{S}, \boldsymbol{\theta} | \mathbf{D}, m) p(\mathbf{D} | m)}{q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S})} d\boldsymbol{\theta} d\mathbf{S} &= \\ \ln p(\mathbf{D} | m) - \int q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S}) \ln p(\mathbf{D} | m) - q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S}) \ln \frac{q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S})}{p(\mathbf{S}, \boldsymbol{\theta} | \mathbf{D}, m)} d\boldsymbol{\theta} d\mathbf{S} &= \\ \int q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S}) \ln \frac{q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S})}{p(\mathbf{S}, \boldsymbol{\theta} | \mathbf{D}, m)} d\boldsymbol{\theta} d\mathbf{S} &= KL(q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) q_{\mathbf{S}}(\mathbf{S}) || p(\mathbf{S}, \boldsymbol{\theta} | \mathbf{D}, m)) \end{aligned} \quad \text{eq. 1.26}$$

The E and M step

From eq. 1.26 it is seen that improving the likelihood corresponds to setting the hidden variables as well as θ equal their posteriors, as this minimizes the Kullback-Leibler divergence $KL(q_{\theta}(\theta)q_S(\mathbf{S})\|p(\mathbf{S}, \theta|\mathbf{D}, m))$.

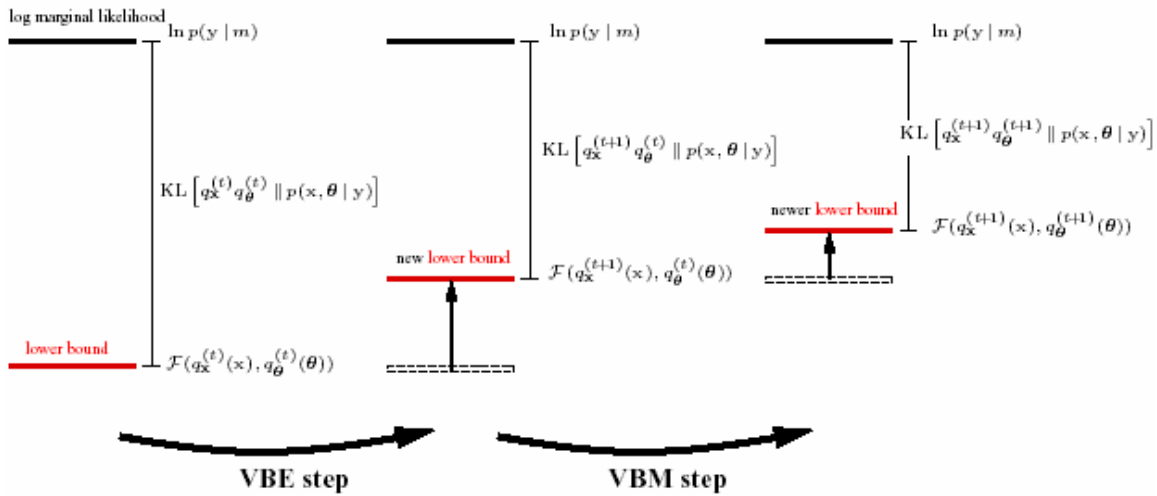


Figure 1.3: The VBEM algorithm. In the VBE step the variational posterior over hidden variable is no longer set to the exact model posterior. However, each VBE and VBM step is assured to improve the lower bound of the likelihood (Figure taken from [3]).

As the marginal likelihood doesn't change, choosing the right model amounts to finding the model having the largest lower bound of the marginal likelihood. Furthermore, hyper parameters can be used to model the various parameters of the model. The hyper parameters can then indicate how many factors to include by how certain the underlying parameter is zero. Estimating the number of factors to include in the model by hyper parameters is called Automatic Relevance Determination (ARD).

1.3 Multi-way arrays

Multi-way arrays in the literature also referred to as tensors, higher-order tensors or multidimensional matrices [19] are simply any set of data for which the elements can be arranged as [5]:

$$x_{ijk\dots} \quad i=1..I, j=1..J, k=1..K, \dots \text{ i.e. } \mathbf{X} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_n}$$

Notice that vectors and matrices are two special cases of multi-way arrays; a 1-way array and a 2-way array. In the following the various ways of a multi-way array will also be referred to as modalities. For a description of the different aspects of multi-way arrays see Appendix B: Multi-way array algebra.

1.3.1 Unfolding

The unfolding operation folds one of the ways of the multi-way data onto another. Consider for example the three-way array \mathbf{X} defined by x_{ijk} , $i=1..I, j=1..J, k=1..K$. Unfolding the third way of \mathbf{X} onto the second way gives:

$$\mathbf{X}^{I \times J \times K} \xrightarrow{\text{unfolding}} \mathbf{X}^{I \times JK}$$

While unfolding the second way of \mathbf{X} onto the third way gives:

$$\mathbf{X}^{I \times J \times K} \xrightarrow{\text{unfolding}} \mathbf{X}^{I \times KJ}$$

For a three-way array there are 6 different options of unfolding \mathbf{X} into a matrix as revealed in Figure 1.4. The unfolding can be performed consecutively turning for instance a four-way array into a vector by three unfolding operations:

$$\mathbf{X}^{I \times J \times K \times L} \xrightarrow{\text{unfolding}} \mathbf{X}^{IL \times J \times K} \xrightarrow{\text{unfolding}} \mathbf{X}^{IL \times JK} \xrightarrow{\text{unfolding}} \mathbf{x}^{ILJK}$$

Unfolding multi-way data enables manipulation of the data using normal vector and matrix calculation.

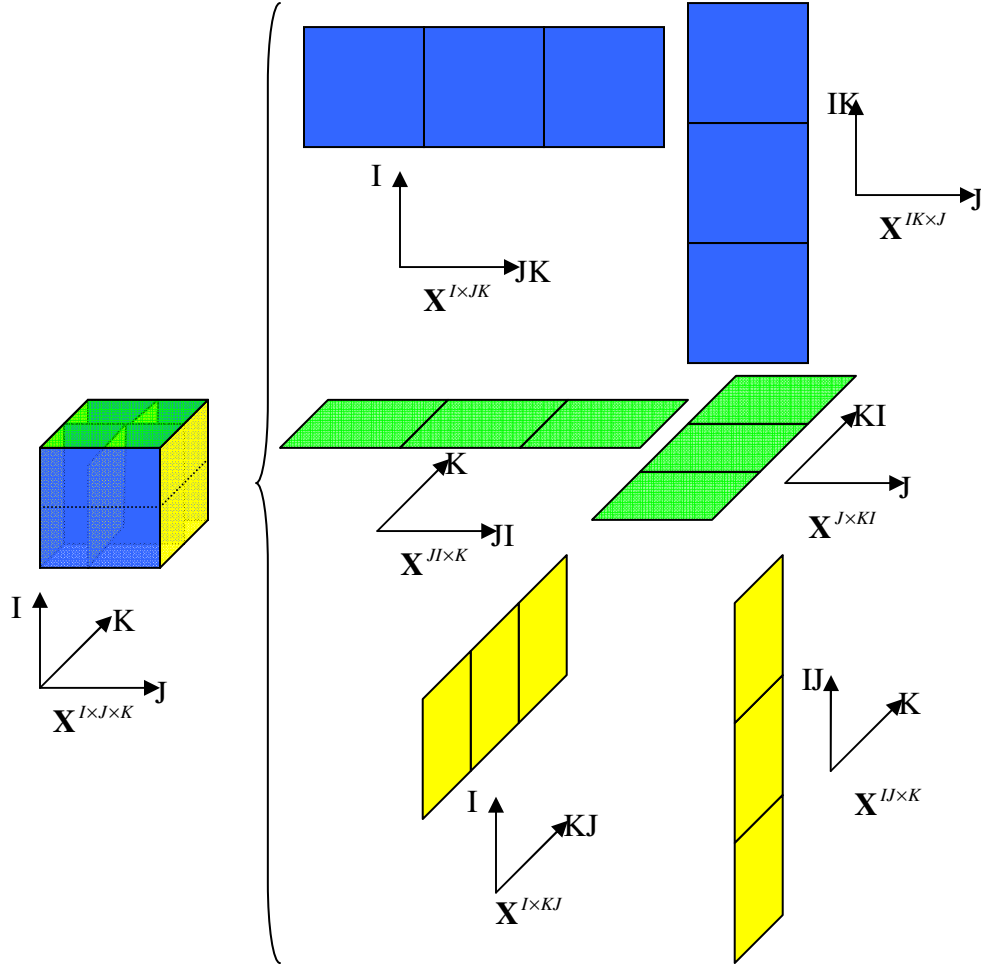


Figure 1.4: The six ways of unfolding the three-way array \mathcal{X} into a matrix.

The n -mode-matricizing $\mathbf{X}_{(n)}$ of the tensor $\mathcal{X} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_N}$ is defined as the $\dim(\mathcal{X})-2$ unfolding giving [19],[23]:

$$\mathbf{X}_{(n)} \in \mathcal{R}^{I_n \times I_1 \dots I_{n-1} I_{n+1} \dots I_N}$$

The inverse operation is denoted:

$$\mathbf{X}_{(n)^{-1}} = \mathcal{X} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_N}.$$

For MATLAB implementation of the described multi-way array manipulation as well as the algebra given in Appendix B: Multi-way array algebra, confer Appendix C: MATLAB implementation of multi-way array manipulations.

1.4 Models

The two most used forms of decomposition of multi-way arrays are the PARAFAC and the TUCKER model [23]. Where the PARAFAC decomposition gives easily interpretable components, the TUCKER model is a convincing multilinear generalization of the SVD concept to higher order [19]. Furthermore, the TUCKER model enables evaluation of the PARAFAC model using the so-called Core Consistency Diagnostic, see also paragraph 1.5.1.

1.4.1 Parallel Factor Analysis (PARAFAC)

The PARAFAC model is intrinsically related to the principle of parallel proportional profiles [5]. Suppose that the matrix $\mathbf{X}^{(1)}$ can be adequately modeled as $\mathbf{A}\mathbf{S}^T$ where the number of columns of \mathbf{A} and \mathbf{S} is the same.

$$\mathbf{X}^{(1)} = \mathbf{A}\mathbf{S}^T = \mathbf{a}_1\mathbf{s}_1^T + \mathbf{a}_2\mathbf{s}_2^T + \dots + \mathbf{a}_F\mathbf{s}_F^T = \mathbf{a}_1\mathbf{s}_1^T d_{11}^{(1)} + \mathbf{a}_2\mathbf{s}_2^T d_{22}^{(1)} + \dots + \mathbf{a}_F\mathbf{s}_F^T d_{FF}^{(1)} = \mathbf{A}\mathbf{D}^{(1)}\mathbf{S}^T \quad \text{eq. 1.27}$$

where $\mathbf{D}^{(1)} = \mathbf{I}$

Suppose another matrix $\mathbf{X}^{(2)}$ can be described by the same matrices \mathbf{A} and \mathbf{S} only in different proportions:

$$\mathbf{X}^{(2)} = \mathbf{a}_1\mathbf{s}_1^T d_{11}^{(2)} + \mathbf{a}_2\mathbf{s}_2^T d_{22}^{(2)} + \dots + \mathbf{a}_F\mathbf{s}_F^T d_{FF}^{(2)} = \mathbf{A}\mathbf{D}^{(2)}\mathbf{S}^T \quad \text{eq. 1.28}$$

where $\mathbf{D}^{(2)}$ is a diagonal matrix

The two models consist of the same (parallel) profiles only in different proportions. Cattell was the first to prove that the presence of parallel proportional profiles would lead to an unambiguous decomposition [5].

The Parallel Factor, PARAFAC, model was independently proposed by Harshman [13] and by Carrol and Chang [5] in 1970. The latter naming it Canonical Decomposition, CANDECOMP. The model can be expressed in several ways:

$$x_{ijk} = \sum_{\lambda=1}^F a_{i\lambda} b_{j\lambda} s_{k\lambda} + e_{ijk} \quad , \text{ where } F \text{ is the number of factors.} \quad \text{eq. 1.29}$$

Due to the symmetry of the components in eq. 1.29 the index order of the components doesn't matter. Another formulation of the model is given by:

$$\mathbf{X}^{(i)} = \mathbf{A}\mathbf{D}^{(i)}\mathbf{S} + \mathbf{E}^{(i)} \quad , \text{ where } \mathbf{X} = \begin{bmatrix} \mathbf{X}^{(1)} \\ \vdots \\ \mathbf{X}^{(M)} \end{bmatrix} , \text{ and } \mathbf{D}^{(i)} \text{ is a diagonal matrix} \quad \text{eq. 1.30}$$

From the formulation of the model in eq. 1.30 the relation of PARAFAC to parallel proportional profiles is evident. Finally, the model can be expressed more compact by the Khatri-Rao product.

$$\mathbf{X}^{I \times JK} = \mathbf{A}(\mathbf{S}|\otimes|\mathbf{B})^T \quad \text{eq. 1.31}$$

Where the i^{th} row of \mathbf{B} corresponds to the diagonal of $\mathbf{D}^{(i)}$. The Khatri-Rao product is given by [5]:

$$\mathbf{A}|\otimes|\mathbf{B} = [\mathbf{a}_1 \otimes \mathbf{b}_1 \quad \dots \quad \mathbf{a}_f \otimes \mathbf{b}_f], \text{ where } \begin{matrix} \mathbf{A} = [\mathbf{a}_1 & \dots & \mathbf{a}_f] \\ \mathbf{B} = [\mathbf{b}_1 & \dots & \mathbf{b}_f] \end{matrix}$$

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \dots & a_{1m}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{n1}\mathbf{B} & \dots & a_{nm}\mathbf{B} \end{bmatrix} \quad \text{eq. 1.32}$$

and the Kronecker product

The PARAFAC model is easily generalized to higher orders. The higher order equivalents are given by:

$$x_{i_1 i_2 \dots i_N} = \sum_{\lambda=1}^F a_{i_1 \lambda}^{(1)} a_{i_2 \lambda}^{(2)} \dots a_{i_N \lambda}^{(N)} + e_{i_1 i_2 \dots i_N}, \text{ where } F \text{ is the number of factors.} \quad \text{eq. 1.33}$$

Expressing the PARAFAC for arrays of more than three dimensions in terms of eq. 1.29 and eq. 1.30 yields:

$$\mathbf{X}^{(I_1 \times I_2 \times i_3 \dots i_N)} = \mathbf{A} \mathbf{D}^{(i_3)} \mathbf{D}^{(i_4)} \dots \mathbf{D}^{(i_N)} \mathbf{S} + \mathbf{E}^{(I_1 \times I_2 \times i_3 \dots i_N)} \quad \text{eq. 1.34}$$

where $\mathbf{D}^{(j)}$ are diagonal matrices

$$\mathbf{X} = \mathbf{A}(\mathbf{B}^{(N-1)}|\otimes|\mathbf{B}^{(N-2)}|\otimes|\dots|\otimes|\mathbf{B}^{(1)})^T + \mathbf{E} \quad \text{eq. 1.35}$$

From eq. 1.29 and eq. 1.33 it is seen that PARAFAC decomposes the multi-way array into a sum of effects pertaining to each dimension. Each factor consists of one vector from each dimension. Consequently, each factor's relation to each dimension can easily be read from the factor vector corresponding to the dimension, see also Figure 1.5.

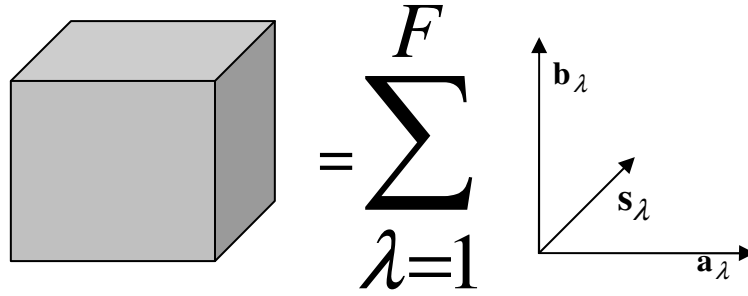


Figure 1.5: Graphical representation of the PARAFAC model as formulated in eq. 1.29. The model decomposes the multi-way array into a sum over factor effects pertaining to each dimension.

The PARAFAC model is however very restricted as the number of free parameters, D , is given by:

$$D = F \sum_{j=1}^N I_j \ll \prod_{j=1}^N I_j \text{ as } F \text{ in general is less than } \max(I_i \forall i) \quad \text{eq. 1.36}$$

Uniqueness

From the formulation of the PARAFAC model given in eq. 1.30 PARAFAC doesn't hold the rotational freedom other factor models such as independent component analysis, ICA and principal component analysis, PCA have.

$$\begin{aligned} \mathbf{X}^{(i)} = \mathbf{A}\mathbf{D}^{(i)}\mathbf{S} = \mathbf{A}\mathbf{P}\mathbf{P}^{-1}\mathbf{D}^{(i)}\mathbf{Q}\mathbf{Q}^{-1}\mathbf{S} &= (\mathbf{A}\mathbf{P})(\mathbf{P}^{-1}\mathbf{D}^{(i)}\mathbf{Q})(\mathbf{Q}^{-1}\mathbf{S}) \\ &\Downarrow \\ (\mathbf{P}^{-1}\mathbf{D}^{(i)}\mathbf{Q}) &\text{ must be a diagonal matrix} \end{aligned} \quad \text{eq. 1.37}$$

According to eq. 1.37 the rotational freedom of PARAFAC requires that the product $(\mathbf{P}^{-1}\mathbf{D}^{(i)}\mathbf{Q})$ must be a diagonal matrix. In practice, this means that \mathbf{P} and \mathbf{Q} can only be scaling and permutation matrices. Consequently, the only indeterminacies are the order of the components and the magnitudes of the loading vectors [5].

The PARAFAC model seems a logic extension of the factor analysis as the generalization to any dimension given in eq. 1.33 yields the well known factor analysis model in the 2-

way array case: $x_{ij} = \sum_{\lambda=1}^F a_{i\lambda} b_{j\lambda}$. It is, however, much more restricted than the normal

factor analysis, as a matrix of $N_1 \times N_2 \cdot N_3$ in a factor analysis with no restrictions would give $F \cdot (N_1 + N_2 \cdot N_3)$ free parameters while the PARAFAC model of the corresponding $N_1 \times N_2 \times N_3$ multi-way array only yields $F(N_1 + N_2 + N_3)$ free parameters.

Sidiropoulos and Bro have extended J. B. Kruskal's result of uniqueness from 1977 to higher orders, for the proof see [31]. The result makes use of the k-rank which is given by the least amount of columns of a matrix that are linearly independent, see also Definition 6 page 112. Let the PARAFAC model be defined as in eq. 1.33. The model is insured to be unique apart from permutations and scaling if:

$$\sum_{i=1}^N k_{\mathbf{A}^{(i)}} \geq 2F + (N - 1) \quad \text{Eq. 1.38}$$

1.4.2 TUCKER and Higher Order Singular Value Decomposition

The generalization of singular value decomposition to multidimensional data has not yet come to one "ideal" form. For a discussion on what might define the "ideal" HOSVD see [23]. However, higher order singular value decomposition, HOSVD, of the multi-way array \mathbf{X} will follow the definition of Lathauwer, Moor and Vandewalle [19].

The TUCKER model is defined by:

$$x_{i_1 i_2 \dots i_N} = \sum_{j_1}^{J_1} \sum_{j_2}^{J_2} \dots \sum_{j_N}^{J_N} s_{j_1 j_2 \dots j_N} u_{i_1 j_1}^{(1)} u_{i_2 j_2}^{(2)} \dots u_{i_N j_N}^{(N)} \quad \text{eq. 1.39}$$

In the 3-Way case the Tucker model can be formulated the following way using the Kronecker product:

$$\mathbf{X}^{I \times J \times K} = \mathbf{A} \mathbf{G} (\mathbf{B} \otimes \mathbf{S})^T \quad \text{eq. 1.40}$$

eq. 1.39 can equivalently be expressed by

$$\mathbf{X} = \mathbf{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \dots \times_N \mathbf{U}^{(N)} \quad \text{eq. 1.41}$$

Where \mathbf{S} is denoted the core multi-way array, and \times_n is the n-mode multiplication see also Definition 1 page 111.

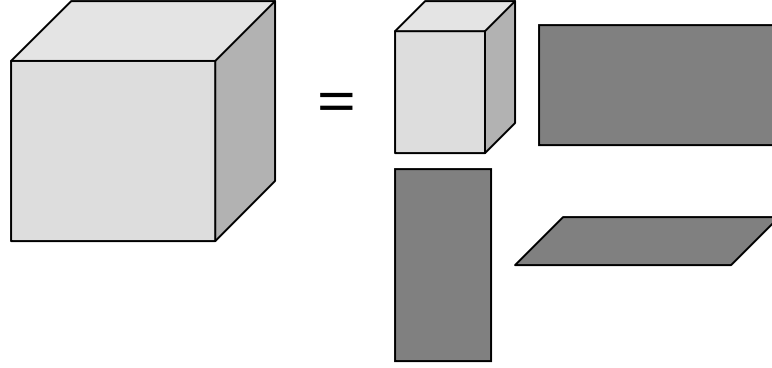


Figure 1.6: Graphical representation of the TUCKER model of a 3-way array. The model decomposes the multi-way array into matrices (dark grey) pertaining to each modality, while the core array relates each modality.

Although the TUCKER model doesn't impose any constraints, to obtain the HOSVD $\mathbf{U}^{(i)}$ has to be an orthonormal ($I_i \times I_i$) matrix, and \mathbf{S} a multi-way array of same size as \mathbf{X} subject to:

$$\begin{aligned} \langle \mathbf{S}_{I_n=\alpha}, \mathbf{S}_{I_n=\beta} \rangle &= 0 \text{ if } \alpha \neq \beta && \text{(all-orthogonality)} \\ \|\mathbf{S}_{I_n=1}\| \geq \|\mathbf{S}_{I_n=2}\| \geq \dots \geq \|\mathbf{S}_{I_n=I_n}\| &\geq 0 && \text{(ordering)} \end{aligned} \quad \text{eq. 1.42}$$

eq. 1.39 and eq. 1.41 can be equivalently expressed in matrix notation as:

$$\mathbf{X}_{(n)} = \mathbf{U}^{(n)} \cdot \mathbf{S}_{(n)} \cdot (\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(n+2)} \otimes \dots \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \mathbf{U}^{(2)} \otimes \dots \otimes \mathbf{U}^{(n-1)}) \quad \text{eq. 1.43}$$

Where $\mathbf{X}_{(n)}$ is the n-mode-matricizing of \mathbf{X} and $\mathbf{S}_{(n)}$ the n-mode-matricizing of \mathbf{S} . The resemblance of HOSVD to SVD becomes evident as the singular value decomposition of a matrix \mathbf{F} can be expressed by [19]:

$$\begin{array}{c} I_2 \\ \boxed{\mathbf{F}} \\ I_1 \end{array} = \begin{array}{c} I_1 \\ \boxed{\mathbf{U}} \\ I_1 \end{array} \begin{array}{c} I_2 \\ \boxed{\mathbf{S}} \\ I_1 \end{array} \begin{array}{c} I_2 \\ \boxed{\mathbf{V}^T} \\ I_2 \end{array}$$

$$\mathbf{F} = \underbrace{\mathbf{U}^{(1)} \mathbf{S} \mathbf{U}^{(2)}}_{\text{eq. 1.43}} = \underbrace{\mathbf{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)}}_{\text{eq. 1.41}} \quad \text{eq. 1.44}$$

Here $\mathbf{S} \in \mathcal{R}^{I_1 \times I_2}$ is an ordered pseudo-diagonal matrix, i.e.

$$\begin{aligned} \mathbf{S} &= \text{diag}(\sigma_1, \sigma_1 \dots \sigma_{\min(I_1, I_2)}) \\ \sigma_1 &\geq \sigma \geq \dots \geq \sigma_{\min(I_1, I_2)} \geq 0 \end{aligned}$$

Evidently, HOSVD becomes SVD for the 2-way array.

The number of free parameters for the HOSVD is given by:

$$D = \underbrace{\prod_{i=1}^N I_i - \sum_{i=1}^N \sum_{j=1}^{I_i-1} j}_{\mathbf{S} \text{ All orthogonal}} + \sum_{i=1}^N \underbrace{\left(I_i^2 - \left(\sum_{j=1}^{I_i-1} j \right) - \underbrace{I_i}_{\text{norm}=1} \right)}_{\mathbf{U}^{(i)} \text{ Orthogonality}} = \text{eq. 1.45}$$

$$\underbrace{\prod_{i=1}^N I_i - \sum_{i=1}^N \frac{I_i^2 - I_i}{2}}_{\mathbf{S}} + \sum_{i=1}^N \underbrace{\left(I_i^2 - \frac{I_i^2 - I_i}{2} - I_i \right)}_{\mathbf{U}^{(i)}} = \prod_{i=1}^N I_i$$

1. Follows as $\sum_{i=1}^n \frac{n(n+1)}{2}$

As can be seen from eq. 1.45 the free parameters of HOSVD exactly match the number of parameters of the multi-way array. Consequently, the HOSVD as a decomposition will reconstruct the multi-way array exact.

Calculating the HOSVD

According to eq. 1.43:

$$\mathbf{X}_{(n)} = \mathbf{U}^{(n)} \cdot \mathbf{S}_{(n)} \cdot \left(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(n+2)} \otimes \dots \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \mathbf{U}^{(2)} \otimes \dots \otimes \mathbf{U}^{(n-1)} \right)^T$$

Defining the following set of matrices [19]:

$$\mathbf{V}^{(n)} = \tilde{\mathbf{S}}_{(n)} \left(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(n+2)} \otimes \dots \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \mathbf{U}^{(2)} \otimes \dots \otimes \mathbf{U}^{(n-1)} \right)^T$$

$$\mathbf{S}_{(n)} = \mathbf{\Sigma}^{(n)} \tilde{\mathbf{S}}_{(n)}$$

$$\mathbf{\Sigma}^{(n)} = \text{diag}([\sigma_1, \sigma_1 \dots \sigma_N])$$

Where $\mathbf{\Sigma}^{(n)}$ is selected so that $\tilde{\mathbf{S}}_{(n)}$ is a normalized version of $\mathbf{S}_{(n)}$ with the rows scaled to unit length. This gives:

$$\mathbf{X}_{(n)} = \mathbf{U}^{(n)} \cdot \mathbf{S}_{(n)} \cdot \left(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(n+2)} \otimes \dots \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \mathbf{U}^{(2)} \otimes \dots \otimes \mathbf{U}^{(n-1)} \right)^T = \mathbf{U}^{(n)} \mathbf{\Sigma}^{(n)} \mathbf{V}^{(n)} \quad \text{eq. 1.46}$$

From eq. 1.46 it is seen that $\mathbf{U}^{(n)}$ can be calculated by the normal singular value decomposition (SVD) of $\mathbf{X}_{(n)}$.

From eq. 1.41 we had:

$$\mathbf{X} = \mathbf{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \dots \times_N \mathbf{U}^{(N)}$$

This gives:

$$\begin{aligned} \mathbf{X} &= \mathbf{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \dots \times_N \mathbf{U}^{(N)} \Leftrightarrow \\ \mathbf{X} \times_N \mathbf{U}^{(N)T} &= \mathbf{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \dots \times_{N-1} \mathbf{U}^{(N-1)} \Leftrightarrow \\ \mathbf{X} \times_N \mathbf{U}^{(N)T} \times_{N-1} \mathbf{U}^{(N-1)T} \dots \times_1 \mathbf{U}^{(1)T} &= \mathbf{S} \Leftrightarrow \\ \mathbf{S} &= \mathbf{X} \times_1 \mathbf{U}^{(1)T} \times_2 \mathbf{U}^{(2)T} \dots \times_N \mathbf{U}^{(N)T} \end{aligned} \quad \text{eq. 1.47}$$

1. Follows as $\mathbf{U}^{(i)}$ is an orthonormal ($I_i \times I_i$) matrix.
2. Result of the end of Definition 1 page 111.

\mathbf{S} can equivalently be calculated by the Kronecker product:

$$\mathbf{S}_{(n)} = \mathbf{U}^{(n)T} \mathbf{X}_{(n)} \left(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(n+2)} \otimes \dots \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \mathbf{U}^{(2)} \otimes \dots \otimes \mathbf{U}^{(n-1)} \right) \quad \text{eq. 1.48}$$

From eq. 1.42 \mathbf{S} had to fulfill $\|\mathbf{s}_{I_n=1}\| \geq \|\mathbf{s}_{I_n=2}\| \geq \dots \geq \|\mathbf{s}_{I_n=I_n}\| \geq 0$ this ordering is ensured by the fact that $\mathbf{U}^{(n)} \forall n$ is ordered. $\langle \mathbf{s}_{I_n=\alpha}, \mathbf{s}_{I_n=\beta} \rangle = 0$ if $\alpha \neq \beta$ can be proven by:

$$\begin{aligned} \mathbf{S}_{(n)} &= \mathbf{U}^{(n)T} \mathbf{X}_{(n)} \left(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \dots \otimes \mathbf{U}^{(n-1)} \right) \\ \mathbf{S}_{(n)} \mathbf{S}_{(n)}^T &= \mathbf{U}^{(n)T} \mathbf{X}_{(n)} \left(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \dots \otimes \mathbf{U}^{(n-1)} \right) \left(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \dots \otimes \mathbf{U}^{(n-1)} \right)^T \mathbf{X}_{(n)}^T \mathbf{U}^{(n)} = \\ &= \mathbf{U}^{(n)T} \mathbf{U}^{(n)} \boldsymbol{\Sigma}^{(n)} \mathbf{V}^{(n)} \mathbf{V}^{(n)T} \boldsymbol{\Sigma}^{(n)T} \mathbf{U}^{(n)T} \mathbf{U}^{(n)} = \boldsymbol{\Sigma}^{(n)} \mathbf{V}^{(n)} \mathbf{V}^{(n)T} \boldsymbol{\Sigma}^{(n)T} = \boldsymbol{\Sigma}^{(n)} \boldsymbol{\Sigma}^{(n)T} \end{aligned}$$

1. Follows as

$$\begin{aligned} &\left(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \dots \otimes \mathbf{U}^{(n-1)} \right) \left(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(N)} \otimes \mathbf{U}^{(1)} \otimes \dots \otimes \mathbf{U}^{(n-1)} \right)^T = \mathbf{I} \text{ and} \\ &\mathbf{X}_{(n)} = \mathbf{U}^{(n)} \boldsymbol{\Sigma}^{(n)} \mathbf{V}^{(n)} \end{aligned}$$

As $\boldsymbol{\Sigma}^{(n)} \boldsymbol{\Sigma}^{(n)T}$ is a diagonal matrix the assumption of $\langle \mathbf{s}_{I_n=\alpha}, \mathbf{s}_{I_n=\beta} \rangle = 0$ if $\alpha \neq \beta$ holds.

$$\begin{aligned}
N &= \dim(\mathbf{X}) \\
\mathbf{U}^{(1)} &\leftarrow \text{SVD}(\mathbf{X}_{(1)}) \\
\mathbf{U}^{(2)} &\leftarrow \text{SVD}(\mathbf{X}_{(2)}) \\
&\vdots \quad \quad \quad \vdots \\
\mathbf{U}^{(N)} &\leftarrow \text{SVD}(\mathbf{X}_{(N)}) \\
\mathbf{S} &= \mathbf{X} \times_1 \mathbf{U}^{(1)T} \times_2 \mathbf{U}^{(2)T} \dots \times_N \mathbf{U}^{(N)T}
\end{aligned}$$

Figure 1.7: The calculation of HOSVD

As shown in Figure 1.7, HOSVD as for SVD doesn't require any iterative operations to be estimated. Furthermore, for SVD $\mathbf{U}^{(1)}$ is the eigenvectors of $\mathbf{X}\mathbf{X}^T = \mathbf{X}_{(1)}\mathbf{X}_{(1)}^T$ and $\mathbf{U}^{(2)}$ is the eigenvectors of $\mathbf{X}^T\mathbf{X} = \mathbf{X}_{(2)}\mathbf{X}_{(2)}^T$. Clearly, HOSVD keeps this characteristic, as $\mathbf{U}^{(n)}$ is the eigenvectors of $\mathbf{X}_{(n)}\mathbf{X}_{(n)}^T$. Consequently, HOSVD seems like a parsimonious multi-dimensional generalization of SVD.

The uniqueness of HOSVD

As $\mathbf{U}^{(n)}$ can be calculated by the normal singular value decomposition (SVD) of $\mathbf{X}_{(n)}$ it follows that $\mathbf{U}^{(n)}$ share the uniqueness properties of SVD:

If the singular values found in eq. 1.46 all are different then $\mathbf{u}_i^{(n)}$ is unique up to sign. Furthermore, the vectors corresponding to the same n-mode singular value can be replaced by multiplication with an orthogonal matrix.

If $\mathbf{U}^{(n)}$ isn't unique there exist an orthogonal matrix \mathbf{Q} so: $\mathbf{V}^{(n)} = \mathbf{U}^{(n)}\mathbf{Q}$ is also a solution of eq. 1.46. This gives a new core given by $\hat{\mathbf{S}} = \mathbf{S} \times_{(n)} \mathbf{Q}^{-1}$. As \mathbf{S} in eq. 1.48 is found by the knowledge of $\mathbf{U}^{(n)} \forall n$ and $\mathbf{X}_{(n)}$ it follows that \mathbf{S} due to the fact that $\mathbf{u}_i^{(n)}$ only maximally is unique up to sign can't be unique. As revealed in eq. 1.48 the change of sign gives a new value of \mathbf{S} . This gives as many different values as there are sign-combinations of the vectors of $\mathbf{U}^{(n)} \forall n$.

1.4.3 Model Relations

According to Definition 7 page 112 a multi-way array is called diagonalizable if the core multi-way array \mathbf{S} of the HOSVD fulfills $\mathbf{S}_{i_1 i_2 \dots i_N} = 0$ unless $i_1 = i_2 = \dots = i_N$.

With this definition the TUCKER, HOSVD and PARAFAC models can be related as shown in Figure 1.8.

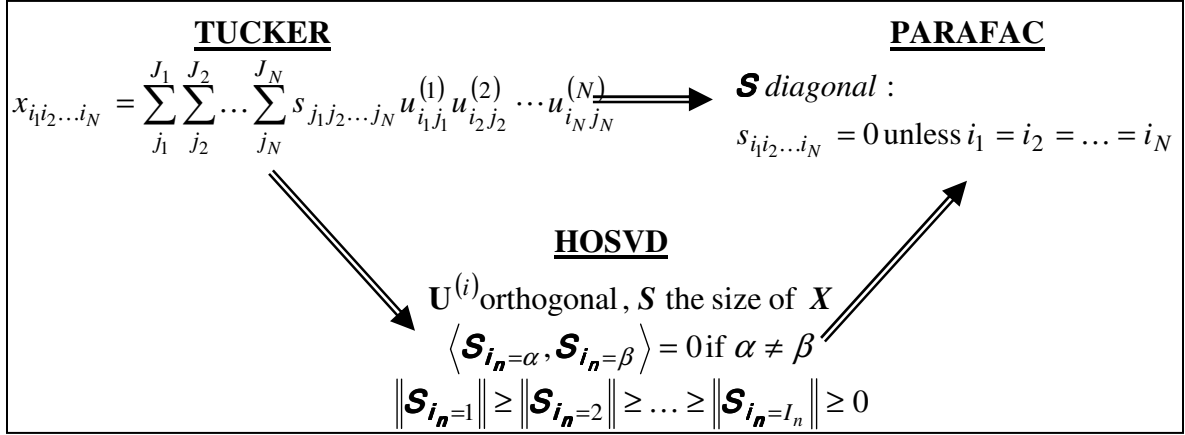


Figure 1.8: The relation between the TUCKER, HOSVD and PARAFAC model.

From Figure 1.8 it follows that if a multi-way array is diagonalizable it can be perfectly represented by the PARAFAC model. Furthermore, if the HOSVD was based on a PARAFAC model this would require the core \mathbf{S} to be diagonal. However, this is too strong a condition greatly reducing the number of free parameters (compare eq. 1.36 to eq. 1.45) disabling the HOSVD to perfectly model the data. As the core \mathbf{S} has same size as \mathbf{X} the results of HOSVD contrary to the PARAFAC model is very hard to interpret.

The PARAFAC model is also related to the rank of a multi-way array. An N-way array \mathbf{A} has rank-1 when it equals the outer product of N vectors, i.e. $\mathbf{A} = \mathbf{u}^{(1)} \circ \mathbf{u}^{(2)} \circ \dots \circ \mathbf{u}^{(N)}$. PARAFAC of a multi-way array is the decomposition of the multi-way array into a minimal sum of rank-1 components [21]. Furthermore, by the definition of the rank of multi-way arrays, see the notice in Definition 4 page 111, the PARAFAC decomposition approximately describes the rank of the multi-way array. Consequently, a PARAFAC model based on finding the best sum of rank-1 components as described in [20],[33] will also be implemented in the following section.

1.5 PARAFAC Algorithms

Algorithms are based on an initialization and an iterative optimization. Before addressing the three iterative optimization approaches; Alternating Least Squares, EM-algorithm and VBEM-algorithm, the problem of initialization will first be addressed. This will be followed by a brief description of the handling of non-negativity and a way of evaluating the PARAFAC model by the so called Core Consistency Diagnostic.

Initialization

The choice of initialization can have significant impact on the time it takes an algorithm to converge, but more serious problems arises in situations where the function to optimize has local extremes. In this situation, the algorithm might also converge to different values depending on where it is initialized, see Figure 1.9.

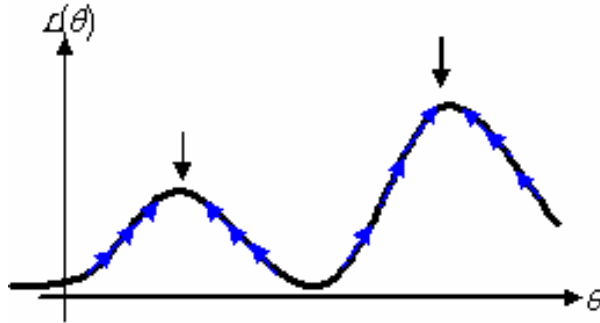


Figure 1.9: A one dimensional optimization situation where convergence depends on initialization. Here, initializing with a low value of θ makes the algorithm converge to a local extreme.

From Figure 1.8 it was seen that if the HOSVD had a diagonal core it could be considered a PARAFAC model. Furthermore, HOSVD has the advantage that $\mathbf{U}^{(i)}$ is ordered making the first F eigenvectors of $\mathbf{U}^{(i)}$ describe the most of the variation of the data. Let \mathbf{k} be the vector containing the diagonal of \mathbf{S} . From the HOSVD an initial guess of the PARAFAC model parameters can be found by taking the F first eigenvectors of $\mathbf{U}^{(i)} \forall i$:

$$x_{i_1 i_2 \dots i_N} = \sum_{\lambda=1}^F k_{\lambda} u_{i_1 \lambda}^{(1)} u_{i_2 \lambda}^{(2)} \dots u_{i_N \lambda}^{(N)} = \sum_{\lambda=1}^F v_{i_1 \lambda} u_{i_2 \lambda}^{(2)} \dots u_{i_N \lambda}^{(N)} \text{ where } v_{i_1 \lambda} = k_{\lambda} u_{i_1 \lambda}^{(1)} \quad \text{eq. 1.49}$$

However, to insure no local optimum is found from the algorithms using different initialization points is recommended.

Non negativity

In the problems at hand, the PARAFAC models are required to yield non-negative results. In the case of the VBEM-algorithm, insuring non-negativity could be achieved by using non-negative distributions such as gamma distributions as priors of the factors,

however only the hidden variables in the EM-algorithm can in this way be assured non negativity. Bro derives in his thesis two simple methods of assuring non-negativity.

One way of assuring non-negativity is simply to find the optimal unconstrained solution and to set all negative elements in every column of the factors to zero, see also Theorem 10 page 109. This method is denoted “Column-wise Non-Negativity” [5].

A more flexible method is to insure the non-negativity constraint using the following algorithm which manipulates iteratively each row of the constrained factors. The algorithm uses the feature that each element in the row of a factor is only affected by the other elements of the same row, consequently the name “Row-wise Non-Negativity” [5].

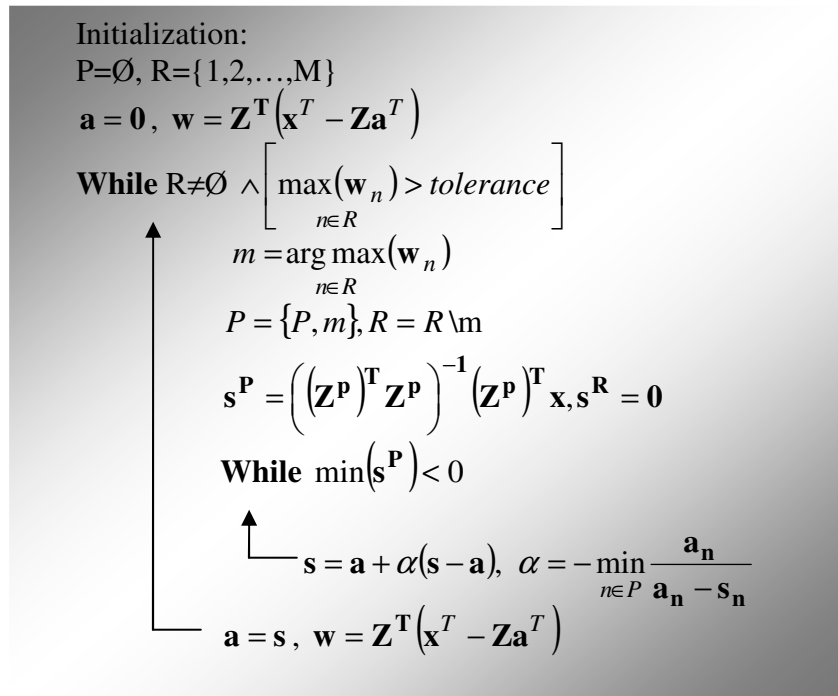


Figure 1.10: Algorithm for “Row-wise Non-Negativity”. x corresponds to a row of X , where X and Z is defined as in Figure 1.11.

The algorithm is optimal in a least square sense as it optimizes the parameters according to which variable contributes the most to the squared error, i.e. $\|x^T - Za^T\|^2$. The “Row-wise” and “Column-wise” implementation of non-negativity is only valid for algorithms based on Least Square optimization, i.e. minimizing the sum of square error.

Even though the data to model aren’t non-negative it can be an advantage to impose non-negativity by adding a constant and fitting the model under the non-negativity constraint. This insures no factor can counteract the effect of any other factor by having opposite sign eliminating the risk of degeneracy in the factors.

1.5.1 Core Consistency Diagnostic

The Core Consistency Diagnostic, CCD, can be applied to any model that can be considered a restricted 3-way TUCKER model [5]. Consider the PARAFAC model given in eq. 1.31. i.e. $\mathbf{X}^{I \times JK} = \mathbf{A}(\mathbf{S}|\otimes|\mathbf{B})^T$. According to Figure 1.8 this can be considered a restricted 3-way TUCKER model as given in eq. 1.40, i.e. $\mathbf{X}^{I \times JK} = \mathbf{A}\mathbf{G}(\mathbf{S}|\otimes|\mathbf{B})^T$ where the core \mathbf{G} is zero apart from along the superdiagonal which has ones. This core is denoted \mathbf{T} . \mathbf{G} of the TUCKER model is now found by inserting \mathbf{A} , \mathbf{B} and \mathbf{S} obtained from the PARAFAC model, i.e [5]:

$$\min \left\| \mathbf{X}^{I \times JK} - \mathbf{A}\mathbf{G}(\mathbf{S}|\otimes|\mathbf{B})^T \right\|^2 = \min \left\| \text{vec}\mathbf{X} - (\mathbf{B} \otimes \mathbf{S}|\otimes|\mathbf{A})\text{vec}\mathbf{G} \right\|^2 \quad \text{eq. 1.50}$$

If the PARAFAC model is valid \mathbf{G} should resemble \mathbf{T} . A measure of resemblance is the core consistency¹:

$$\text{Core Consistency} = 100 \cdot \left(1 - \frac{\sum_{d=1}^F \sum_{e=1}^F \sum_{f=1}^F (t_{def} - g_{def})^2}{\sum_{d=1}^F \sum_{e=1}^F \sum_{f=1}^F g_{def}^2} \right) \quad \text{eq. 1.51}$$

From eq. 1.51 it's seen that if the PARAFAC model is perfect, the nominator becomes zero giving a 100% consistency. If the PARAFAC model isn't correct the percentage of \mathbf{G} not consistent with \mathbf{T} reduces the Core Consistency. A core consistency well below 70-90% indicates that either too many components are used or the model otherwise is mis-specified [1].

Although the Core Consistency is an effective measure of how many factors to include, Bro emphasizes that other measures such as sum of squared residuals versus number of factors, inspection of the parameters and cross validation also should be taken into consideration.

¹ Bro has unclearly defined the Core Consistency in his thesis [5] and several other papers, he defines:

$$\text{Core Consistency} = 100 \cdot \left(\frac{1 - \sum_{d=1}^F \sum_{e=1}^F \sum_{f=1}^F (t_{def} - g_{def})^2}{\sum_{d=1}^F \sum_{e=1}^F \sum_{f=1}^F t_{def}^2} \right)$$

However, in his implementation he uses the definition given in eq. 1.51.

The core consistency generalized to higher orders yields:

$$\text{Core Consistency} = 100 \cdot \left(1 - \frac{\sum_{i_1=1}^F \sum_{i_2=1}^F \cdots \sum_{i_N=1}^F (t_{i_1 i_2 \dots i_N} - g_{i_1 i_2 \dots i_N})^2}{\sum_{i_1=1}^F \sum_{i_2=1}^F \cdots \sum_{i_N=1}^F g_{i_1 i_2 \dots i_N}^2} \right) \quad \text{eq. 1.52}$$

1.5.2 PARAFAC by Alternating Least Squares

Rasmus Bro and Claus Andersson have created a multi-way toolbox for Matlab [1]. Their implementation of PARAFAC is based on the technique of alternating least squares (ALS).

The principle of ALS is quite simple; initialize all model parameters for example randomly. Update each parameter by minimizing a cost-function with respect to the parameter while holding all other parameters fixed.

Consider the PARAFAC model as defined in eq. 1.31. Giving the cost function $\min \|\mathbf{X}^{I \times JK} - \mathbf{A}(\mathbf{S}|\otimes|\mathbf{B})^T\|^2$, the Alternating Least Squares algorithm for PARAFAC is then defined by [5]:

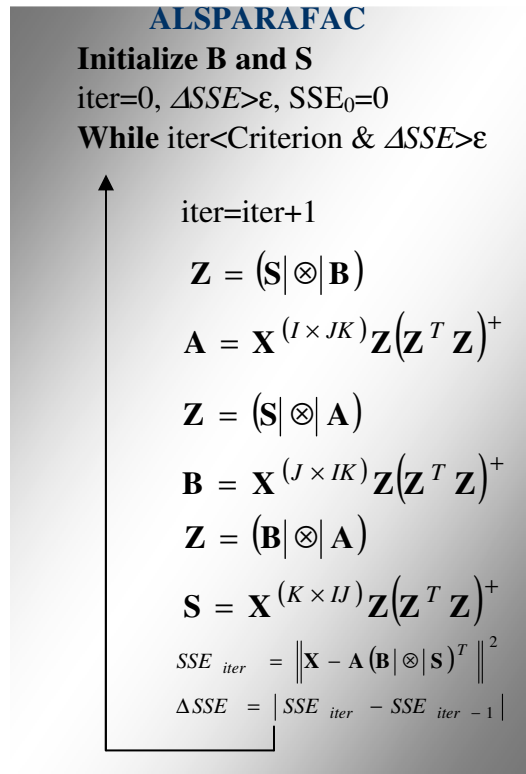


Figure 1.11: The ALSPARAFAC algorithm.

Notice how the algorithm makes use of the interchangeability of the model parameters as revealed in eq. 1.29 changing the order of the parameters by changing the unfolding of \mathbf{X} .

1.5.3 PARAFAC by multi-way rank one decomposition

As the PARAFAC model can be formulated as a sum of rank one components a PARAFAC algorithm can be defined in the framework of the HOSVD as explained in [19], [20]. The algorithm finds the best rank one decomposition by an alternating least square approach where each dimensions factor can be directly found from the n-mode multiplication. Each consecutive factor explains the most of the remaining variation in the multi-way array.

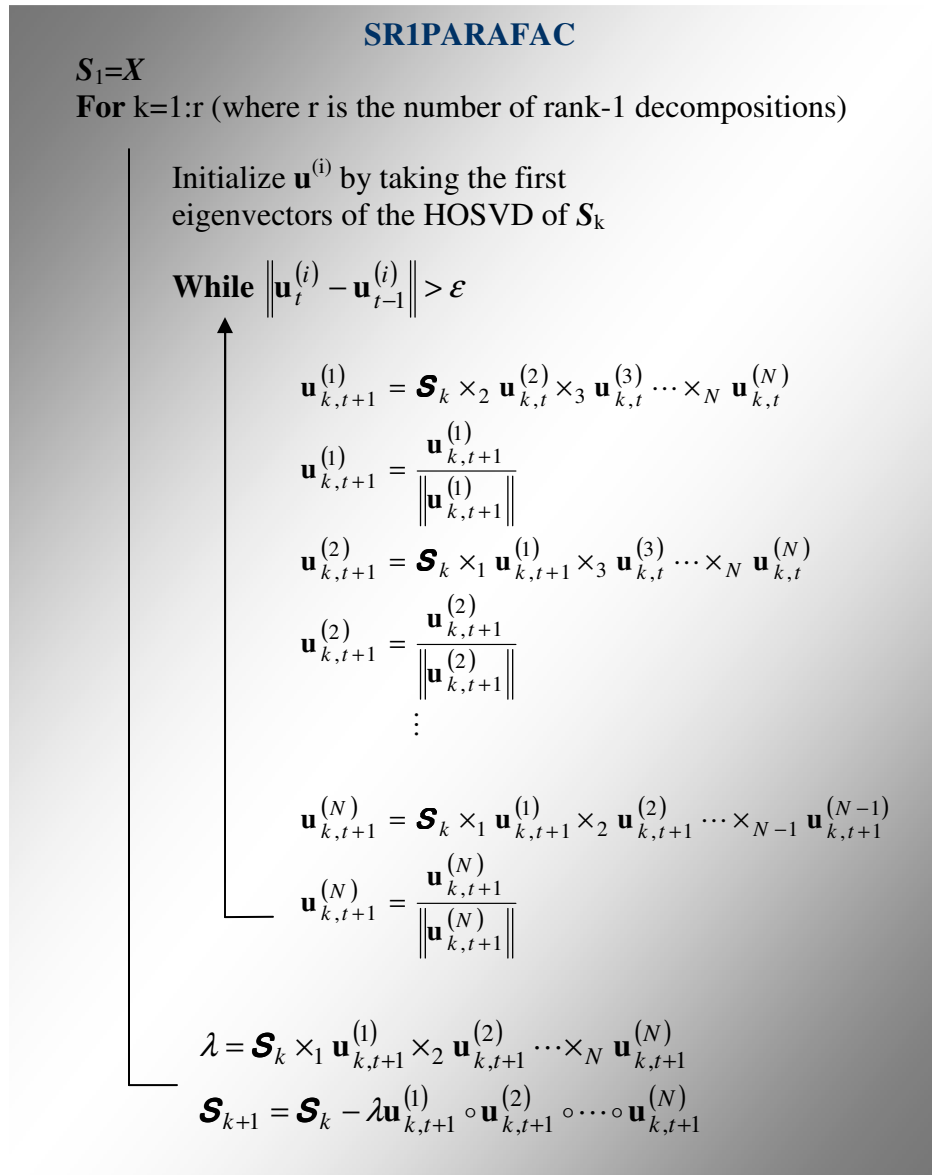


Figure 1.12: PARAFAC based on a sum of multi-way rank one decompositions, SR1PARAFAC.

1.5.4 PARAFAC by EM and VBEM

Frederik Brink Nielsen seems to be the first to derive the PARAFAC model in a statistical framework using the EM and VBEM algorithms [27]. His derivation is based on the assumption of normal distributed factors. Although non-negativity could be insured for example by insuring gamma priors it turns out to be unnecessary in the case of the EM algorithm as the row wise non negativity constraint can be implemented on all but the hidden variable \mathbf{S} . This practically always insures \mathbf{S} to be positive. Gamma priors could have been used in the VBEM algorithm to ensure the non-negativity. In order to implement the gamma priors each element of the k^{th} factor would depend on the other elements of its row. The algorithm would have to iterate over each dimensions own factors slowing down the already very slow algorithm. Consequently, the factors are assumed normal distributed as given by [27]. However, if the VBEM algorithm finds the true factors, non-negativity becomes just a matter of choosing the correct sign of each factor.

PARAFAC by EM

The expectation maximization of PARAFAC will be based on the following assumptions:

Model :	$\mathbf{X}^{(i)} = \mathbf{A}\mathbf{D}^{(i)}\mathbf{S} + \mathbf{E}^{(i)}$
Assumptions :	$\mathcal{N}(\mathbf{e}^{(i)} \mid \mathbf{0}, \mathbf{\Psi}^{(i)}) \quad (1)$
	$\mathcal{N}(\mathbf{s}_j \mid \mathbf{0}, \mathbf{I}) \quad (2)$
	$\mathbf{\Psi}^{(i)}, \mathbf{D}_i \text{ diagonal} \quad (3)$

Due to assumption 1 the following is valid:

$$\begin{aligned}
 p\left(\{\mathbf{x}_i^{(m)}\}_{m=1}^M \mid \mathbf{s}_i, \mathbf{A}, \{\mathbf{D}^{(m)}, \mathbf{\Psi}_j^{(m)}\}_{m=1}^M\right) &= \prod_{m=1}^M (2\pi)^{-N/2} \left| \mathbf{\Psi}^{(m)} \right|^{-\frac{1}{2}} \exp\left[-\frac{1}{2} \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i \right)^T \mathbf{\Psi}^{(m)-1} \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i \right)\right] \\
 &\propto \exp\left[-\frac{1}{2} \sum_{m=1}^M \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i \right)^T \mathbf{\Psi}^{(m)-1} \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i \right)\right]
 \end{aligned}
 \tag{eq. 1.53}$$

E Step

As derived in Theorem 2 page 95 the following holds:

$$p\left(\mathbf{s}_i | \mathbf{A}, \{\mathbf{x}_i^{(m)}, \mathbf{D}^{(m)}, \boldsymbol{\Psi}^{(m)}\}_{m=1}^M\right) \propto \exp\left(\mathbf{s}_i^T \mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} - \frac{1}{2} \mathbf{s}_i^T \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)}\right) \mathbf{s}_i\right) \quad \text{eq. 1.54}$$

This yields the update rule below, as derived in Theorem 3 page 96 :

$$\begin{aligned} \boldsymbol{\Sigma}_{\mathbf{S}} &= \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \right)^{-1} \\ \langle \mathbf{S} \rangle &= \boldsymbol{\Sigma}_{\mathbf{S}} \left(\sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{X}_m \right) \\ \langle \mathbf{S} \mathbf{S}^T \rangle &= N \boldsymbol{\Sigma}_{\mathbf{S}} + \langle \mathbf{S} \rangle \langle \mathbf{S} \rangle^T \end{aligned}$$

M Step

In the M Step the likelihood $\mathcal{L}(\boldsymbol{\theta})$ is maximized. From eq. 1.15 and eq. 1.53 the following is derived as revealed in Theorem 4 page 98:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= \\ &= -\frac{N}{2} \sum_{m=1}^M \ln |\boldsymbol{\Psi}^{(m)}| - \frac{1}{2} \sum_{i=1}^N \sum_{m=1}^M \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} + \text{tr} \left(\mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s} \mathbf{s}^T \rangle \right) - 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle + \text{const} \end{aligned} \quad \text{eq. 1.55}$$

This yields the update rule derived in Theorem 5 page 99:

$$\begin{aligned} \mathbf{a}_k &= \left(\sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \mathbf{D}^{(m)} \right)_k \left(\sum_{m=1}^M \left(\boldsymbol{\Psi}^{(m)-1} \right)_{kk} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle^T \mathbf{D}^{(m)} \right)^{-1} \\ \mathbf{d}_m &= \left(\langle \mathbf{S} \mathbf{S}^T \rangle \bullet \left(\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \right) \right)^{-1} \text{vec} \left[\text{diag} \left[\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \right] \right] \\ \boldsymbol{\Psi}^{(m)} &= \frac{1}{N} \text{diag} \left[\mathbf{X}^{(m)} \mathbf{X}^{(m)T} + \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle \mathbf{D}^{(m)} \mathbf{A}^T - 2 \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \rangle \mathbf{X}^{(m)T} \right] \end{aligned}$$

Consequently the EMPARAFAC algorithm can be stated as shown in Figure 1.13:

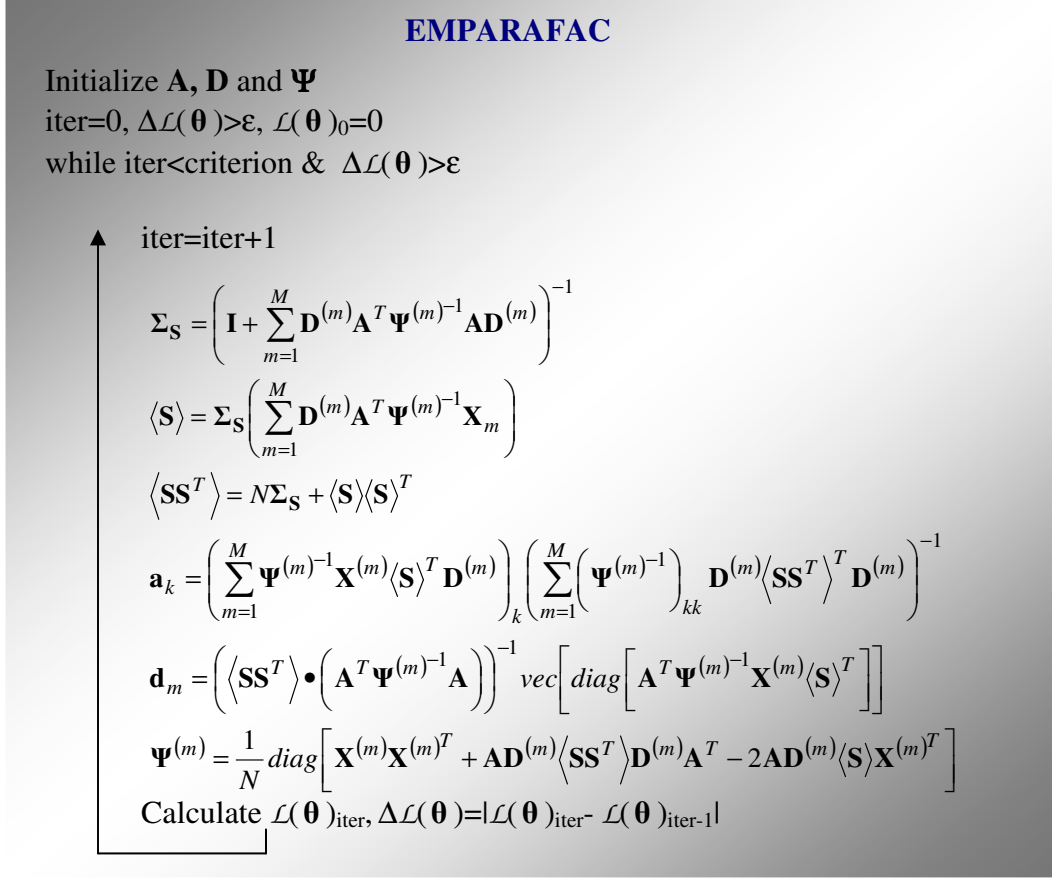


Figure 1.13: The EMPARAFAC algorithm.

PARAFAC by VBEM

For the VBEM algorithm the following assumptions are made:

$$p(\mathbf{E}^{(m)}) = \prod_{i=1}^N \mathcal{N}(\mathbf{e}_i^{(m)} | \boldsymbol{\theta}, \text{diag}(\boldsymbol{\varphi}_m))$$

$$p(\mathbf{S}) = \prod_{i=1}^N \mathcal{N}(\mathbf{s}_i | \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{A} | \boldsymbol{\alpha}) = \prod_{f=1}^F \mathcal{N}(\mathbf{a}_f | \mathbf{0}, \frac{1}{\alpha_f} \mathbf{I})$$

$$p(\boldsymbol{\alpha} | a^\alpha, b^\alpha) = \prod_{f=1}^F \mathcal{G}(\alpha_f | a_f^\alpha, b_f^\alpha)$$

$$p(\mathbf{d}_m | \boldsymbol{\mu}_m, \frac{1}{v_m}) = \mathcal{N}(\mathbf{d}_m | \boldsymbol{\mu}_m, \frac{1}{v_m} \mathbf{I})$$

$$p(\boldsymbol{\varphi}_m | \mathbf{a}_m^\varphi, \mathbf{b}_m^\varphi) = \prod_{p=1}^P \mathcal{G}(\varphi_{mp} | a_{mp}^\varphi, b_{mp}^\varphi)$$

For the derivation of the VBEM algorithm consult Theorem 9 page 104. The derivation yields the algorithm shown in Figure 1.14:

VBPARAFAC

Initialize $\mathbf{A}, \mathbf{D}, \boldsymbol{\alpha}, \boldsymbol{\varphi}, \mathbf{a}^\alpha, \mathbf{b}^\alpha, \mathbf{a}^\varphi, \mathbf{b}^\varphi$
iter=0, $\Delta \mathcal{F} > \epsilon$, $\mathcal{F}_0 = 0$
while iter < Criterion & $\Delta \mathcal{F} > \epsilon$

↑ iter=iter+1

$$\boldsymbol{\Sigma}_{\mathbf{S}} = \left(\mathbf{I} + \sum_{m=1}^M \langle \mathbf{d}_m^T \mathbf{d}_m \rangle \bullet \sum_{p=1}^P \langle \boldsymbol{\varphi}_{mp} \rangle \langle \mathbf{a}_p^T \mathbf{a}_p \rangle \right)^{-1}$$

$$\boldsymbol{\mu}_{\mathbf{s}_i} = \boldsymbol{\Sigma}_{\mathbf{S}} \sum_{m=1}^M \langle \mathbf{d}_m^T \rangle \langle \mathbf{A}^T \rangle \langle \boldsymbol{\varphi}_m \rangle \mathbf{x}_i$$

$$\boldsymbol{\Sigma}_{\mathbf{a}^p} = \left(\langle \boldsymbol{\alpha} \rangle \mathbf{I} + \sum_{m=1}^M \langle \boldsymbol{\varphi}^{(m)} \rangle \langle \mathbf{d}_m \mathbf{d}_m^T \rangle \bullet \sum_{i=1}^N \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \right)^{-1}$$

$$\boldsymbol{\mu}_{\mathbf{a}^p} = \boldsymbol{\Sigma}_{\mathbf{a}^p} \left\langle \sum_{m=1}^M \langle \boldsymbol{\varphi}_p^{(m)} \rangle \langle \mathbf{D}^{(m)} \rangle \sum_{i=1}^N \mathbf{x}_i^{(m)T} \langle \mathbf{s}_i \rangle \right\rangle$$

$$\boldsymbol{\Sigma}_{\mathbf{C}^{(m)}} = \left(v_m \mathbf{I} + \sum_{j=1}^P \langle \boldsymbol{\varphi}_{mj} \rangle \langle \mathbf{a}_{:,j}^T \mathbf{a}_{:,j} \rangle \bullet \sum_{i=1}^N \langle \mathbf{s}_i^T \mathbf{s}_i \rangle \right)^{-1}$$

$$\mathbf{m}_{\mathbf{C}^{(m)}} = \boldsymbol{\Sigma}_{\mathbf{C}^{(m)}} \left(v_m \boldsymbol{\mu}_m + \sum_{i=1}^N \sum_{j=1}^D \langle \mathbf{x}_{ij}^{(m)} \rangle \boldsymbol{\varphi}_{mj} \langle \mathbf{s}_i \rangle \bullet \langle \mathbf{a}_{:,j} \rangle \right)$$

$$\hat{a}_f^\alpha = a_f^\alpha + \frac{P}{2}$$

$$\hat{b}_f^\alpha = b_f^\alpha + \frac{\langle \mathbf{a}_f^T \mathbf{a}_f \rangle}{2}$$

$$\langle \boldsymbol{\alpha}_f \rangle = \frac{\hat{a}_f^\alpha}{\hat{b}_f^\alpha}$$

$$\hat{a}_{mp}^\varphi = a_{mp}^\varphi + \frac{N}{2}$$

$$\hat{b}_{mp}^\varphi = b_{mp}^\varphi + \frac{1}{2} \sum_{i=1}^N \left\langle \left(\mathbf{x}_{ip}^{(m)} - \mathbf{a}_{:,p} \mathbf{D}^{(m)} \mathbf{s}_i \right)^2 \right\rangle$$

$$\langle \boldsymbol{\varphi}_{mp} \rangle = \frac{\hat{a}_{mp}^\varphi}{\hat{b}_{mp}^\varphi}$$

Calculate $\mathcal{F}(\boldsymbol{\theta})_{\text{iter}}, \Delta \mathcal{F}(\boldsymbol{\theta}) = |\mathcal{F}(\boldsymbol{\theta})_{\text{iter}} - \mathcal{F}(\boldsymbol{\theta})_{\text{iter-1}}|$

Figure 1.14: The VBPARAFAC algorithm.

1.5.5 PARAFAC combined with ICA

Many factor analysis models exist for the two-way array analysis of matrices. Among these models Independent Component Analysis has gotten much attention as, in many situations, it has proven efficient in finding the relevant components in the data [18].

Here I show how any two-dimensional factor analysis can be applied to create a PARAFAC decomposition using features from the two dimensional factor analyses. The method will focus on the independent component analysis as revealed in the ICAPARAFAC algorithm. Furthermore, I show that for three-way arrays and some four-way arrays the decomposition can become non-iterative.

ICAPARAFAC

First I'll define what I call Combined Independence, CI. Consider the multi-way array $\mathbf{X} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_N}$. That \mathbf{X} is combined independent in the modalities $n, n+1, \dots, N$, i.e. $CI_{n, n+1, \dots, N}$ means that the matricizing of \mathbf{X} into $\mathbf{X} \in \mathcal{R}^{I_1 I_2 \dots I_{n-1} \times I_n I_{n+1} \dots I_N}$ can be described by the model $\mathbf{X} = \mathbf{A}\mathbf{S}$ where the rows of \mathbf{S} are mutually independent, but this is not possible including less modalities in the CI. Notice, the CI applies to any combination of modalities as the order of the multi-way array modalities is only a matter of permutation. In the two-way array case, CI_2 correspond to the normal ICA model.

In ICAPARAFAC the combined independence is first identified. \mathbf{A} and \mathbf{S} is then found using an ICA algorithm from matricizing the multi-way array so that the columns of \mathbf{X} constitutes the modalities of combined independence. Consequently, $\mathbf{A} \in \mathcal{R}^{I_1 I_2 \dots I_{n-1} \times F}$, $\mathbf{S} \in \mathcal{R}^{F \times I_n I_{n+1} \dots I_N}$. To find the λ^{th} factors corresponding to each of the dimensions $I_n I_{n+1} \dots I_N$ unmatricize the λ^{th} row of \mathbf{S} , i.e. $\mathbf{S}^{(\lambda)} \in \mathcal{R}^{I_n \times I_{n+1} \times \dots \times I_N}$. Each dimensions λ^{th} factors can now be found by the best rank-one decomposition of $\mathbf{S}^{(\lambda)}$ using for example the SR1PARAFAC algorithm with one factor. To find the factors of each dimension in \mathbf{A} corresponding to the $n-1$ first modalities, two approaches can be used. The first approach is similar to finding the factors underlying \mathbf{S} ; unmatricize the λ^{th} column of \mathbf{A} to give $\mathbf{A}^{(\lambda)} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_{n-1}}$. Again find each dimensions λ^{th} factor by the best rank-one decomposition of $\mathbf{A}^{(\lambda)}$. The second approach is to find the remaining modalities factors using ALSPARAFAC on \mathbf{X} while holding the factors found underlying \mathbf{S} , i.e. the $n, n+1, \dots, N$ modalities fixed. Where the first approach gives the decomposition that the best describe \mathbf{A} found by the ICA, the second approach gives a better approximation to \mathbf{X} . The method is described in Figure 1.15.

ICAPARAFAC

Consider the model $x_{i_1 i_2 \dots i_N} = \sum_{\lambda=1}^F u_{i_1 \lambda}^{(1)} u_{i_2 \lambda}^{(2)} \dots u_{i_N \lambda}^{(N)}$

Identify the CI

Matricize \mathbf{X} so the columns of \mathbf{X} correspond to the modalities of the CI.

Solve $\mathbf{X}=\mathbf{A}\mathbf{S}$ using ICA

Find each $\mathbf{u}_\lambda^{(i)}$ for the CI dimensions by finding the best rank one decomposition of the unmatricized multi-way array corresponding to the λ^{th} row of \mathbf{S} .

Find \mathbf{u}_λ for the dimensions not in the CI, by either finding the best rank one decomposition of the unmatricized multi-way array corresponding to \mathbf{a}_λ , or by ALSPARAFAC on \mathbf{X} where $\mathbf{U}^{(i)}$ of each of the CI dimensions are kept fixed.

Figure 1.15: A PARAFAC model based on ICA.

The algorithm becomes very simple when the multi-way array \mathbf{X} is of few modalities. If \mathbf{A} or \mathbf{S} only holds one modality the factor of this modality is given directly by \mathbf{A} or \mathbf{S} . Furthermore, if \mathbf{A} or \mathbf{S} holds two modalities, the factors of each of these two modalities can be found from the first eigenvectors corresponding to the SVD solution as this is the same as the rank one decomposition. Finally, if \mathbf{X} is a three-way array and the number of modalities of CI is one, the two approaches to find \mathbf{A} yields the same results, see also Theorem 11 page 110.

Non-iterative methods for ICA such as the Molgedey-Schuster algorithm exists [18]. Consequently, decomposing any three-way arrays and four-way arrays where the number of modalities of the CI equals 2 can be done completely non-iterative when combined with the SVD method. As a result, the calculations needed to estimate the ICAPARAFAC parameters can be considerably reduced making the algorithm much faster than the ALSPARAFAC algorithm. Especially in multi-way analysis speed is an important issue as the amount of data tend to be tremendous. Therefore, the ICAPARAFAC algorithm just described seems very promising as long as combined independence can be assumed present in the data.

1.5.6 Algorithm relations

The rank one algorithm, SR1PARAFAC, corresponds to successively running an ALSPARAFAC algorithm with one factor, subtracting the found factors from the data, and finding the next factor from this subtracted dataset. The ALSPARAFAC algorithm corresponds to an EMPARAFAC model where the prior on the hidden variable \mathbf{S} is a delta function, making the E step resemble the M step of the algorithm. The weight Φ in the EM-PARAFAC model could also have been implemented by weighted regression in the ALSPARAFAC model. Finally, the EMPARAFAC algorithm is the special case of the VBPARAFAC algorithm where the parameters priors are assumed delta functions, see also Figure 1.16.

Whereas the ALSPARAFAC algorithm seeks to find an optimal solution in terms of explaining the most variance, i.e. reducing the sum of square error, the SR1PARAFAC algorithm seeks to consecutively find factors explaining the most of the variation. The EM method and VBEM method, however, seeks to optimize the likelihood of the observed data. This does not necessarily optimize the sum of square error as the priors affect the solution. However, the EM and VBEM method is expected to generalize well by not over fit the model to the data due to the priors restricting the variables. Although the main interest in this thesis will be a PARAFAC model that well explains the observed data in favor of the ALSPARAFAC, SR1PARAFAC and ICAPARAFAC rather than finding a model that generalizes well on new data in favor of the EM and VBEM algorithm, the latter models have only been included for completeness of PARAFAC methods. Furthermore, the statistical framework enables the evaluations of questions concerning the number of factors to include in the models by the ARD and Bayesian Information Criterion, rather than just relying on the Core Consistency Diagnostic. The Bayesian Information Criterion, BIC, has been derived for the EM-algorithm in Theorem 6 page 101, while BIC has been derived for any least square optimization algorithm in Theorem 7 page 102.

Finally, the ICAPARAFAC algorithm has been developed here to handle data that can be considered CI. Whereas the ALSPARAFAC isn't optimized to yield solutions insuring CI, this is achieved by the ICAPARAFAC algorithm.

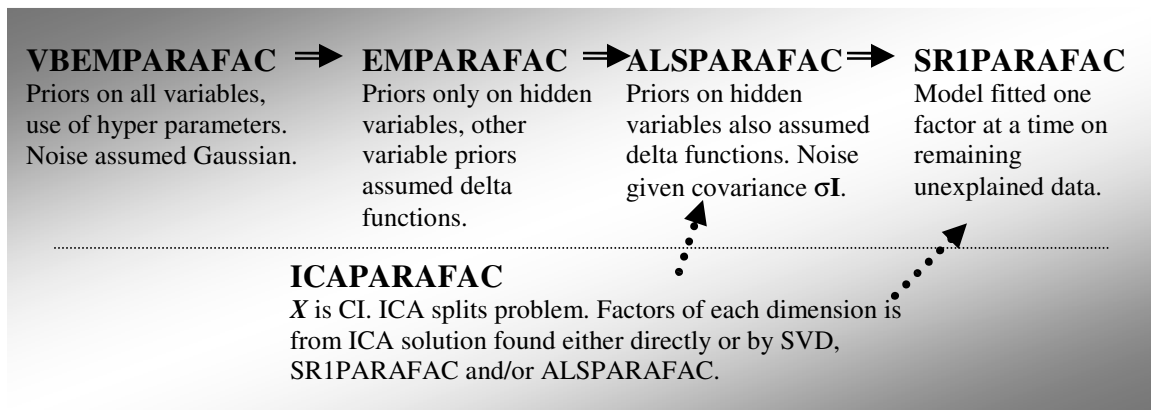


Figure 1.16: The relationship between the developed PARAFAC algorithms of this thesis.

As quoted in the beginning of this chapter, Abraham H. Maslow said: “If the only tool you have is a hammer, you tend to see every problem as a nail“. Hopefully, the tools described above will be adequate to the problems at hand.

2 EEG

“The human brain is the last, and greatest, scientific frontier. It is truly an internal cosmos that lies contained within our skulls. The more than 100 billion nerve cells and trillion supporting cells that make up your brain and mine constitute the most elaborate structure in the known universe.”

Joel Davis

The term electroencephalography (EEG), as we commonly use it refers to electrical activity measured at the scalp that arises from neurons in the brain. This includes activities that arise spontaneously or in response to sensory stimuli although the latter are more commonly known as ‘evoked response potentials’ (EP) [29]. Finally, the EEG of sensory stimuli timed to an event is referred to as event related potentials (ERP).

As it hasn’t been possible to experimentally identify the sources of the EEG signals for certain, many theories as to what constitutes the signal has arisen spanning from membrane quantum dynamical effects to K^+ fluctuations within the extracellular space. However, in this thesis only the theory which is given the most recognition in the literature will be described. This theory is primarily based on the theoretical framework explained by Paul Nunez [28], [29].

2.1 Dipoles

The electric force between two charges is defined by the well known coulombs law:

$$F = \frac{1}{4\pi\epsilon_0} \frac{q_1 q_2}{R^2} \quad \text{eq. 2.1}$$

From coulombs law the electric field at a point \mathbf{r}_1 due to a point charge q located at \mathbf{r}_2 can be derived:

$$\mathbf{E}(\mathbf{r}_1, \mathbf{r}_2) = \frac{1}{4\pi\epsilon_0} \frac{q}{\|\mathbf{r} - \mathbf{r}_i\|^2} \frac{(\mathbf{r}_1 - \mathbf{r}_2)}{\|\mathbf{r} - \mathbf{r}_i\|} = \frac{1}{4\pi\epsilon_0} \frac{q(\mathbf{r}_1 - \mathbf{r}_2)}{\|\mathbf{r} - \mathbf{r}_i\|^3} \quad \text{eq. 2.2}$$

Consequently, if there are n charges q_i located at various positions \mathbf{r}_i they produce an electric field at \mathbf{r} given by:

$$\vec{\mathbf{E}}(\mathbf{r}) = \sum_{i=1}^n \vec{\mathbf{E}}(\mathbf{r}, \mathbf{r}_i) = \frac{1}{4\pi\epsilon_0} \sum_{i=1}^n \frac{q(\mathbf{r} - \mathbf{r}_i)}{\|\mathbf{r} - \mathbf{r}_i\|^3} \quad \text{eq. 2.3}$$

The cornerstone of the understanding of electric field behavior comes from Maxwells equations [12]:

Maxwells equations

$$\nabla \cdot \vec{\mathbf{E}} = \rho_t / \epsilon_0 \quad (1)$$

$$\nabla \times \vec{\mathbf{E}} = -\partial \vec{\mathbf{B}} / \partial t \quad (2)$$

$$\nabla \cdot \vec{\mathbf{B}} = 0 \quad (3)$$

$$\nabla \times \vec{\mathbf{B}} = \frac{\vec{\mathbf{J}}_t}{c^2 \epsilon_0} + \frac{1}{c^2} \frac{\partial \vec{\mathbf{E}}}{\partial t} \quad (4)$$

eq. 2.4

$\vec{\mathbf{B}}$ is the magnetic field, ρ_t is the total charge density and $\vec{\mathbf{J}}$ is the current density.

From Maxwell's equation 2 and 4 it is seen that a change in the magnetic field results in a change in the electric field and vice versa. However, when field frequencies in the brain are less than in the order of MHz the effect of the interaction between magnetic and electric field becomes negligible [28] and the electrical potential Φ can solely be determined by the electric field:

$$\vec{\mathbf{E}} = -\nabla \Phi \Rightarrow$$

$$\Phi(\mathbf{r}) = \frac{1}{4\pi\epsilon_0} \sum_{i=1}^n \frac{q_i}{\|\mathbf{r} - \mathbf{r}_i\|} \quad \text{eq. 2.5}$$

The potential of a monopole, dipole and quadrupole is shown in Figure 2.1. In general the potential due to all charges or current sources can be expressed as the following series of terms called a multipole expansion [28]:

$$\Phi = (\text{monopole contribution}, r^{-1}) + (\text{dipole contribution}, r^{-2}) + (\text{quadrupole contribution}, r^{-3}) + \dots = \sum_{i=0}^{\infty} (2^i - \text{pole contribution}, r^{-(i+1)}) \quad \text{eq. 2.6}$$

Where r is the distance to the center of the pole source

In Figure 2.1 the potential for a monopole, dipole and quadrupole is drawn.

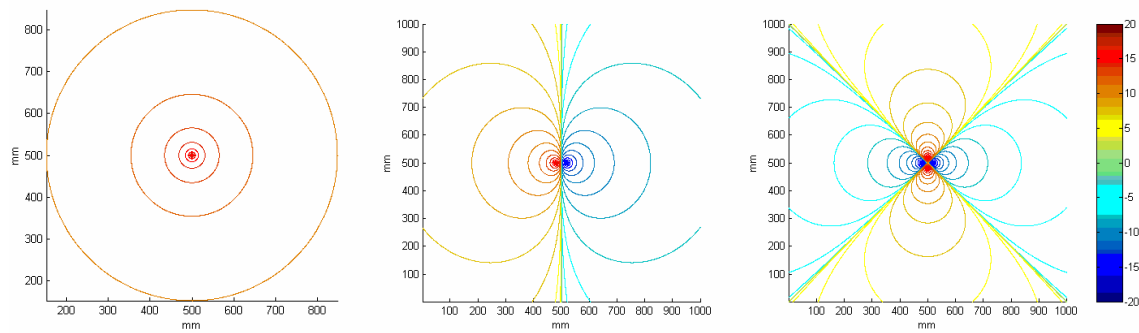


Figure 2.1: The potential for a monopole, dipole and quadrupole (units $Volt^{0.1}$)

Combining Maxwell's equation 1 and 4 yield;

$$\nabla \cdot \bar{\mathbf{J}} + \frac{\partial \rho}{\partial t} = 0 \quad \text{eq. 2.7}$$

i.e. charge is neither created nor destroyed – charge is conserved. As the net charge is conserved during brain activity the monopole contribution vanishes. As all other contributions than the dipole quickly drops to zero with distance to the sources, the only contribution believed to significantly contribute to the EEG is that of the dipole. The potential of a dipole can be approximated to be, see also the derivation in Theorem 8 page 103:

$$\Phi_{dipole}(\mathbf{r}, \mathbf{d}) = \frac{1}{4\pi\epsilon_0} \frac{q\|\mathbf{d}\|\cos\theta}{\|\mathbf{r}\|^2} \quad \text{eq. 2.8}$$

Where \mathbf{d} is the vector going from the negative charge to the positive of the dipole and θ is the angle between the vectors \mathbf{r} and \mathbf{d} .

As seen from the potential lines of the dipole on Figure 2.1 no or little potential is found oblique to the dipole, this is confirmed by eq. 2.8 as $\cos(90^\circ) = 0$. Consequently, the EEG can only pick up signals from dipoles radial to the electrodes, as seen on Figure 2.2.

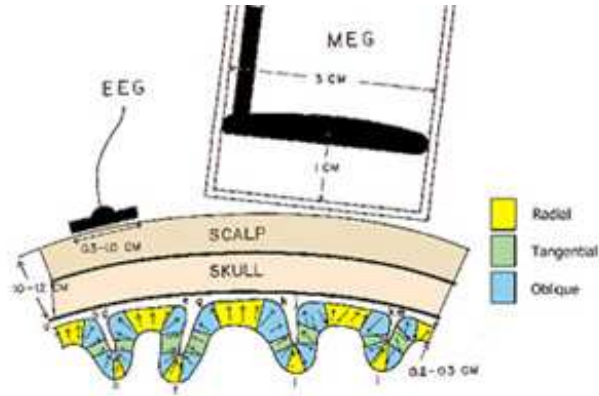


Figure 2.2: EEG can only pick up potentials from radial dipoles, whereas MEG pick up potentials from oblique dipoles (induced magnetic field is angular to current flow) (taken from [36]).

Two effects are seen with the generation of potential; capacitive current and resistive current. Nunez gives a proof that if a media is linear in both the dielectric and conductive sense i.e. the capacitor (polarization) and current is proportional to the electric field the following approximation is in general valid [28]:

$$\frac{\text{Capacitive current}}{\text{Resistive current}} \approx 0.02 \quad \text{eq. 2.9}$$

Consequently, the capacitive current is in the following considered negligible. As capacitive current is ignored eq. 2.7 reduces to:

$$\nabla \cdot \vec{\mathbf{J}} = 0 \quad \text{eq. 2.10}$$

Imagine two regions where all current sources are located in the first region. Expressing the current density at the boundary between the two regions as the sum of an Ohmic current (i.e. Ohms law states that current is linear to potential also makes current linear with electric field strength) and a source current $\vec{\mathbf{J}}_s$, yields:

$$\vec{\mathbf{J}} = \sigma \vec{\mathbf{E}} + \vec{\mathbf{J}}_s \quad \text{eq. 2.11}$$

Where σ is the conductivity

Using first part of eq. 2.5, eq. 2.9 and eq. 2.10 this gives:

$$\nabla \cdot \sigma \nabla \Phi = \nabla \cdot \vec{\mathbf{J}}_s \equiv \mathbf{j}_s \quad \text{eq. 2.12}$$

Where \mathbf{j}_s has the dimension amperes pr. m³

Using first part of eq. 2.5, and Maxwell's first equation we also find:

$$\nabla \cdot \sigma \nabla \Phi = \mathbf{j}_s \Leftrightarrow \nabla \cdot \nabla \Phi = \nabla \cdot \vec{\mathbf{E}} = \mathbf{j}_s / \sigma \quad \text{eq. 2.13}$$

As \mathbf{j}_s is the source density, using last part of eq. 2.5 with Maxwell's first equation gives:

$$\Phi(\mathbf{r}, t) = \frac{1}{4\pi\sigma} \sum_{i=1}^n \frac{I_i(t)}{\|\mathbf{r} - \mathbf{r}_i\|} \quad \text{eq. 2.14}$$

Where $I_i(t)$ is the i^{th} current source at time t at position \mathbf{r}_i . eq. 2.14 is very important as it translates current sources into EEG measurable potentials. eq. 2.13 can also be written:

$$\nabla^2 \Phi = \mathbf{j}_s / \sigma \quad \text{eq. 2.15}$$

From eq. 2.15 it is seen that the current source densities can be estimated by taking the Laplacian of the measured EEG potential. Consequently, taking the Laplacian of the potential is believed to improve the spatial resolution of the EEG signal as it estimates the current sources and sinks. In practice, calculating the surface Laplacian of the EEG requires some form of interpolation of the EEG to estimate the signals between the electrodes. Often splines in the form of Legendre polynomials are used for these interpolations.

2.2 EEG and coherence

Coherence is a measure of the synchrony between sources. Let the i^{th} source signal at time l be defined as E_{il} . Furthermore, let the spectral density function $q_{ii}(f)$, also called the power spectrum, and cross spectral density function $q_{ij}(f)$, also called the cross power spectrum, for two sources be given by their Fast Fourier Transform (FFT). The coherence function between the sources i and j , $\gamma_{ij}^2(f)$, is then defined by:

$$q_{ij}(f) = \frac{2\Delta t}{N} \left(\sum_{l=0}^{N-1} E_{il} \exp(-i2\pi f \Delta t) \sum_{l=0}^{N-1} E_{jl} \exp(i2\pi f \Delta t) \right) \Rightarrow \quad \text{eq. 2.16}$$

$$\gamma_{ij}^2(f) = \frac{|q_{ij}(f)|^2}{|q_{ii}(f)| \cdot |q_{jj}(f)|}$$

If the sources are coherent at the frequency f then $\gamma_{12}^2(f) = 1$, the sources are considered incoherent if $\gamma_{12}(f) = 0$. Notice how the coherence is defined as the squared cross spectral correlation between the two sources.

Let M be the amount of coherent sources, N the amount of incoherent sources. The relative contribution of coherent to incoherent sources in the EEG is estimated to be M / \sqrt{N} [28].

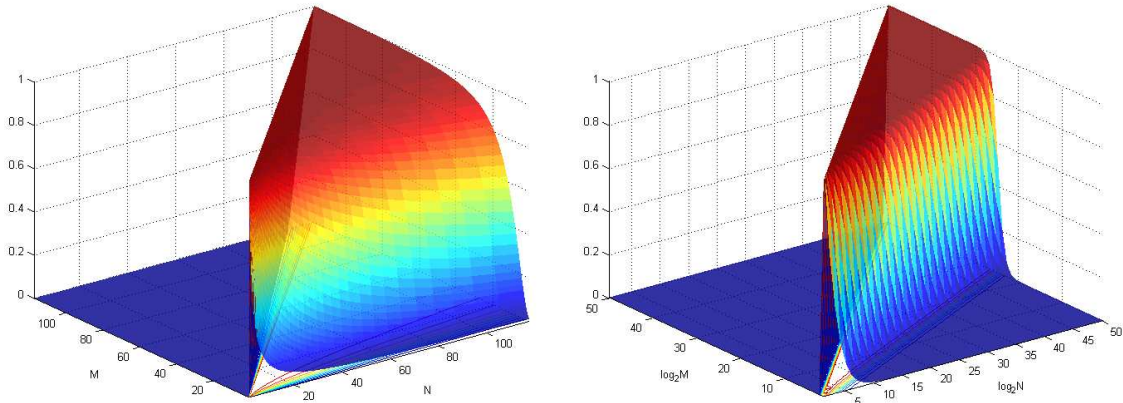


Figure 2.3: The percentage of measured signal picked up from the EEG of coherent sources m versus active sources N . Clearly, the coherent sources dominates the recorded EEG-signal even when only few of the active sources are coherent. (Notice figure not valid for $M > N$).

As seen on Figure 2.3 the coherent sources dramatically dominate the recorded EEG signal even though the coherent neurons are far less numerous than the incoherent neurons. The EEG signal is therefore believed to originate from the synchronous firing of parallel oriented neurons. 65-75% of the cortical neurons are oriented perpendicular to the cortical surface. These pyramidal cells have large amounts of interconnections, so it seems as if a relatively high degree of synchrony can be obtained from these neurons. Furthermore, the pyramidal neurons of the neocortex are the neurons closest to the scalp surface. According to eq. 2.6 the distance to the scalp causes pyramidal neurons to be the least reduced. Therefore, the EEG is believed to mostly originate from pyramidal cells in the neocortex [29]. Amplitude changes in the EEG with physiological state becomes, in this framework, a result of changes in the number of synchronously active neurons [28]. This is supported by current/source density studies that indicate that the pyramidal cells of layer III and V are the principal source of the EEG [35].

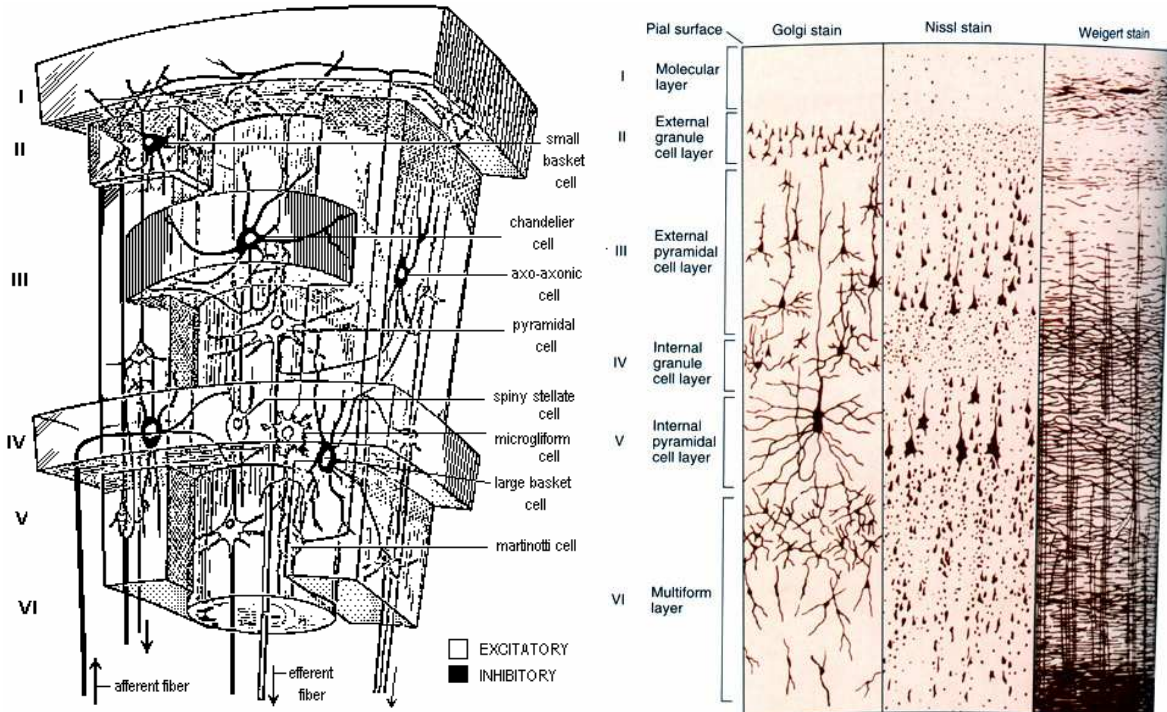


Figure 2.4: Model of the layers of neocortex and findings using Golgi, Nissl and Weigert stains. From the stains the parallel orientation of the pyramidal cell perpendicular to the cortex surface becomes evident. The pyramidal neurons constitute 70% of the neurons of cortex and posses 10^3 to 10^5 synaptic contacts. Nunez believes that each electrode pick up signal from around 30 to 40 macro-columns where each macro-column contains 10^5 - 10^6 neurons [29]. (Taken from [16],[37]).

The possible roles of cortical coherence are believed to be

1. **Blocking**
Focus on a particular modality
2. **Matching**
Transfer of data from one group of neurons to the next
3. **Binding**
Synchronous activity between neuron groups
4. **Plasticity**
Transfer of function from one neuron group to another

When analyzing the ERP the following measures are very useful [8]:

$$\begin{aligned}
 ERSP(f, t) &= \frac{1}{n} \sum_{k=1}^n |F_k(f, t)|^2 && \text{Event - related spectral perturbation} \\
 ITPC(f, t) &= \frac{1}{n} \sum_{k=1}^n \frac{F_k(f, t)}{|F_k(f, t)|} && \text{intertrial phase coherence} \\
 ERPCOH^{a,b}(f, t) &= \frac{1}{n} \sum_{k=1}^n \frac{F_k^a(f, t) F_k^b(f, t)^*}{|F_k^a(f, t) F_k^b(f, t)|} && \text{Event related cross coherence}
 \end{aligned}
 \tag{eq. 2.17}$$

Where n is the number of trials. Notice that the square of the ERPCOH corresponds to the coherence measure as defined in eq. 2.16.

Although coherence may appear to be an ideal measure of brain function, interpretations of experimental EEG coherence are often confounded by technical limitations. Scalp coherence between electrode sites closer than about 8 to 10 cm is typically large or moderate due only to passive current spread, volume conduction and reference electrode effects, even when the underlying cortical sources are uncorrelated [29]. This problem can to some degree be circumvented by taking the surface Laplacian as described in eq. 2.15.

2.3 Synaptic potentials and action potentials

eq. 2.8 can also be expressed in terms of current sources and sinks. Let the dipole consist of a current source at time t , $I(t)$ and a current sink $-I(t)$. Let \mathbf{d} be the vector from the sink to the source and \mathbf{r} the vector to the center of the dipole, then:

$$\Phi_{dipole}(\mathbf{r}, \mathbf{d}, t) = \frac{1}{4\pi\epsilon_0} \frac{I(t) \|\mathbf{d}\| \cos \theta}{\|\mathbf{r}\|^2} \quad \text{eq. 2.18}$$

Potential differences recorded in the EEG are therefore believed to derive from coherent dipoles constituting current sources and sinks. These are thought to originate from two different processes [28]; synaptic potentials and action potentials, the latter also referred to as sodium-potassium spikes.

2.3.1 Synaptic potentials

An action potential in the presynaptic axon activates a chemical agent (transmitter) which diffuses across the synaptic cleft into the subsynaptic membrane. If the synapse is excitatory, the effect of the chemical transmitter is to increase (excitatory post synaptic potential, EPSP) or decrease (inhibitory post synaptic potential, IPSP) the permeability of the subsynaptic membrane to positive/negative ions which flow through the local surface of the membrane. The current flows across the membrane, through the intracellular fluid, back across the membrane at more distant locations, and finally back to the synapse to complete a closed loop, see Figure 2.5. This loop acts as a current source-sink dipole [28].

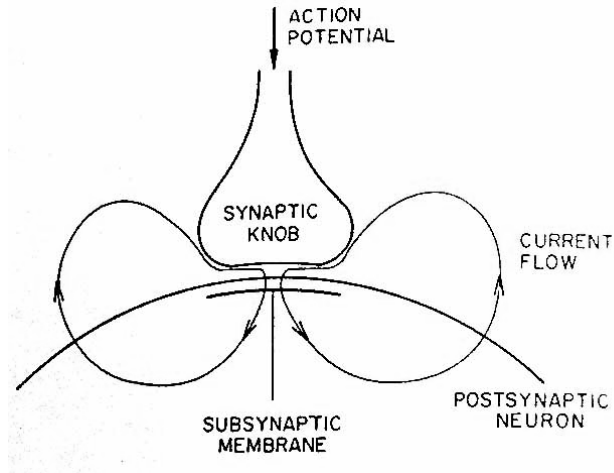


Figure 2.5: The current flow of a synaptic potential. Positive/negative ions flow through the local surface of the membrane. The current flows across the membrane, through the intracellular fluid, back across the membrane at more distant locations, and finally back to the synapse to complete a closed loop. (taken from [28]).

2.3.2 Action potentials

The action potential arises when a stimulus opens a few sodium channels. As a result, a net influx of sodium starts due to the concentration gradient of sodium and accelerates as the depolarization makes more sodium channels open. Eventually the increasing depolarization causes potassium channels to open while the sodium channels at this point closes. An outflow of potassium is caused by potassium's concentration gradient repolarizing the cell. Eventually the potassium channel closes and the resting potential is restored due to the sodium-potassium-pump, see also Figure 2.6.

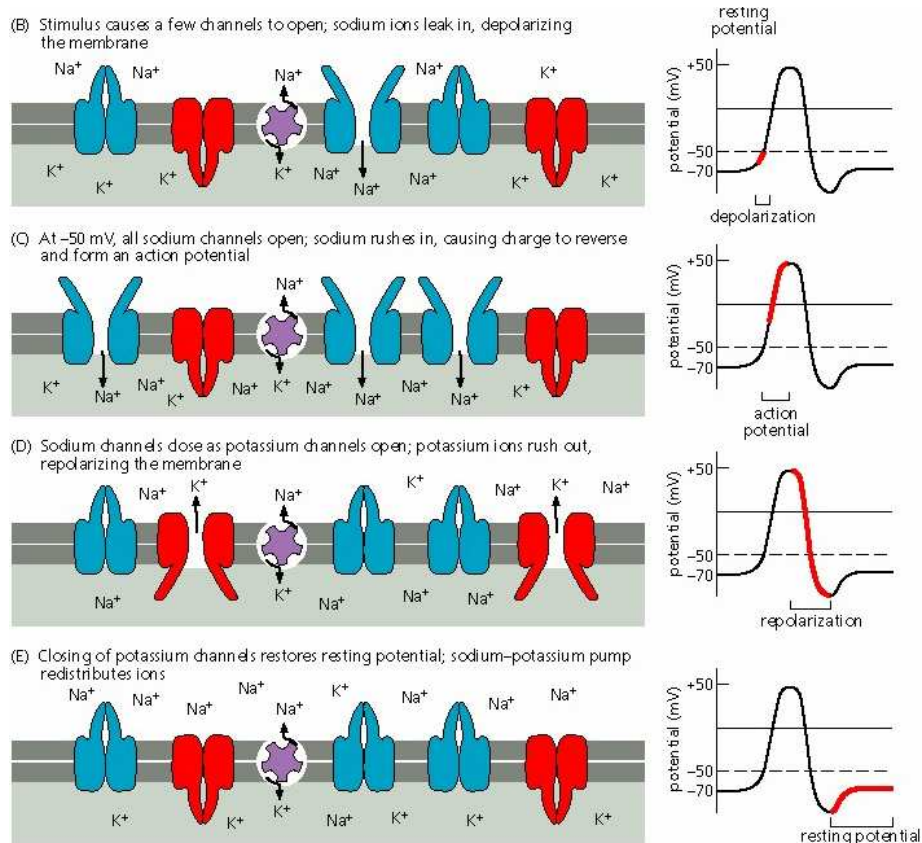


Figure 2.6: The generation of an action potential (adapted from [39])

From Figure 2.6 it is seen that the sodium currents are followed by opposite delayed potassium currents. Nonetheless, this potassium current is unable to counteract the potential generated from the sodium current as $\sigma_{Na^+} < \sigma_{K^+}$ both interior and exterior the cell as revealed on Figure 2.7.

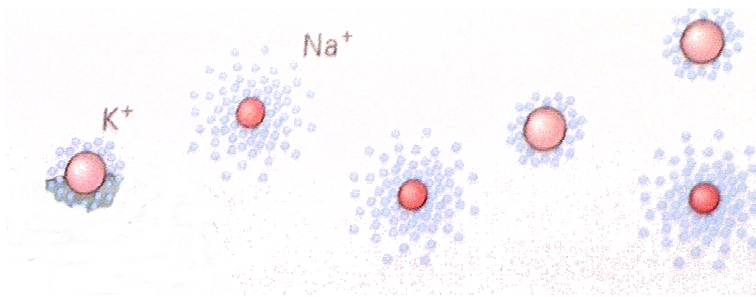


Figure 2.7: The smaller an ion is the more highly localized is its charge and the stronger its effective electric field. As a result, smaller ions attract water more strongly. Consequently, because of its larger water shell, Na^+ behaves as if it is larger than K^+ making it less mobile, i.e. giving Na^+ lower conductivity than K^+ [16]. (Figure adapted from [16]).

2.3.3 Synaptic Potentials versus Action Potentials

From eq. 2.18 it is seen that the dipole potential is linear to the distance between the current source and sink $\|\mathbf{d}\|$. As myelination of neurons dramatically increases the distance between the current sources and sinks, most EEG signals from action potentials is believed to originate from pyramidal cells of heavy myelinated axons [28]. However, as the propagation of current is much slower in unmyelinated neurons coherence is more easily achieved in the unmyelinated regions in favor of the signal originating mostly from the synaptic potentials from the dendritic parts of the pyramidal cell [11]. On the other hand, the dendrites in general aren't as well aligned as the axon's - compare in Figure 2.4 the Golgi stain emphasizing dendritic trees with the weigert stain emphasizing myelinated axonal fibers. However, the synaptic potential has duration between 5 ms. – 20 min. whereas the action potential only has duration of 1-10 ms. Consequently, coherence is more easily achieved for the synaptic potential. As a result the greatest contribution to the EEG is believed to be that of the synaptic potential [9],[16],[28].

Table 1: Processes of potential generation favoring/ disfavoring coherence.

Measured EEG	Action Potential	Synaptic Potential
Advantage	Myelinated axons → Large distance between current sources	Long duration (5 ms-20 min) → Coherence easy Both excitatory and inhibitory effects.
Disadvantage	Short duration (1-10 ms.) → Coherence difficult	Dendrites not so well aligned → Coherence difficult

2.4 Features of the EEG and ERP

The EEG has been decomposed into a series of fixed broad spectral bands based, unfortunately, on history and discovery more than on a theoretical framework [24]. These bands are described in Figure 2.8. In general, the frequency of brain oscillations is negatively correlated with their amplitude [28],[30].

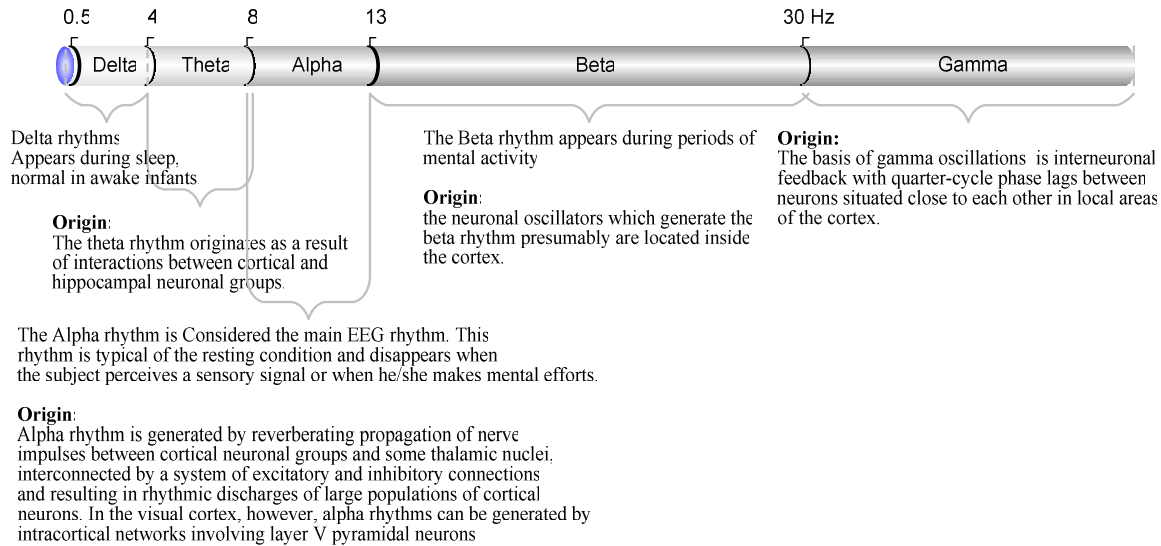


Figure 2.8: EEG rhythms and their believed origins (summary of [9],[16], [28], [34])

2.4.1 Event Related Potentials, ERP

Event Related Potentials (ERP) is the measured EEG signal timed to sensory stimuli. The EEG signal can be split into induced activity, evoked activity and noise. The induced and evoked activities are both believed to be generated from thalamic relay cells. An internally or externally paced event results not only in the ‘evoked’ generation of an evoked response potential (EP) but also in an ‘induced’ change in the ongoing EEG/MEG in form of an event-related desynchronization (ERD) or event-related synchronization (ERS) [30]. The EP represent the responses of cortical neurons due to changes in afferent activity, while ERD/ERS reflect changes in the activity of local interactions between main neurons and interneurons that control the frequency components of the ongoing EEG [30].

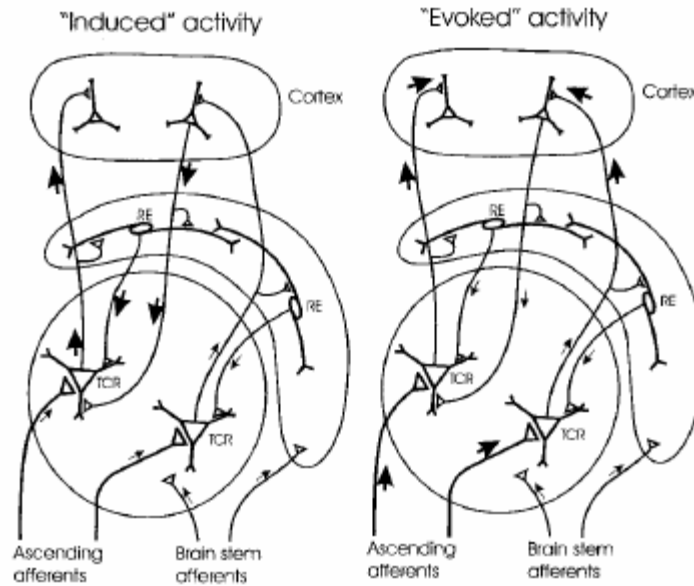


Figure 2.9: Schema for the generation of the highly frequency specific induced (ERD/ERS) and the evoked (EP) activity. TCR: thalamic relay cells, RE: reticular thalamic nucleus [30]. Where evoked activity is mainly controlled by the ascending afferents the induced is controlled by the reticular thalamic nucleus.

As revealed in Figure 2.9 of principal concern is the reticular activating system, a complex and diffuse system projecting from the brainstem to the cortex, which provides both inhibitory and excitatory inputs. This subsystem is itself under control from the cortex, as well as from collaterals of sensory pathways. The activating system is most crucially concerned with maintenance of the waking state, desynchronization of the EEG, direction of attention and governance of motivation [35]. Although EEG is believed primarily to stem from dipoles of synaptic potentials firing in synchrony mostly due to the reticular activation system, where and what exactly generates the signals remains unclear.

Evoked oscillations exhibit a strict phase-locking to the experimental event (e.g. stimulus presentation) across trials. Hence, they can be extracted from the averaged ERP, e.g. by filtering or by the ITPC of a wavelet analysis. On the other hand induced oscillations are (by definition) not at all phase-coupled to a stimulus, and show a certain degree of phase-jittering. Therefore, by averaging across trials these oscillations will cancel out completely and hence are only detectable by appropriate ways of analysis, e.g. by a single trial based wavelet analysis with subsequent averaging [6], such as the ERSP.

The components of the ERP are labeled by latency and polarity. A positive component at 100 ms is called 'P100' and a negative deflection at 200 ms is a 'N200', see Figure 2.10. Furthermore, the event related potentials that are determined by physical aspects of the stimulus are labeled, 'exogenous' whereas higher order processing are labeled 'cognitive'.

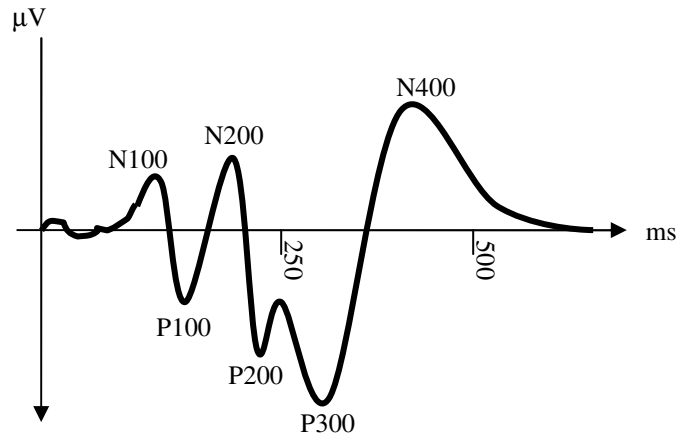


Figure 2.10: Typical significant components of the ERP

(notice: negativity is due to historical reasons directed upwards so that current flowing out of the scalp is up)

Many ERP experiments are based on the oddball paradigm: Two stimuli are presented, one of which is an infrequent target. When the target is discriminated from the other stimuli with attention, the P300 is of greater amplitude [9]. This has made especially the P300 a very interesting component of the ERP. The various most common ERP components are described in Table 2. The table is only a guideline. Experts do not agree on how to interpret all components and the components are difficult to generalize over different experimental paradigms.

Table 2: Features of ERP components

	Spike	Description
Exogenous process	P100	<p>Generation Appear to be generated in the primary receiving areas of the brain.</p> <p>Identification Both have been found to occur 100 ms after presentation of a visual stimulus. Found in the receiving areas of the brain.</p>
	N100	
	N200	<p>Generation Increases in amplitude to task deviant stimuli.</p> <p>Identification Located Posterior</p>
	P200	<p>Generation Increases to novel stimuli</p>
Cognitive process	P300b	<p>Generation Attention dependent. Monotonous inverse relationship between amplitude and stimulus probability. Negatively correlated with speed of information processing as indexed by reaction times – the faster speed of processing the earlier latency. The latency increases with the time the subject needs to distinguish the rare stimulus. The amplitude increases with rarity of the stimulus and to some extent with stimulus intensity. The amplitude of the P300 has been shown to be inversely proportional to stimulus presentation probability and directly to task complexity. P300 amplitude has been shown to be inversely related to a prior probability and is influenced by the sequence of immediately preceding events, i.e., sequential event structure. Unrelated to the specific sensory areas but probably related to the parietotemporal association cortex and subcortical structures such as hippocampus and thalamus. Occurs to infrequent non-target stimuli</p> <p>Identification: In normal young adults a positive wave over the Centro-parietal scalp is seen. Can occur at any point between 280-800 ms. Posterior scalp distribution with maximum at Pz, P300 latency is age dependent, being longer in children, progressively decreasing until 18 before increasing by 1.25 ms per year In adult. Aging increases the latency, decrease the amplitude and cause forward shift in the distribution of P300</p>
	P300a	<p>Generation Mostly concerned with novel stimuli. The more the stimulus is known the more P300a approaches P300b. Its features are Similar to P300b.</p> <p>Identification Shorter latency (\approx250 ms) than P300b. Situated more frontal-central than P300b and habituates rapidly. In different modalities it has also been reported as Centro-parietal.</p>
	N400	<p>Generation Found in numerous studies that have employed a task in which items were presented sequentially and subjects were asked to respond if the stimulus was unmatched (incongruent) or matched (congruent) with the preceding items. It has been proposed that it represent the associated activation of neural networks basic to stimulus integration.</p> <p>Identification Located in frontal area, and is larger in non-matched items than matched items.</p>

2.4.2 Noise

EEG potentials are measured as the difference between two points, one on the scalp where EEG effects are strong and one (the reference electrode) hopefully isolated from these effects. Some commonly used reference sites are Cz, earlobes, mastoids (bone right behind ear), tip of nose and average reference ("reference free"). Earlobes or mastoids are generally linked either physically or mathematically in order to maintain symmetry. The average reference uses the constraint that the sum of the potentials over a spherical surface is zero and requires fairly high density recording (~128 channels). It can be improved by estimating potentials for the inferior spherical area. Spherical spline interpolation is sometimes used for these estimates. The 10-20 international system is a standardized system to place the electrodes. It relies on taking measurements between certain fixed points on the head. The name 10-20 refers to the fact that the electrodes used to be placed at points 10% or 20% of these distances [40]. Today other fractions are used, but the name 10-20 has been kept. How the recorded signal is referred can have crucial impact on the noise in the EEG as a very noise full reference can mess up the signal of all recorded channels.

Many sources of noise disturb the EEG/ERP signal. Of primary concern is muscular action due to eye movement, hearth beat etc. Furthermore, noise from electronic devices can have a great impact especially in the 50 Hz range where most electronic devices, at least in Europe, operate. Finally, volume conductance is another huge problem when dealing with the EEG/ERP signal. Both heartbeat and eye movement can be reduced by correcting the signals from recording sites by the eye (EOG) and heart (ECG). Instead of specific recordings from these sights, recordings of the influence of eye and cardiac activity can also be identified prior to the ERP experiment. As previously mentioned, volume conductance can be reduced by taking the surface Laplacian.

2.4.3 EEG/ERP and PARAFAC

Extracting the correct features of the EEG is of crucial importance. Well used methods of analyzing EEG data are; Wavelet Analysis, Neural Networks Analysis, Blind Deconvolution and Source Separation Methods such as Principal Component Analysis (PCA) and Independent Component Analysis (ICA).

In this thesis, an approach described by Martinez-Montes et al. [24] and Miwakeichi et al. [25] where sources are separated using parallel factor analysis, PARAFAC, will be used. Although Harshmann in his original paper in 1970 [13] suggested the use of the PARAFAC model on EEG the use has been very limited. In 1988 Möcks [26] and in 1991 Field [10] used the model on the ERP to decompose the space-time-subject. In 1985 Cole et al.[7] used it on the ongoing EEG in a way similar to that of Miwakeichi et al.

Martinez-Montes and Miwakeichi et al. used the PARAFAC model to extract features of the ongoing EEG. They proved that PARAFAC was capable of successfully identifying the theta and alpha atoms of a cognitive task and showed furthermore the algorithms ability to identify eye blinks. In this thesis the PARAFAC decomposition will instead be applied to the ERP. The PARAFAC model seems plausible as the EEG/ERP-signals in several ways can be considered multi-way arrays as seen on Figure 2.11.

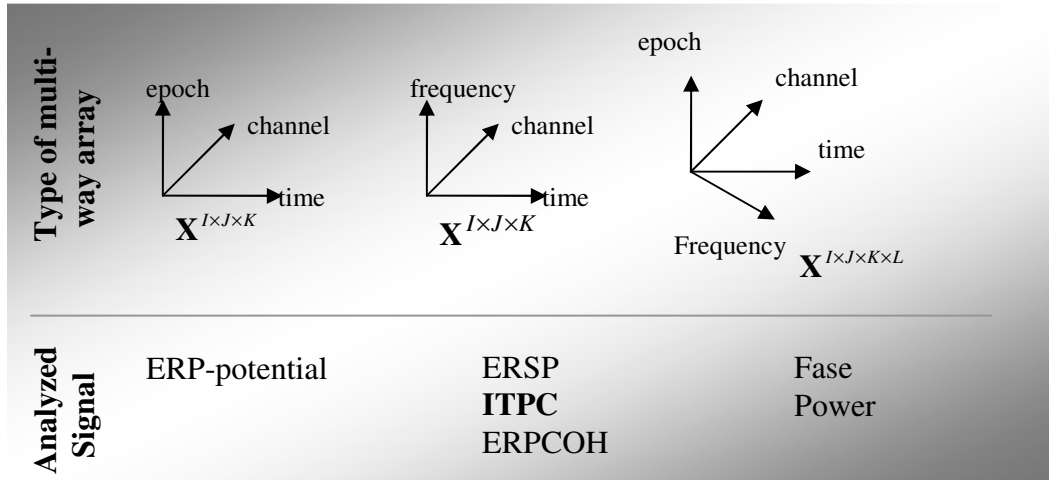


Figure 2.11: Different forms of multi-way arrays arising from the EEG.

More modalities could easily be added to Figure 2.11 denoting for example the analyzed subjects or the various conditions under which the data has been recorded. The PARAFAC analysis of the epoch averaged ERP-potential given by $channel \times time \times subject$ and the ITPC is in the following of great interest.

From eq. 1.29 we had $x_{ijk} = \sum_{\lambda=1}^F a_{i\lambda} b_{j\lambda} s_{k\lambda}$, in terms of the ITPC; \mathbf{a} denotes the component pertaining to the topography, i.e. channels, \mathbf{b} the frequency component and

s the component describing the temporal development. The model seems plausible as the ERP is believed to be a non-stationary process requiring a change in time of the factor proportions, which by the model is insured by $s_{k\lambda}$. eq. 1.29 can be restated as

$$x_{it} = \sum_{\lambda=1}^F a_{i\lambda} c_{t\lambda}, \text{ where } c_{t\lambda} = b_{j\lambda} s_{k\lambda} \text{ and } L = J \cdot K. \text{ This corresponds to the normal two-}$$

dimensional factor analysis model. Performing ICA on this model makes the assumption that the combined frequency-time components $d_{j\lambda} s_{k\lambda}$ are mutually independent, i.e.

$CI_{2,3}$. The goal in the analysis of the ITPC is both to separate the multi-way array into factors that relate to different time-frequency components favoring the ICAPARAFAC algorithm but also to do data exploration favoring ALSPARAFAC as it keeps the most of the variation.

3 DATA Analysis

"Absence of evidence is *not* evidence of absence!"

-unknown

3.1 Simulated data

To evaluate the ability of the PARAFAC algorithms to find the components of real data, the developed methods were tested on simulated data. A 32 channel EEG sampled at 500 Hz was generated and added 50 Hz oscillations of amplitude 0.8 on all channels mimicking electronic noise. Two burst of 35 Hz sinusoidal oscillations with an amplitude of 1.0 were placed in channel 30,31 and 32 at the posterior areas resembling occipital gamma activity while one burst of 25 Hz oscillation with amplitude 1.5 were generated simultaneously at each ear at channel 11 and 15. Finally, normal distributed random noise of power 1.0 was added to all channels. The data was transformed using a complex Morlet wavelet with bandwidth parameter 2 and center frequency 1, and the power of this wavelet transformed signal analyzed. The three PARAFAC factors shown in Figure 3.3 were expected to be found from the data. On Figure 3.1 the simulated EEG data is revealed and the corresponding power of the wavelet transform is seen on Figure 3.2.

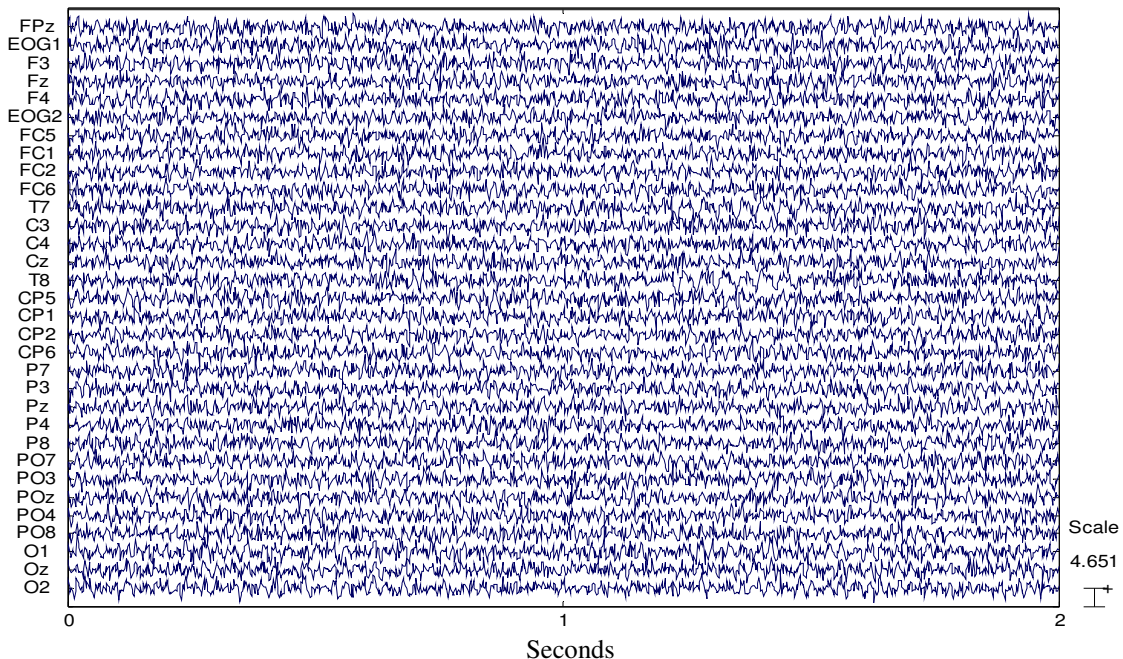


Figure 3.1: The simulated EEG-data.

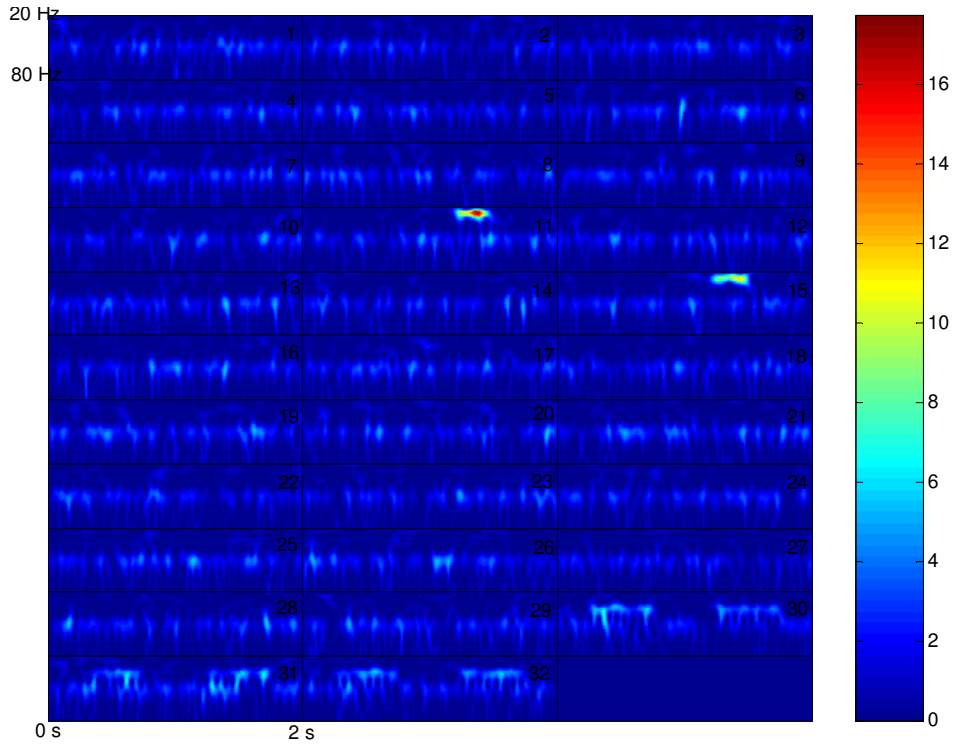


Figure 3.2: The power of the complex Morlet wavelet transform on each of the 32 channels of the simulated data.

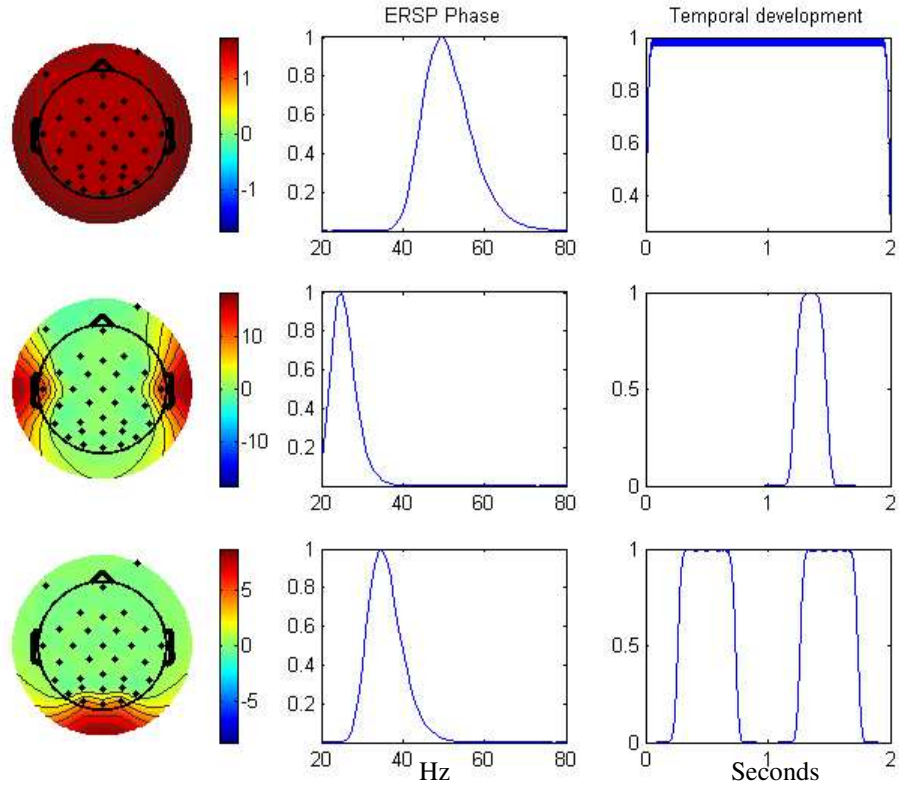


Figure 3.3: The true factors of the simulated data.

The raw simulated data was first analyzed using Independent Component analysis by the ‘runica’ standard method in EEGLAB [8]. As revealed in Figure 3.4 none of the independent components solely captures any of the underlying three factors. Especially the same 50 Hz oscillation present in all channels was split into individual components. Consequently, the independent component analysis did not seem efficient in accessing the various factors present in the data. Furthermore, no clear indication of the time points at which the factors were present was given by the ICA-decomposition as it is irresolvable from the EEG of the components. Thus, without any frequency information the ICA algorithm was unable to identify the factors.

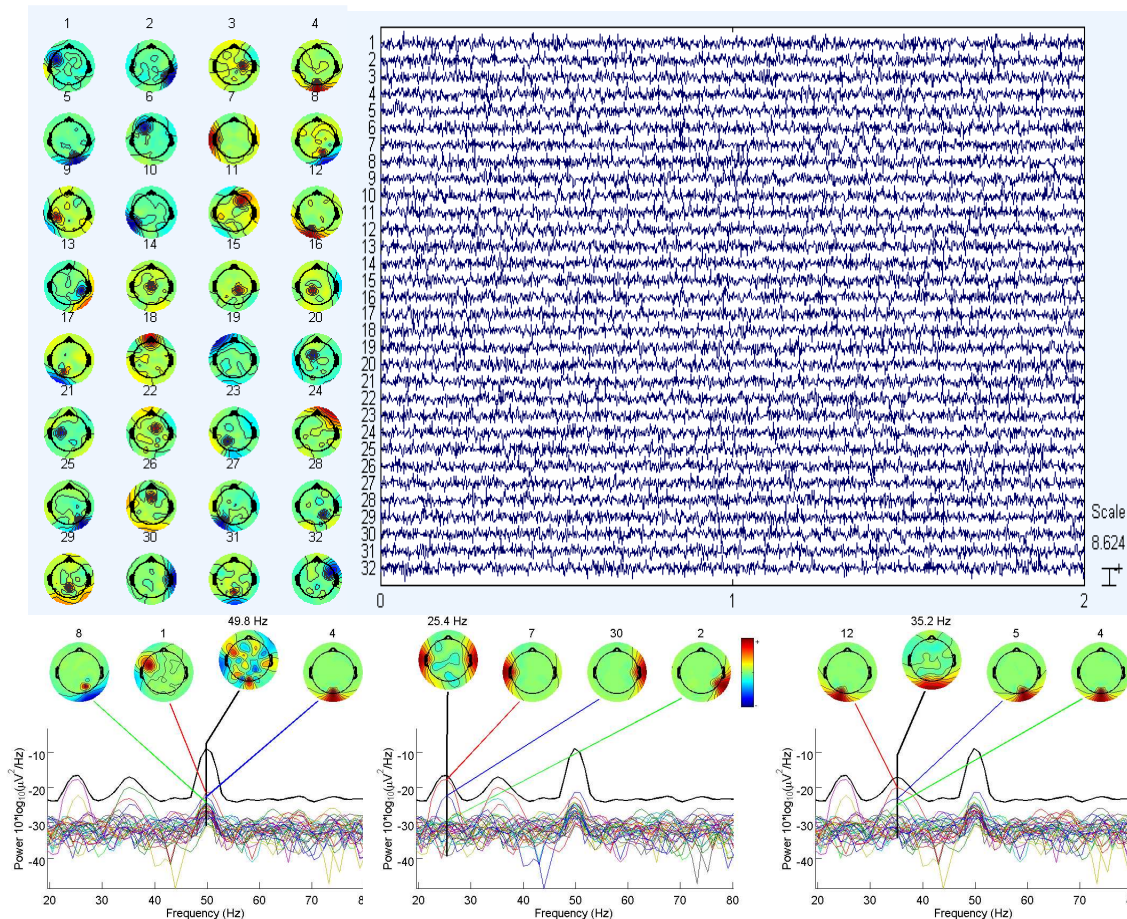


Figure 3.4: Top panel; the component map and time series of all 32 independent components. Lower panel; the maps of the three independent components contributing the most at 50 Hz, 25 Hz and 35 Hz to the specter of the EEG along with the summed map of the three components. Clearly, the ICA decomposition hasn’t been able to identify the true components of the data.

The developed PARAFAC models were then tested in their ability to access the components. In the ICAPARAFAC model $CI_{2,3}$ was assumed, i.e. a combination of the time and frequency dimensions were thought independent. As non-negative solutions were desired, a non-negative matrix factorization (NMF) was compared to a decomposition based on SVD for the rank one decomposition. The NMF was optimized in a least square sense as described in [22]. Although, non-negative ICA algorithms would

be more correct to use due to the non-negative nature of the data, an ICA algorithm based on maximum likelihood described in [18] was used as it gave approximately non-negative results. For the derivation of the Bayesian Information Criteria used consult Theorem 6 and Theorem 7 page 101-102. The number of observation in the BIC measures was defined as the number of time points in the data. Furthermore, BIC was normalized by the number of observations. Although only the ALSPARAFAC corresponded to a least-square optimization the SR1PARAFAC and ICAPARAFAC algorithms also used the BIC given for a least square solution. This was done since the factors found of SR1PARAFAC and ICAPARAFAC were believed to be close to a pure least square solution.

ALSPARAFAC

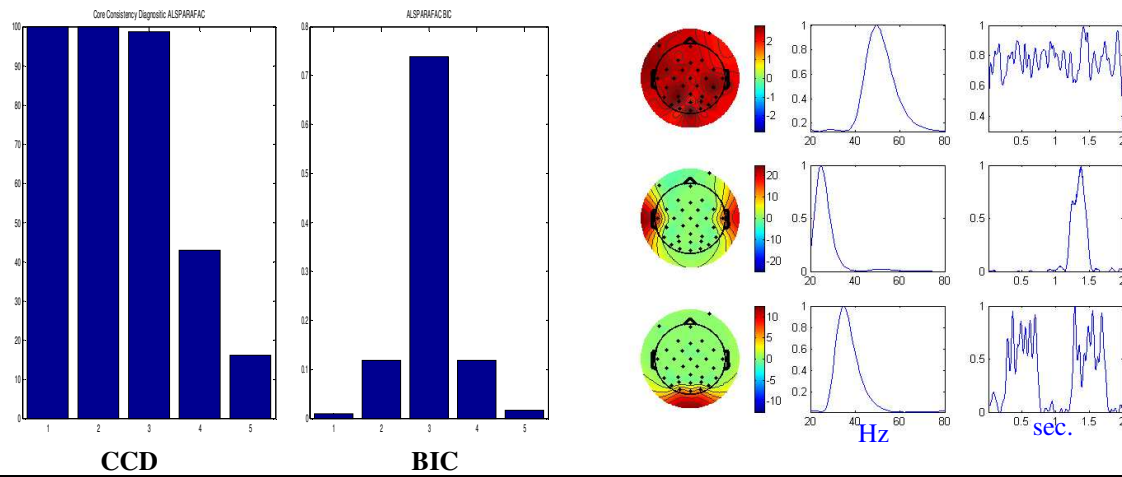


Figure 3.5: To the left; the determination of the number of factors present using ALSPARAFAC, given by the Core Consistency Diagnostic, CCD and BIC. Both the CCD and BIC clearly indicate a three component model. To the right; the factors found when fitting a three component model. Obviously, ALSPARAFAC has positively identified all three factors.

SR1PARAFAC

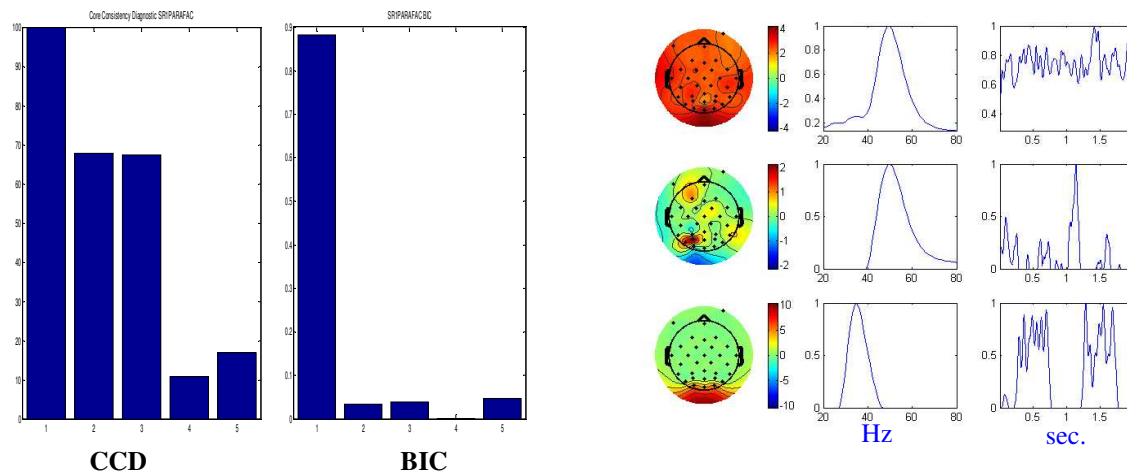


Figure 3.6: To the left; the determination of the number of factors present using the SR1PARAFAC, given by CCD and BIC. The CCD uncertainly indicates one to three components present whereas BIC give strong indication of a one component model. To the right; the factors found when fitting a three component model. The SR1PARAFAC only correctly identifies two of the three factors.

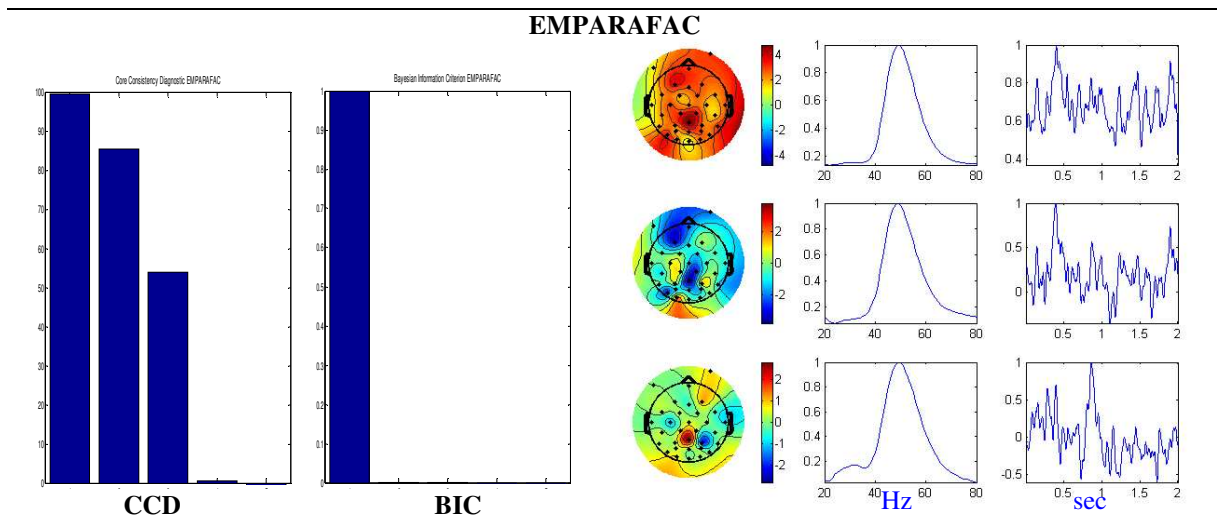


Figure 3.7: To the left; the determination of the number of factors present using EMPARAFAC, given by CCD and BIC. The CCD indicate a model having two factors whereas BIC gives sign of only one factor. To the right; the factors found when fitting a three component model. The EMPARAFAC only identify as indicated by BIC one component - the 50 Hz noise present in all channels.

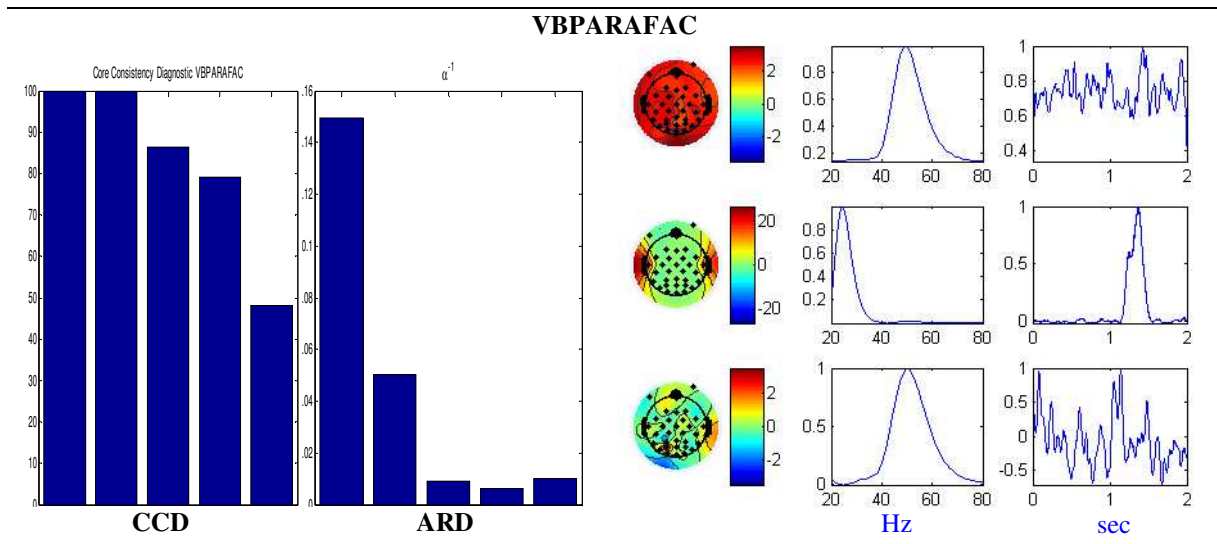


Figure 3.8: To the left; the determination of the number of factors present using VBPARAFAC, given by the CCD and ARD. The CCD is very unclear but indicate that up to four factors are present. The ARD however only reveal that one or two factors are present. To the right; the factors found when fitting a three component model. The VBPARAFAC correctly identifies as indicated by the ARD two components - the 50 Hz noise present in all channels and the 25 Hz ear activity. It is however unable to find the occipital activity².

² The VBPARAFAC ran for 10,000 iterations as suggested by [27], priors were set to be non-informative. However, for both VBPARAFAC and EMPARAFAC it was difficult to determine whether the algorithms had converged.

ICAPARAFAC

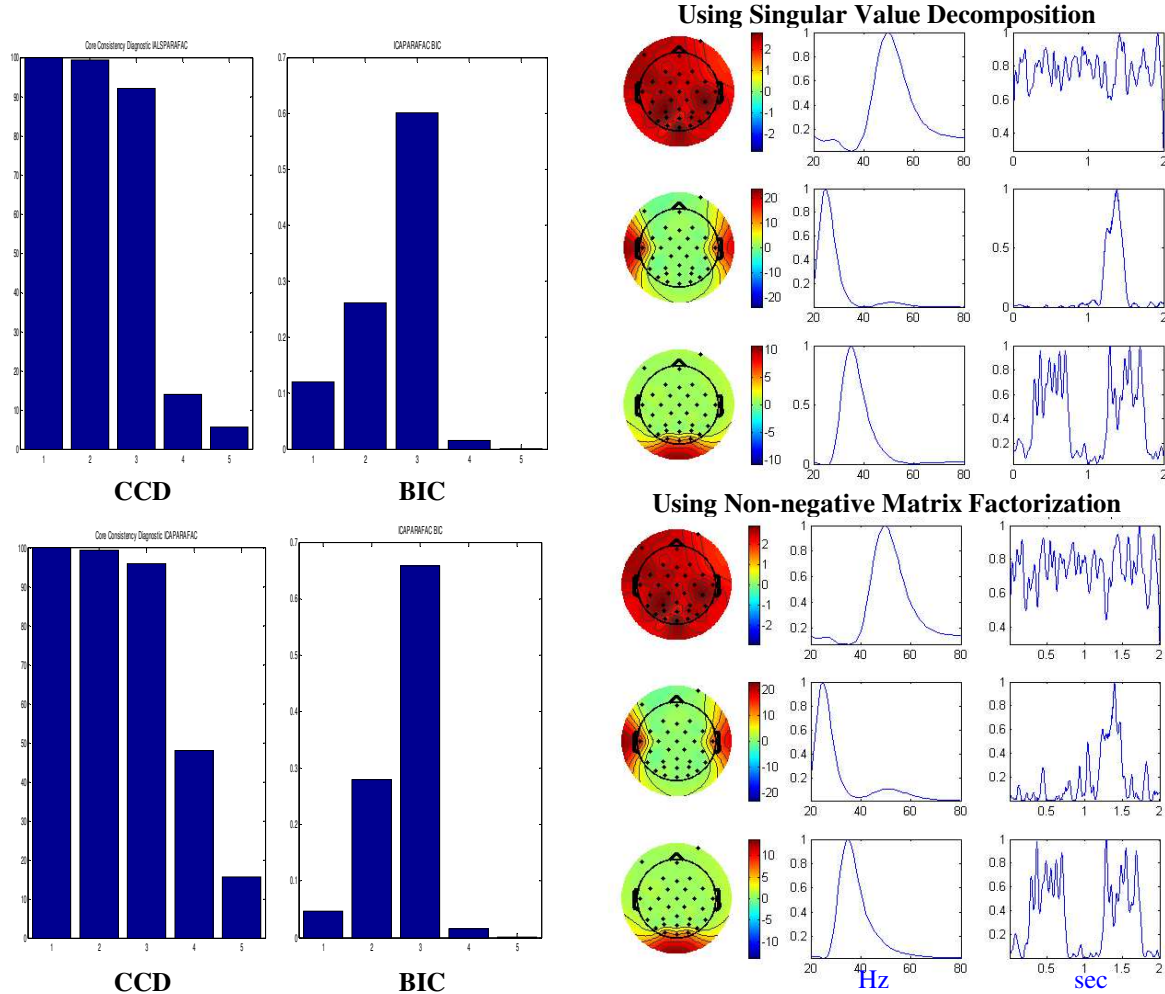


Figure 3.9: To the left; the determination of the number of factors present using ICAPARAFAC, given by CCD and BIC for an ICAPARAFAC algorithm using SVD and one implemented with NMF. The CCD and BIC of both algorithms clearly indicate a three factor model. Both methods are also able to correctly identify the three factors. The frequencies and temporal signatures are however slightly different from each other.

As seen on Figure 3.5 the Core Consistency Diagnostic clearly indicates that three factors are present in the ALSPARAFAC, this is confirmed by the Bayesian Information Criterion. The method is also able to correctly identify all three factors. A much weaker indication of a three component model is given by the CCD for the SR1PARAFAC, see Figure 3.6. The BIC for SR1PARAFAC indicate however that only one factor is present. The SR1PARAFAC is able to correctly identify two components, the 50 Hz activity in all channels and the 25 Hz ear activity, but the 35 Hz occipital activity is lost. The EMPARAFAC algorithm as seen on Figure 3.7 is only able to identify the 50 Hz activity in all channels, from BIC it is also seen that only one factor is indicated to be present in the data. From the automatic relevance determination (ARD) of the VBPARAFAC on Figure 3.8, it is seen that one to two factors are found to be present in the data whereas the CCD is very unclear but indicate that up to four factors are present. The

VBPARAFAC method correctly finds the 50 Hz and 25 Hz activity. Finally, the ICAPARAFAC based on SVD and NMF both clearly indicate from the CCD and BIC as seen on Figure 3.9 that three factors are present in the data. Both methods also correctly identify all three factors. From the results of the ICAPARAFAC algorithm using SVD or NMF didn't change the CCD or BIC. However, the temporal signatures as well as the frequency signatures were slightly altered. Notice how the SVD solution is very close to the ALSPARAFAC solution.

From the simulated data only the ALSPARAFAC algorithm and ICAPARAFAC algorithm successfully identified all the factors. For these two methods the CCD and BIC both worked well, as they strongly indicated three factors were present. The two algorithms were then compared in their ability at different noise level to identify the 25 Hz activity at the ears and 35 Hz activity at the occipital region. For each level of noise fifty ALSPARAFAC and ICAPARAFAC models were fitted to the data. The ICAPARAFAC was based on the Non-negative Matrix Factorization. Both algorithms were evaluated by how much their found factors correlated with the true underlying factors. The correlation was calculated as the average correlation taken over each of the three factor-components, i.e. as the average correlation of the topographic, frequency and temporal signatures between the real and found factors.

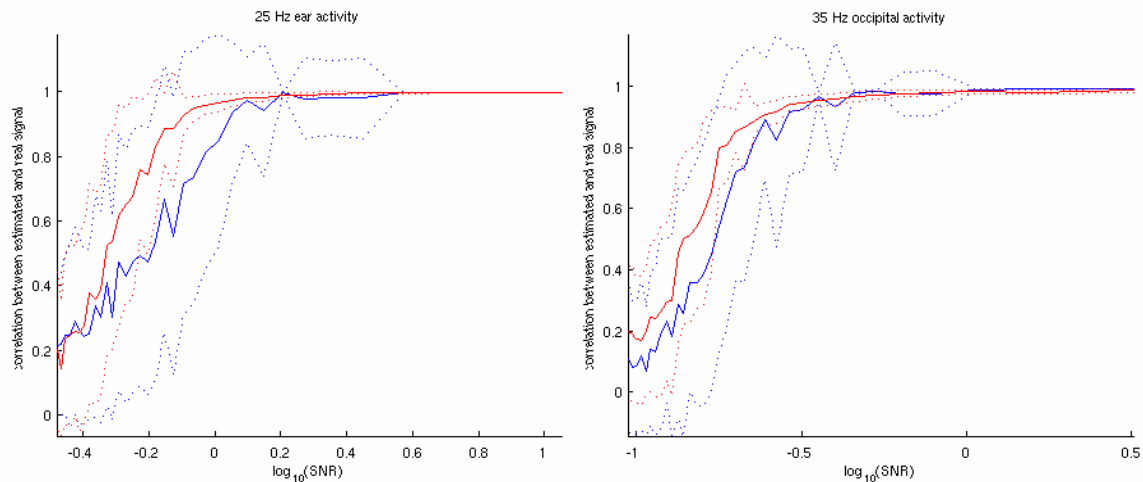


Figure 3.10: The correlation between the true and estimated factors for different signal to noise ratios (SNR). Blue corresponds to ALSPARAFAC, red to ICAPARAFAC. Dashed lines correspond to one standard deviation from the solid lines. Clearly, the ICAPARAFAC is better at finding the true components and more stable than the ALSPARAFAC method as the SNR drops.

From Figure 3.10 it is seen that the ICAPARAFAC algorithm is better at estimating both the 25 Hz ear and 35 Hz occipital activity as the signal to noise ratio drops. Both methods have more problems finding the ear activity than the occipital activity when the signal to noise ratio decreases. Whereas both methods correctly identified the occipital activity down to a $\text{SNR}=10^{-0.5}=0.32$, already around a $\text{SNR}=10^{0.1}=1.26$ the ALSPARAFAC methods have problems finding the ear activity. This stems from the fact that the occipital activity is present longer and in more channels than the ear activity making it easier to detect. Furthermore, ICAPARAFAC is more stable than ALSPARAFAC as the standard deviation of ICAPARAFAC is smaller than that of ALSPARAFAC and ALSPARAFAC

begins to be unstable earlier around a $SNR=10^{0.5}=3.16$ for the ear activity and at $SNR=10^0=1$ for the occipital activity.

Finally, the ICAPARAFAC method was compared to the ALSPARAFAC on the chemometric data set ‘‘Claus’’ described in [41]. The analysis is shown in Appendix D: ICA- and ALSPARAFAC on Chemometric Data. Also on this dataset ICAPARAFAC performed well.

3.2 Real data

The real data is generated from an experiment by Herrmann et al. [15] regarding gamma oscillations in the visual system. Gamma oscillations have been shown to correlate with perceptual binding, attention, arousal, object recognition and language perception. A mechanism which underlies many of the above mentioned cognitive functions is the match of sensory information with memory contents. Herrmann and colleagues argue that the so-called ‘early’ gamma-band activity, occurring in EEG before 150 ms after stimulus presentation, reflects a match with memory. In addition, they argue that ‘late’ gamma activity, which typically emerges with a latency of more than 200 ms, is a temporal signature of utilization processes such as response selection or context updating [14]. This has led to the Match and Utilization Model, MUM, explained in Figure 3.11.

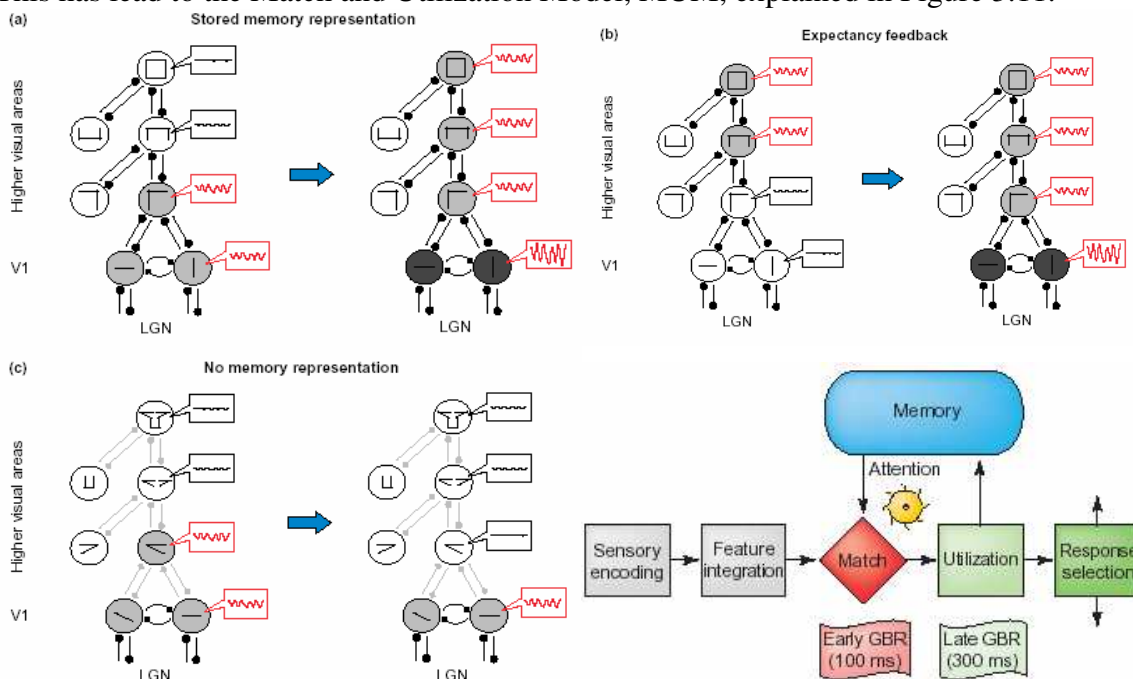


Figure 3.11: Black connections represents memory connection, gray the lack of memory connection. According to Herrmann et al. a stored memory representation will result in an enhanced gamma oscillations and synchrony as the features are matched with the memory content. Furthermore, as revealed in b the expectation of a known visual stimuli can result in enhanced gamma oscillations and synchrony as the neurons expecting the visual feature are closer to threshold. As revealed in c, when no memory representation is present, no enhancement takes place. This concept has been expanded to the Match and Utilization model, MUM, (lower right figure): Sensory coding is integrated into features, and these features are matched with memory contents around 100 ms. At around 300 ms a process denoted utilization takes place. Here updating of memory contents, selection of different behavioural responses and the reallocation of attention is believed to take place. Whereas the match at 100 ms is evoked the utilization at 300 ms is believed to be induced. Figures adapted from [14].

Herrmann et al. find that under visual stimulation a strong increase in evoked oscillations near 40 Hz over posterior areas with a latency of approximately 100 ms and a later increase in induced activity with a latency around 300 ms can be observed [6],[15]. As evoked activity is phase locked to the stimuli, the ITPC will be analyzed. Coherence in the posterior regions is expected to be found. In the following, coherence is defined as the ITPC.

Eleven healthy subjects with mean age 25.7 ± 1.7 years participated in the experiment. All subjects had normal or corrected to normal vision. The experiment was done by Sidse Arnfred at Cognitive Research Unit, Department of Psychiatry, Hvidovre Hospital. The subjects were asked to classify objects as round or edgy by right or left clicking a computer mouse. Some of the objects had a long term memory representation (object) whereas other objects consisted of the same atoms but randomly placed not to make sense (non-objects), see also Figure 3.12. To insure the subjects were naïve to the experiment the task of classifying the objects as edgy or round was given even though no such clear interpretation of the objects was always present.

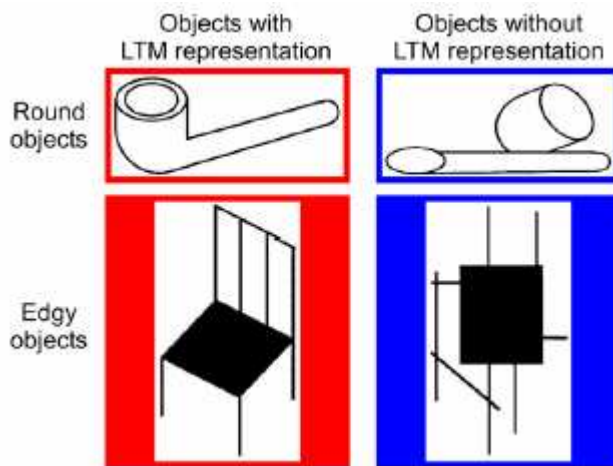


Figure 3.12: Example of stimuli with long term memory representation (object) and without (non-object) , taken from [15].

The subjects were recorded using a BIOSEMI 64 channel active electrode system, see also http://www.biosemi.com/active_electrode.htm. The EEG was referenced to the average of two channels placed at each ear, i.e. channel 65 and 66. Data was sampled at 512 Hz. The epochs were extracted from the data taking -250 to 1000 ms. from stimuli onset. Baseline activity from -250 to -100 ms. was subtracted each epoch. A total of between 102 and 105 epochs were present in both the object and non-object condition for each subject. A complex Morlet wavelet with center frequency 1 and bandwidth parameter 2 was used. Although Herrmann et al. suggest removing epochs having standard deviations more than $50 \mu\text{V}$ [15], we compared this rejection criterion with an extensive rejection analysis of the epochs in EEGLAB using independent component analysis as suggested by Makeig et al. [8]. However, we realized that since we were analyzing the ITPC, the number of epochs used was more important than the quality of

each epoch. Although some epochs were very noisy they still had the correct phase. Since more epochs reduced the noise as averages could be taken over more trials, see also Figure 3.13, we ended up accepting all epochs in the data.

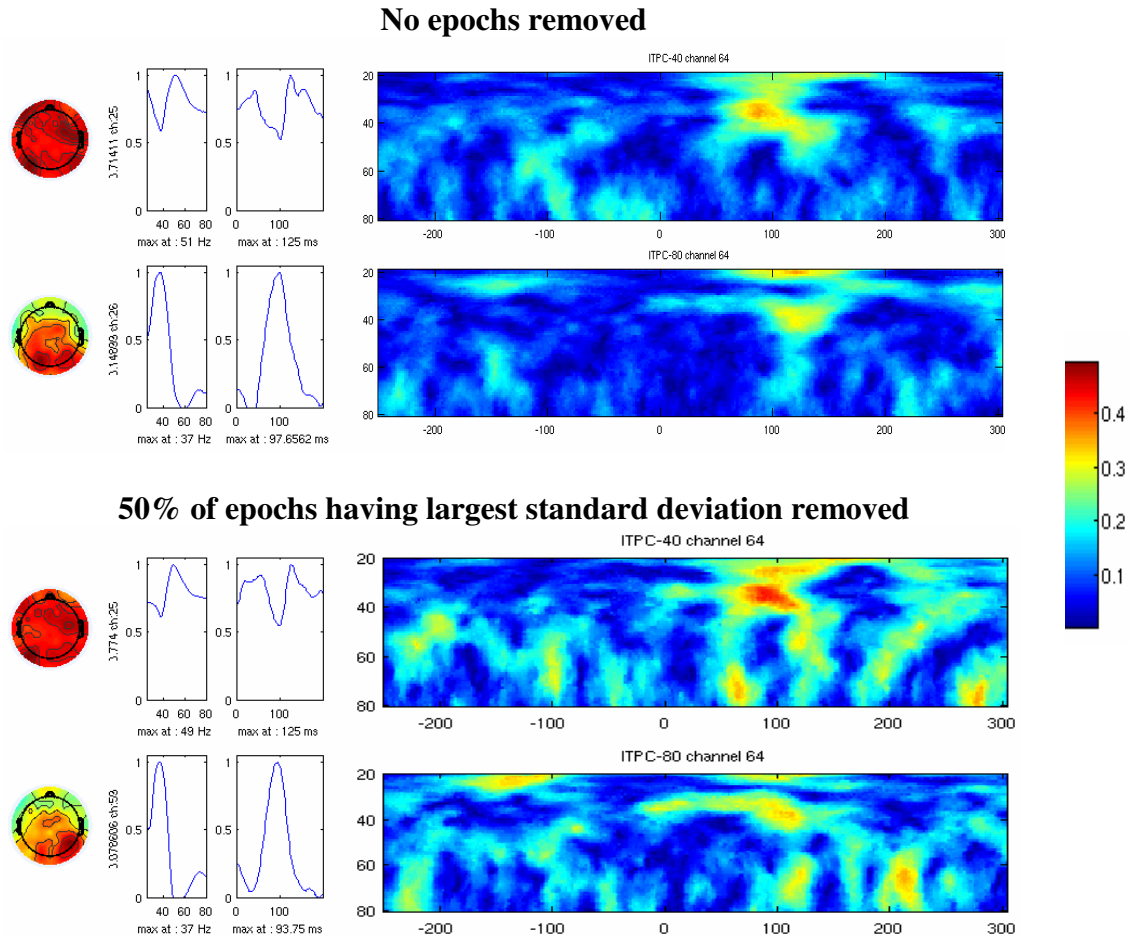


Figure 3.13: Left panel; example of a two component ALSPARAFAC analysis of the ITPC performed on all epochs of a subject and where 50% of the epochs having largest standard deviation within a 200 ms timeframe were removed. Activity at the posterior region is evident from the topographic image with all epochs whereas the removal of 50% of the epochs dramatically removes the coherence in the left occipital region. Right panel; the ITPC found in the object (40) and non-object (80) condition. Clearly the ITPC is less noise full when using all epochs; see top images, compared to removing 50%; bottom images (color scale given to the right).

Finally, the wavelet chosen also to some extent impacted the coherence found as revealed on Figure 3.14.

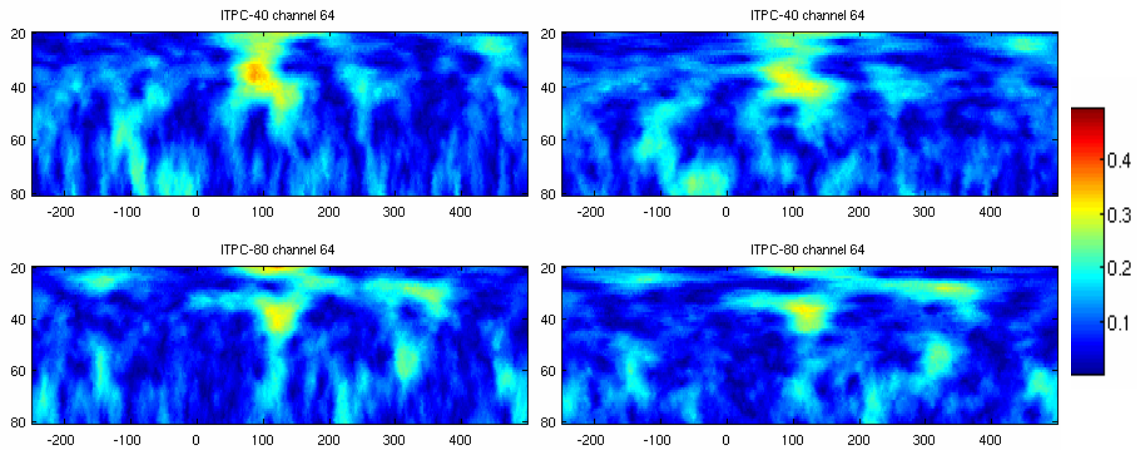


Figure 3.14: Taking the wavelet transform of the data using a complex Morlet wavelet having center frequency 1 and bandwidth parameter 2 (left figure) and bandwidth parameter 4 (right figure). The two wavelet transforms yield slightly different results (x-axis in ms, y-axis in Hz).

Prior to analyzing the data using PARAFAC, the data was analyzed similar to the way Herrmann et al. analyzed their data [15].

Analysis by Herrmann and colleagues

In their analysis, Herrmann and colleagues find the time and frequency corresponding to the coherence peak at channel 64 (equivalent to O2, placed at the center of the right hemispheres occipital lobe). However, as the whole occipital region is affected, we also analyzed the mean of the occipital region corresponding to channel 20-31 and 57-64.

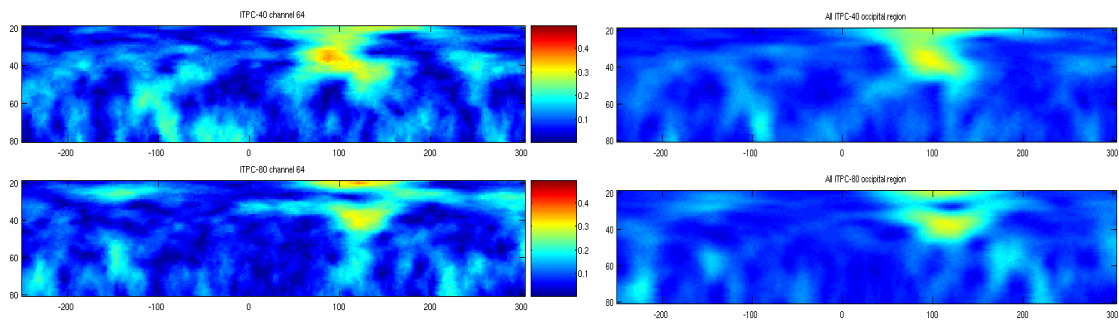


Figure 3.15: Left panel; an example of a subjects ITPC of channel 64 for the object condition top image and non-object condition bottom image. Right panel; same figure, but the average of the whole occipital region. Both panels clearly reveal gamma activity around 100 ms (x-axis in ms, y-axis in Hz).

As revealed on Figure 3.15 there is a strong coherence at around 37 Hz and 100 ms. Furthermore, the coherence seems stronger in the object situation than the non-object. Herrmann et al. find the time and frequency corresponding to each subject's peak in the

object situation and compares this coherence value with the corresponding value at same time-frequency for the non-object condition.

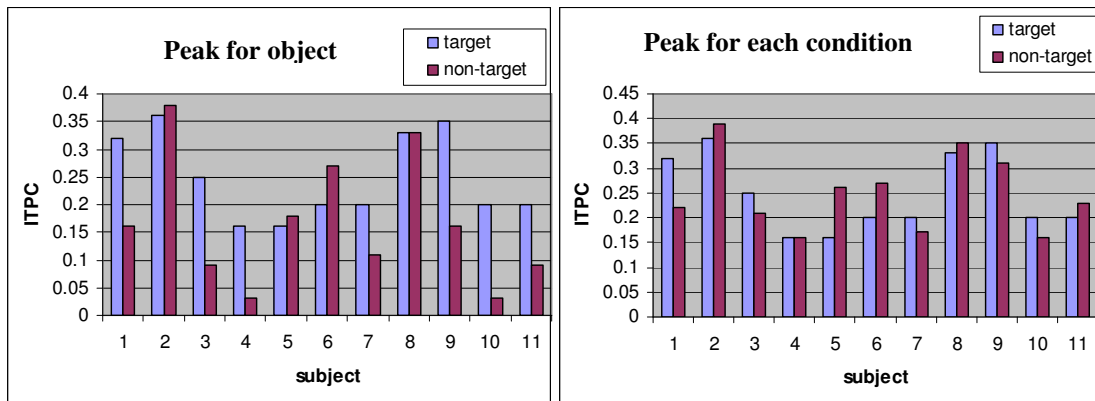


Figure 3.16: The coherence values for object and non object where object peaks and the coherence values for object and non-object where each condition has its peaks. No significant difference is in the two situations found between the conditions (target=object, non-target=non-object).

As seen on Figure 3.16 although the object condition results in higher coherence values in most situations when comparing the coherence at the peak of the object condition, it is not significant as Herrmann et al. find it to be. As using the peak of the object condition favors object we also compared the coherence at the peak of object with that of the non-object condition. Here a difference in degree of coherence seemed very random. Consequently, the finding of Herrmann et al. that the object condition is more coherent than the non-object condition seems very questionable. The same analysis performed on the whole occipital region yielded similar results.

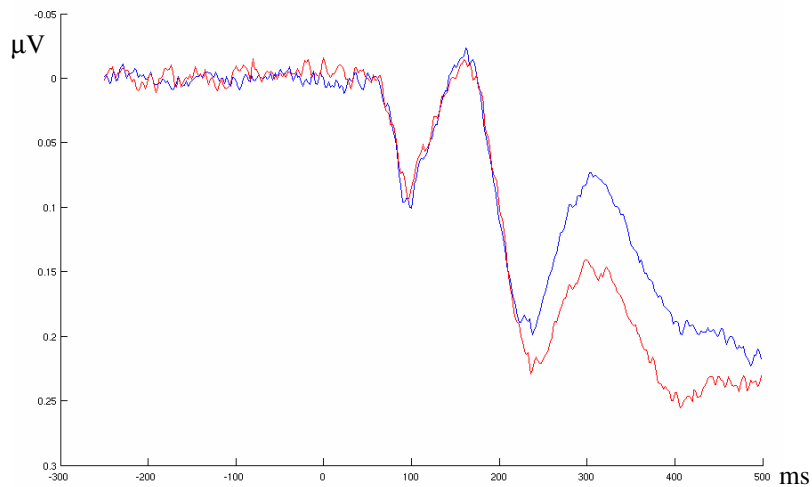


Figure 3.17: The ERP of the grand average of all subjects taken over the whole occipital region, i.e. channel 20 to 31 and 57 to 64, blue is object, red is non-object. 20% of the epoch having largest standard deviation within a 200 ms time-frame was removed. Clearly there is a difference in the ERP of object and non-object from 200-500 ms (Notice; negative is up).

From Figure 3.17 the ERP of the grand average in the occipital region reveals a P100 followed by an N100 after which a P200, N200 and a late P300 appear (around 400 ms). While no difference is present in the ERP from 0-200 ms there seems to be a difference between the ERP of object and non-object from 200-500 ms, this difference hasn't been explained by Herrmann et al. As seen on Figure 3.17 the difference between object and non-object mainly stems from the P200 and N200. In Table 2 page 54 it was explained that the P200 is known to increase with novel stimuli. As non-objects contrary to objects represent a novel stimuli every time this might explain the larger P200 for non-object. Furthermore, N200 is known to increase to task deviant stimuli. As the task was only present to keep the subjects naïve to the experiment, it is difficult to evaluate task effects as they ideally shouldn't be correlated with the condition type. However, as non-objects are probably easier to classify as edgy or round than objects (i.e. the pipe of Figure 3.12 is both edgy and round) this could explain the larger N200 for object.

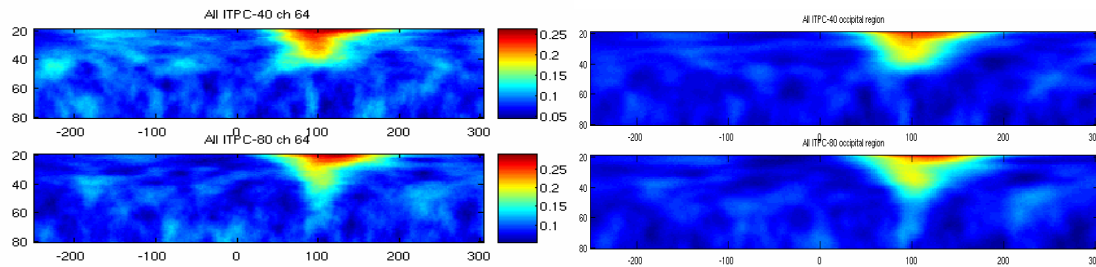


Figure 3.18: Grand average of the ITPC, object and non-object for channel 64 and the average of the whole occipital region; object in top images, non-object in bottom images. Clearly there is a strong coherence between 20-40 Hz around 100 ms. (x-axis in ms, y-axis in Hz).

Figure 3.18 shows the grand average of the ITPC for all 11 subjects at channel 64 and in the whole occipital region. It seems as if slightly more coherence is present in the object situation. Within the gamma band, the grand average at channel 64 peaks at 37 Hz and 107 ms.

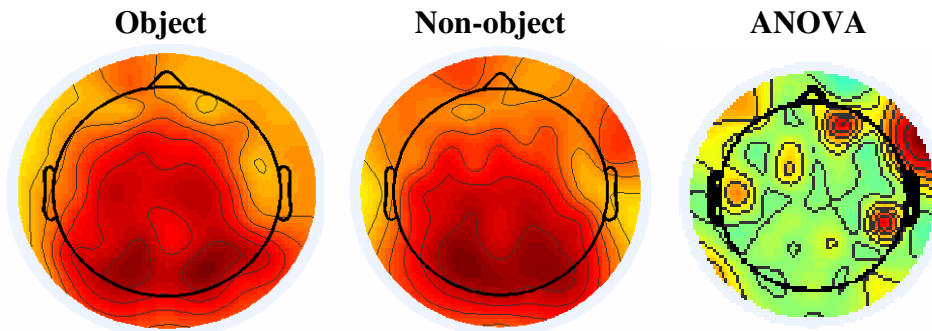


Figure 3.19: The mean coherence in all channels of the 11 subjects at 37 Hz, 107 ms for the object and non object condition (color scale is the same), and the ANOVA of the analysis of difference between object and non-object at this time-frequency point. From the ANOVA no difference between object and non-object is found in the occipital region.

From Figure 3.19 it is seen that the average coherence is more or less the same for object and non-object at the peak of the grand average for object. A test of difference between object and non-object also reveal that no significant difference is present. The largest difference is found to the frontal right where no difference is theoretically justified.

Summary of the analysis by Herrmann and colleagues

From the analysis corresponding to Herrmann and colleagues no significant difference between the object and non-object condition was found. However, coherence was clearly present in the gamma band around 50-150 ms as explained by Herrmann et al. Furthermore, a difference in the ERP between the object and non-object condition seemed to be present from 200-500 ms.

Analysis by PARAFAC

In the following, the analysis, if not otherwise stated, is performed by the ALSPARAFAC algorithm with 'row-wise' non-negativity constraints on all modes. As PARAFAC is a data exploratory tool the analysis was performed without any prior assumptions of what to expect to see from the data. First, an overall 4 way analysis was performed defined by *channel* \times *frequency* \times *time* \times *subject* from 0-200 ms from stimuli onset. A Core Consistency Diagnostic was only possible to access when analyzing three-way arrays as the diagnostic was too memory consuming for MATLAB, even for a computer having 2 GB of RAM. The factors were ordered in accordance to the amount of variation they explained. Each analysis was run several times to assure stable solutions.

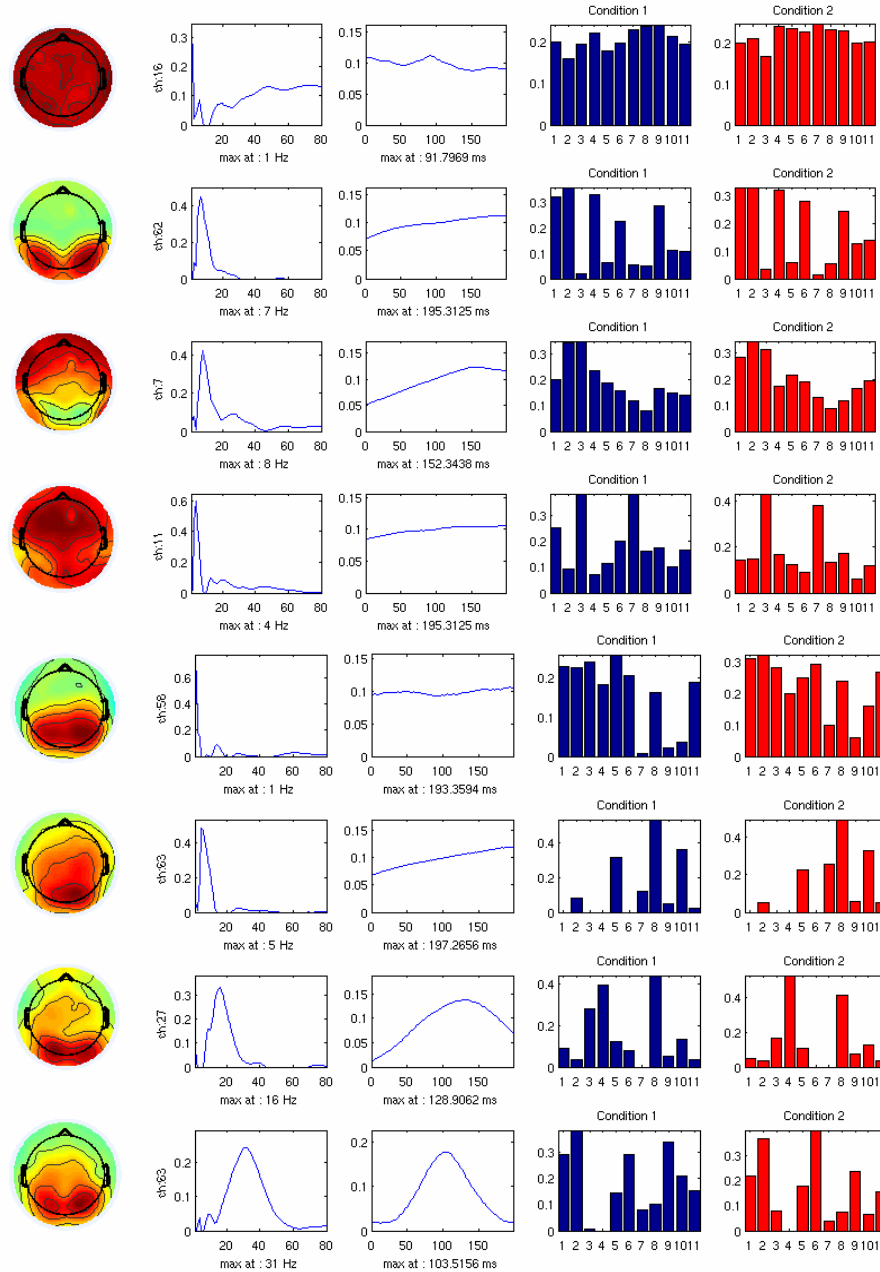


Figure 3.20: An eight component PARAFAC model fitted to the data, blue bars corresponds to object, red to non-object. Factor 2, 5, 6, 7 and 8 all indicate occipital activity. Especially factor 8 pertains to the Gamma activity around 100 ms as described by Herrmann and colleagues.

As seen on Figure 3.20 the first factor models some average activity. The second, fifth and sixth factor correspond to low frequent occipital activity relating to the ERP. For all these factors no significant difference is found between the two conditions. The eighth factor however, reveals a gamma activity in the occipital region around 100 ms corresponding to Herrmann et al.'s findings. The subjects' activities during the two conditions reveal that the last 5 subjects have more gamma activity in this factor during object than the non-object condition. Also, subject three and four seems to almost

completely lack this activity. Furthermore, factor seven reveals some beta activity around 130 ms.

An ANOVA was performed to look for differences between the two conditions for the eleven subjects. This gave an F-test value multi-way array given by $channel \times frequency \times time$ as revealed in Figure 3.21.

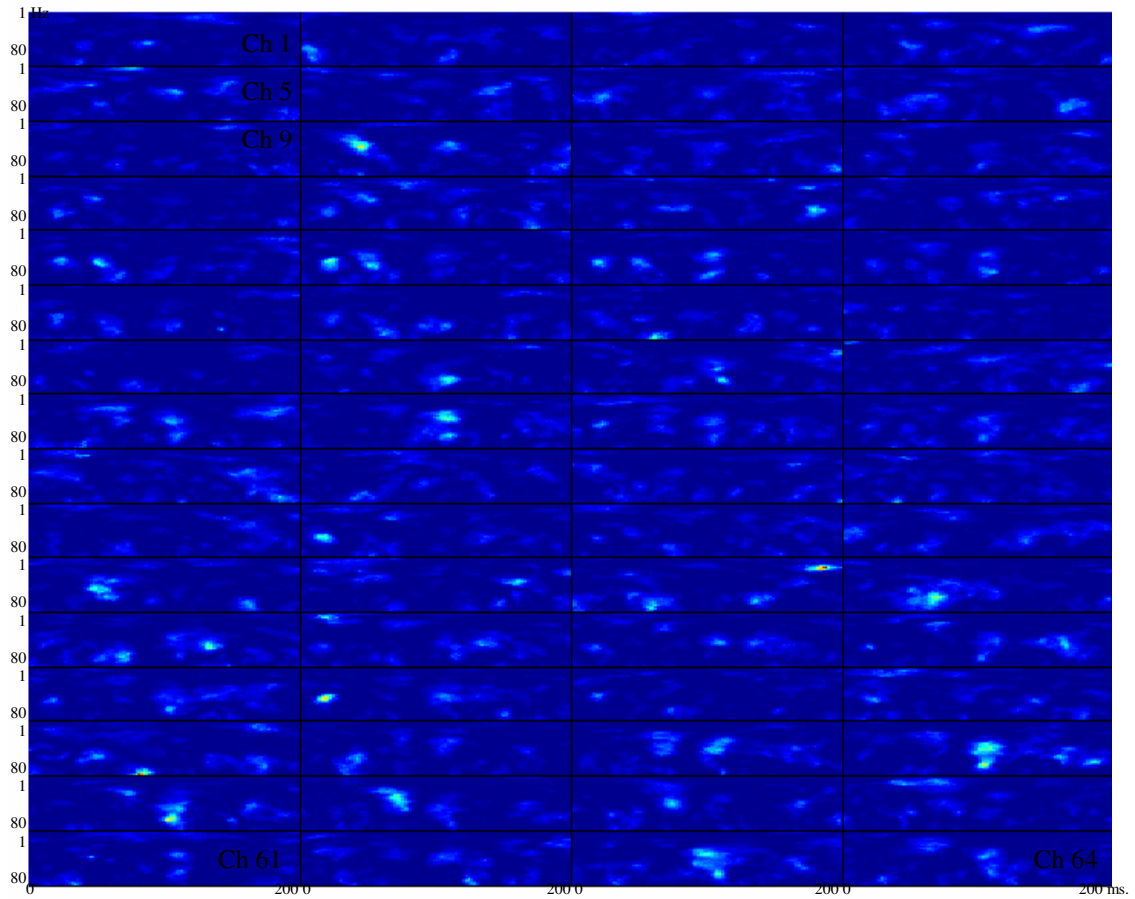
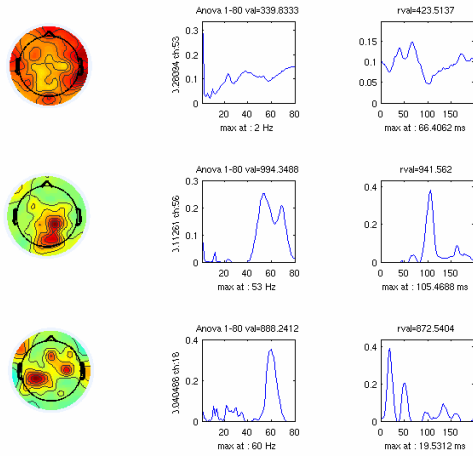


Figure 3.21: ANOVA test of difference between object and non-object in the 11 subjects, shown in a 16×4 array where each array represent a channels F-test value to given frequency-time point. From the F-values in the array it is difficult to grasp where the differences between the two conditions are present.

ALSPARAFAC



ICAPARAFAC

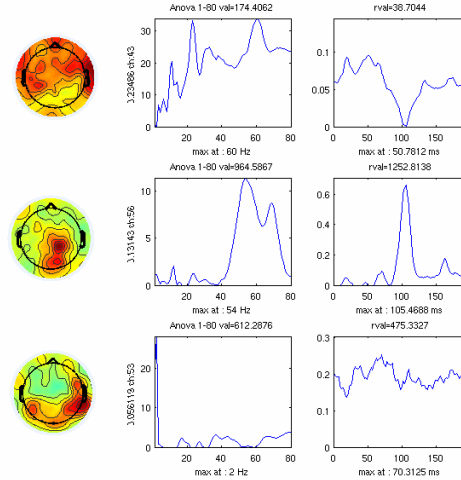


Figure 3.22: A PARAFAC model based on ALSPARAFAC and ICAPARAFAC fitted to the F-test multi-way array. Where first factor models some background activity, the second factor of both methods indicates a difference around 100 ms in the Gamma band in accordance with Herrmann and colleagues findings.

Figure 3.22 shows an ALSPARAFAC and ICAPARAFAC model fitted to the F-test multi-way array. The first factor of both algorithms models some background activity. The second factor shows that the difference between object and non-object primarily is in the occipital region in the gamma band of 30-80 Hz. It is difficult to explain what the last factor of ALSPARAFAC pertains to, but the third factor of the ICAPARAFAC model reveals a 2 Hz difference in the occipital region between the two groups corresponding to a difference in the ERP. As a result, the ANOVA clearly indicate that the difference between the two groups is as Herrmann et al. found in the Gamma band around 100 ms. As a result; the PARAFAC model is capable without any prior knowledge to identify the interesting features of the data.

Analyzing the Gamma band (30-80 Hz) by PARAFAC

To analyze the Gamma range a PARAFAC model was fitted to the data in the frequency range 30-80 Hz. Where the first and second factor of Figure 3.23 models some background activity the third factor shows the occipital gamma activity and the fourth factor reveals a central gamma activity. No systematic difference is found in the factors between the two conditions.

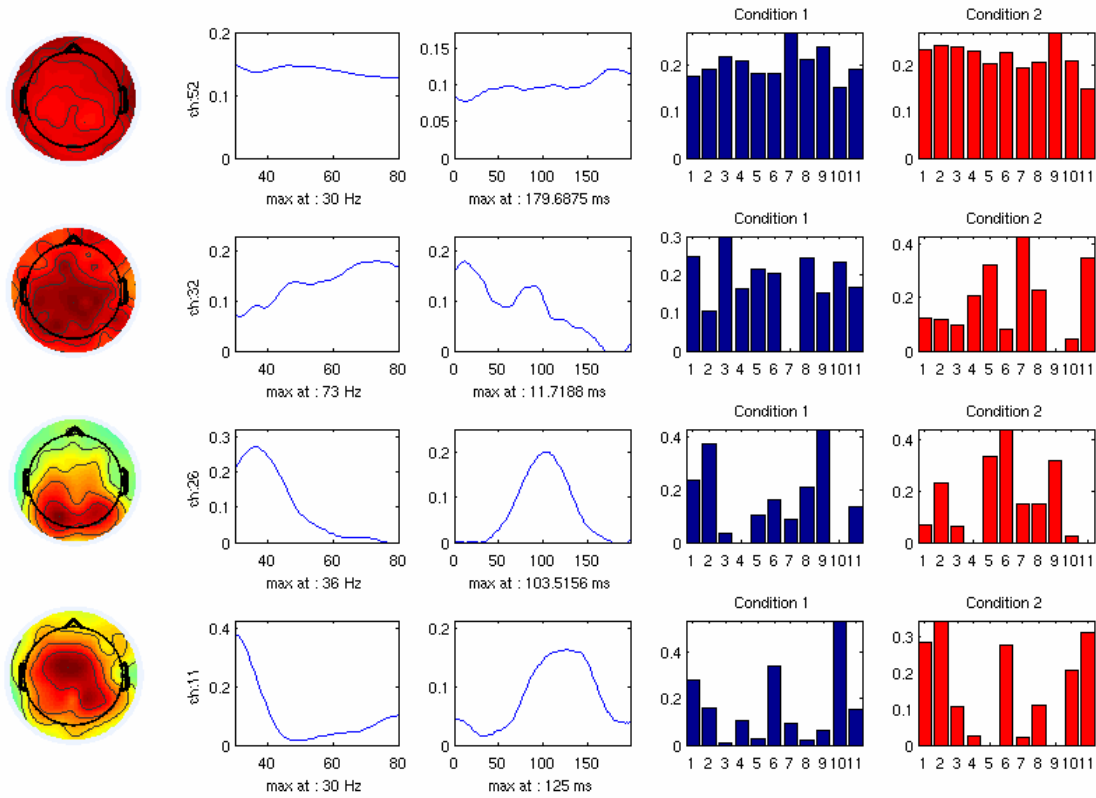


Figure 3.23: A four component PARAFAC model fitted to the data at the frequency range 30-80 Hz. Where first two factors model some average background activity, the third factor clearly reveal an occipital Gamma activity around 36 Hz at 104 ms. Finally, the last factor is more central, delayed and lower frequent.

The condition was also taken into the PARAFAC model yielding the 5-way model given in Figure 3.24. The first factor of this analysis clearly reveals some occipital gamma activity slightly more present in the object (1) than non-object (2) condition. The second factor pertains only to the non-object condition. It models a slightly more frontal, higher frequency activity around 100 ms.

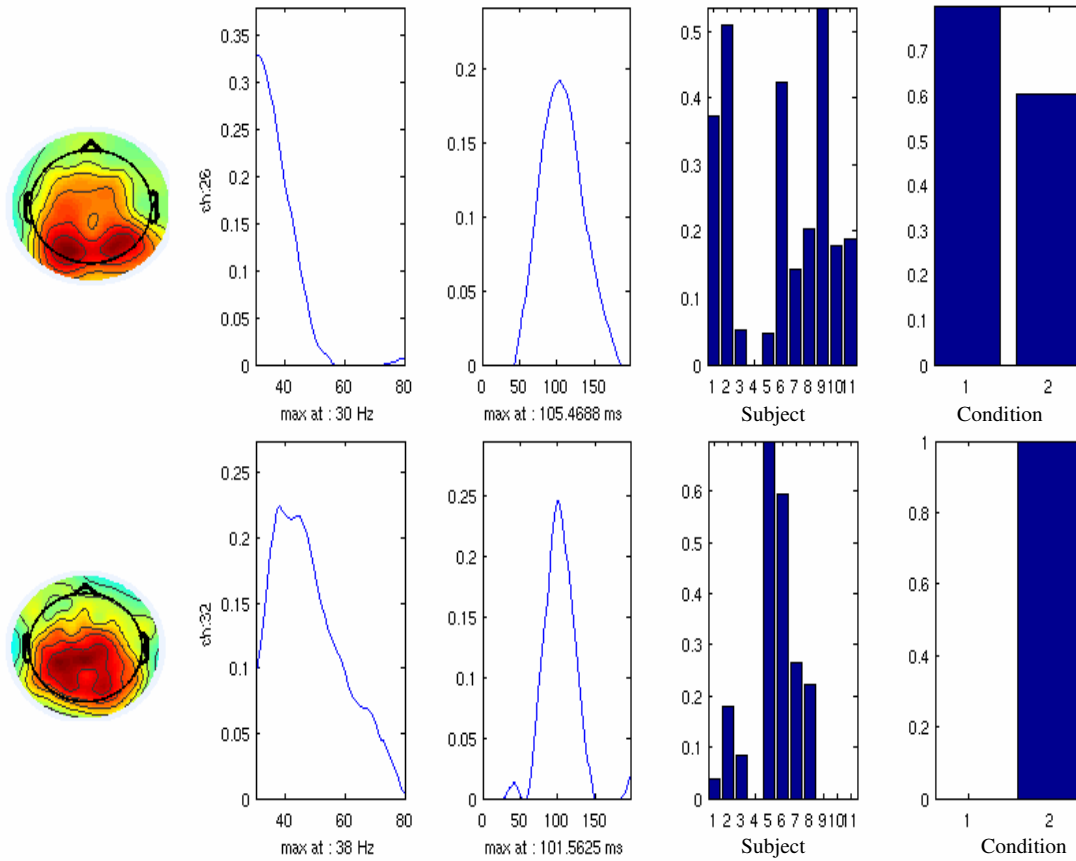


Figure 3.24: A PARAFAC model fitted to the data where condition was taken as an extra modality, 1 is object, 2 is non-object. As only two components could be found due to the limitation of only two conditions baseline activity was subtracted before fitting the PARAFAC. The first factor clearly represents the occipital gamma activity around 100 ms. This factor is mostly present in the object condition but weak in subject 3, 4 and 5. The second factor is higher frequent, slightly more central and pertains only to the non-object condition. The two factors indicate that the object condition is lower frequent whereas the non-object is slightly higher frequent and more central.

Furthermore, a PARAFAC model was fitted to each condition.

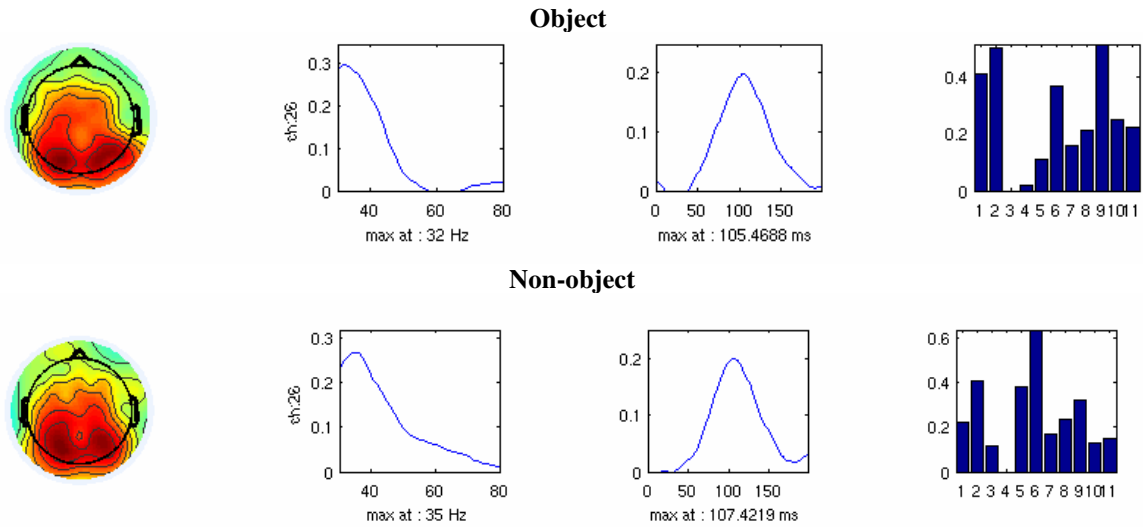


Figure 3.25: Top panel; a PARAFAC model fitted to the object condition. Bottom panel; a PARAFAC model fitted to the non-object condition. Baseline activity subtracted. Again it is revealed that both conditions have clear gamma activity around 100 ms. However, subject 3 and 4 seem to lack the activity in both conditions. Comparing the object with the non-object condition it is seen that the non-object is slightly higher frequent and delayed.

As seen on Figure 3.25 both object and non-object have clear gamma activity around 100 ms in the occipital region. However, the object condition peaks at 32 Hz, 105 ms whereas the non-object peaks at 35 Hz 107 ms. In both conditions subject 3 and 4 have practically no gamma activity in the occipital region.

In addition, a PARAFAC model was fitted to the ANOVA of the gamma band.

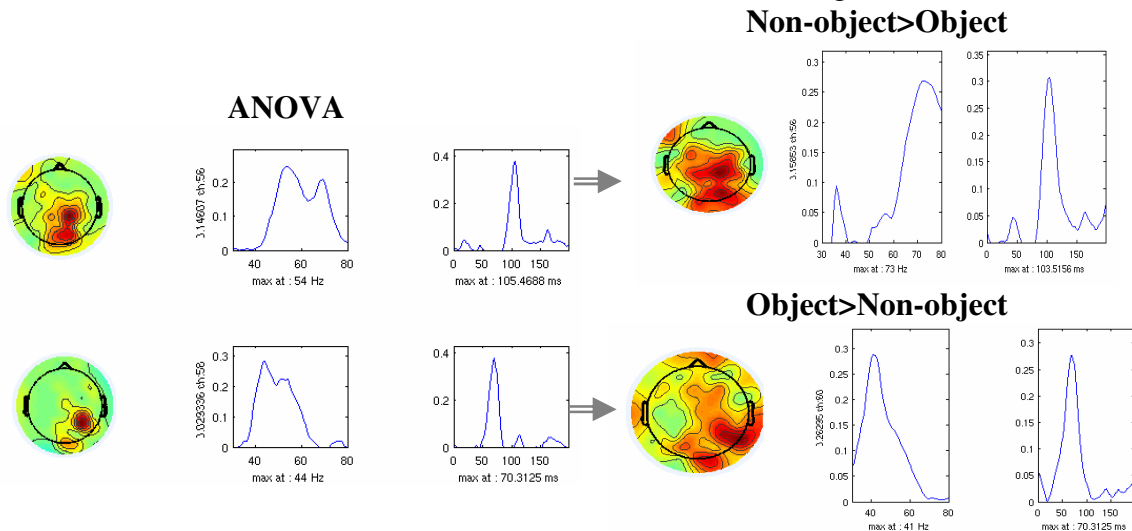


Figure 3.26: Left figure; a PARAFAC based on the F-test value of the gamma band. Top, right figure; a PARAFAC model fitted to regions where non-object is more coherent than object. Bottom, right; a PARAFAC model fitted to regions where object is more coherent than non-object. As seen from the first factor of the ANOVA this factor pertains to the situation where non-target on average is more coherent than target whereas the second factor of the ANOVA corresponds to a situation where object on average is more coherent than non-object. Consequently, object is more coherent early and at lower frequencies whereas non-object is more coherent later and at higher frequencies (baseline activity removed from the data).

From the ANOVA of Figure 3.26 the second factor corresponds to the second factor found in Figure 3.22. Furthermore, the third factor found in the ANOVA of Figure 3.26 also reveals the presence of an earlier and less high frequent difference between the two groups. Analyzing when object is larger than non-object and when non-object is larger than the object condition, it is seen that the first factor of the ANOVA corresponds to the factor where non-object is more coherent than the object condition, whereas the second factor of the ANOVA matches the situation where object is more coherent than the non-object condition. Consequently, the difference between the object and non-object condition is mainly due to the fact that object is coherent earlier and at lower frequencies than non-object.

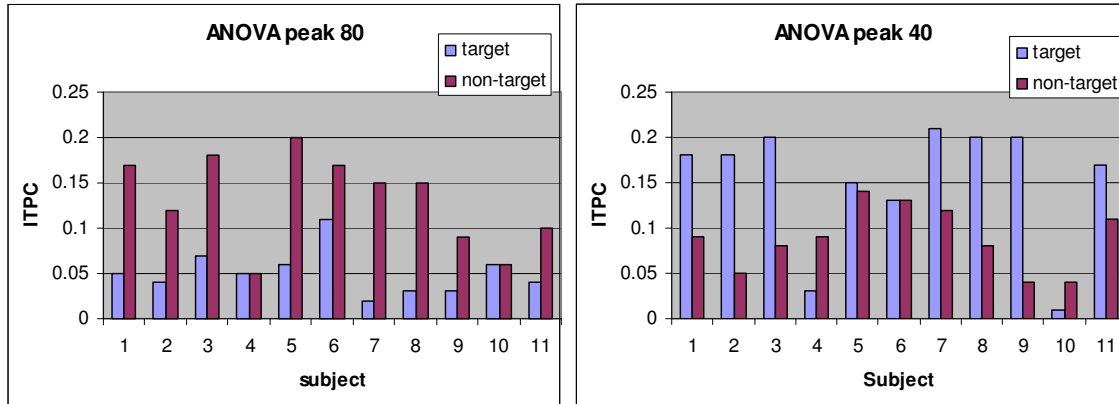


Figure 3.27: The coherence value at the peak of the ANOVA of factor 1 denoted peak 80 and of factor 2 denoted peak 40 found in Figure 3.26. As seen on Figure 3.26 the first factor of the ANOVA corresponds to a situation where non-object is more coherent than object (target) whereas the second factor of the ANOVA pertains to a situation where object is more coherent than non-object.

In Figure 3.27 the same pattern reveals itself. The first factor of the ANOVA in Figure 3.26 corresponds to the situation where non-object is more coherent than object whereas the second factor corresponds to the situation where object is more coherent than the non-object condition.

The PARAFAC model was also fitted to the ITPC of each subject given by the multi-way array $channel \times frequency \times time$.

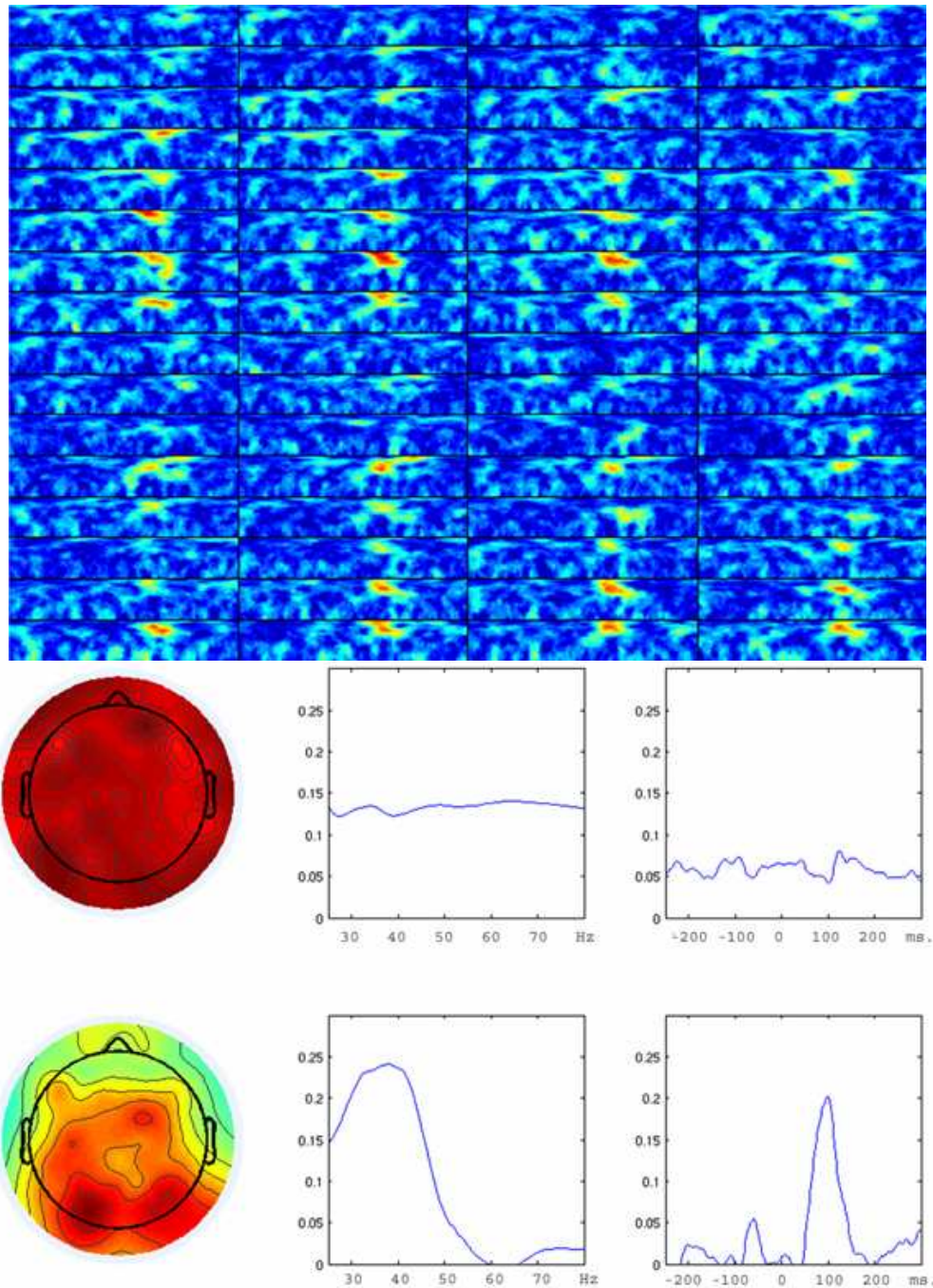


Figure 3.28:Top panel; the ITPC multi-way array given by $channel \times frequency \times time$ of a subject shown in a 16×4 array of channels where x-axis corresponds to time from -250-300 ms and y-axis frequency from 20-80 Hz. Bottom panel; a PARAFAC model fitted to this ITPC. Where the first factor shows some background activity the second factor clearly reveals the occipital Gamma activity around 100 ms.

For each subject, the gamma peak in the occipital region at 50-150 ms was identified in time and frequency by the factor corresponding to the second factor in the PARAFAC

decomposition revealed on Figure 3.28. Notice that the first factor corresponds to some baseline activity.

Having identified the peak of each subject, the coherence at each subject's frequency-time point was found for all channels. Finally, the mean of these topographic maps were calculated as revealed on Figure 3.29.

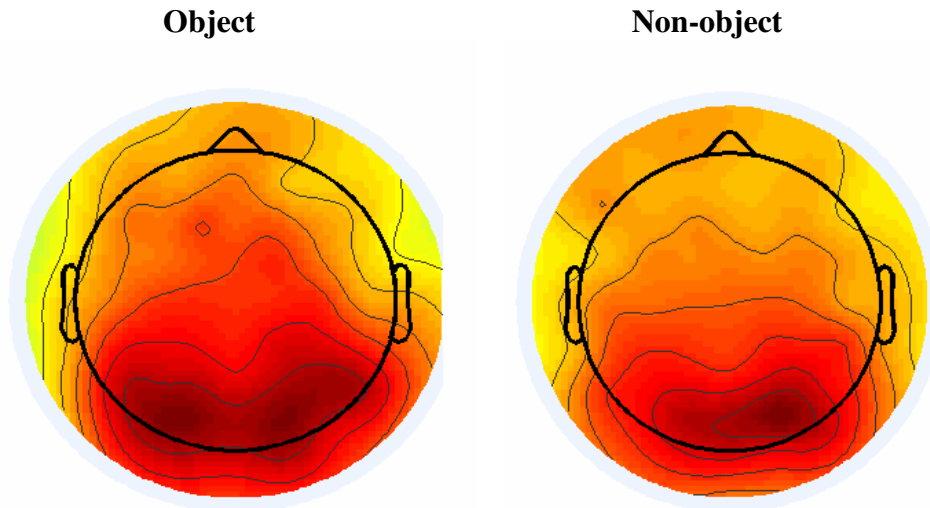


Figure 3.29: The mean coherence for all subjects at their gamma-peak for object and non-object. Coherence seems to be present in a larger region for the object condition than the non-object condition.

Figure 3.29 indicate that the coherence at the peak is high at a much larger region for object than for non-object. As channel 64 which was the basis of Herrmann et al.'s analysis lies right at the peak of both object and the non-object in Figure 3.29 this might be why the findings of difference between the two conditions in Herrmann and colleagues' analysis was poor. Had Herrmann and colleagues' analysis been based on a channel in the left hemisphere, the difference in coherence between object and non-object might have been stronger.

Finally, the PARAFAC model was used to analyze the ERP. This has been done previously by Field et al [10]. Field and colleagues found that a problem of degeneracy arose when fitting the PARAFAC model to the ERP. As a solution they proposed introducing an orthogonality constraint on the dimension representing the temporal development of the ERP. The orthogonality constraint will here be compared to imposing non-negativity as we suggest. The non-negativity can simply be assured by adding a positive constant to the ERP. Prior to analyzing the ERP, 20 % of the epochs having largest standard deviation within a 200 ms time window were removed to get rid of eye and muscle artifacts.

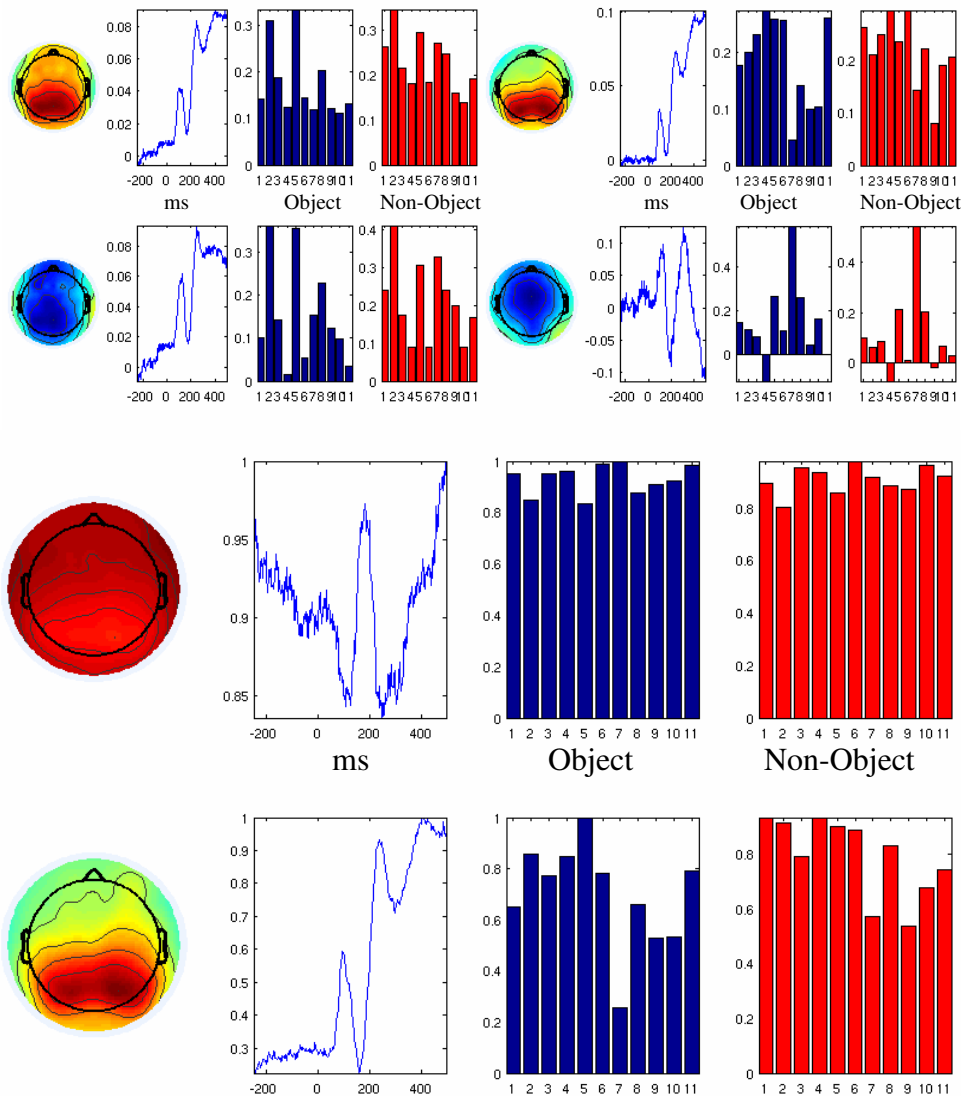


Figure 3.30: Top left panel; analyzing the ERP unconstrained. Top right panel; imposing an orthogonality constraint on the ERP. Bottom panel; imposing non-negativity by addition of a constant. Blue bars correspond to the object condition, red bars to non-object, notice; positive is here upward on the ERP. Neither the unconstrained nor the orthogonality constrained PARAFAC models are able to find the true ERP. This is however found for the non-negativity constrained model where the ERP correctly is split into a frontal and an occipital part.

As seen on Figure 3.30 the unconstrained solution yields highly degenerate factors. The ERP of the second factor is almost identical to the ERP of the first but with opposite sign as revealed in the topographic maps. Imposing the orthogonality constraint insures no degeneracy in the ERP. However, a few subjects have negative coefficients and the justification for the two ERP's to be orthogonal in reality is very questionable. Imposing non-negativity however yields excellent results. The non-negative PARAFAC algorithm has split the ERP into two easy interpretable components. The first component models a mostly frontal ERP whereas the second component beautifully models the ERP of the occipital region. This occipital ERP seems to be more present in the non-object than the

object condition. This corresponds to the findings of Figure 3.17 where the P200 was stronger for the non-object condition.

A PARAFAC model was then fitted to the ERP of each condition.

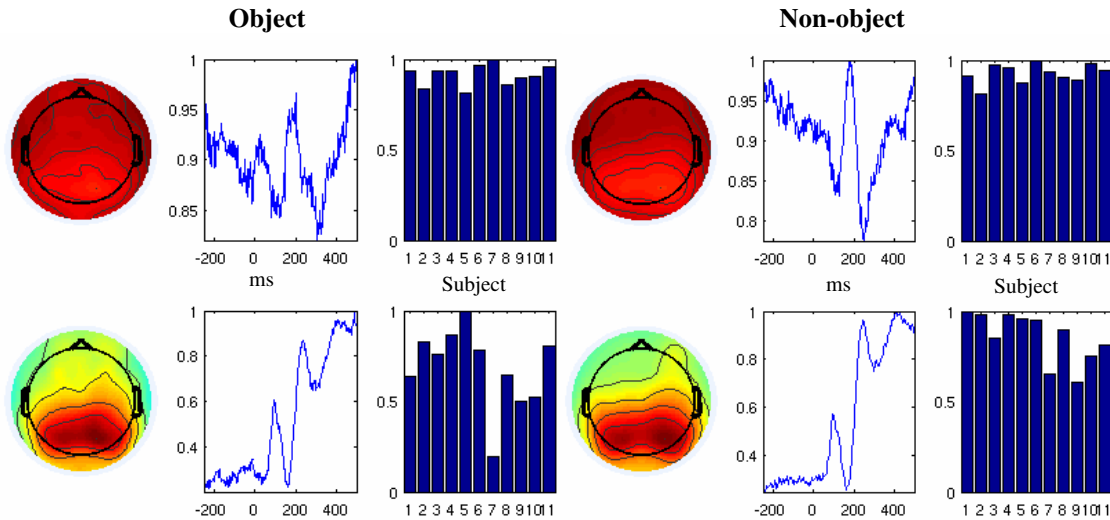


Figure 3.31: The ERP for object, left figure, and non-object, right figure. Where the occipital ERP is very similar in the two conditions the frontal ERP has a stronger P200 in the non-object condition than the object condition.

From the first factor of Figure 3.31 it is seen that the P200 of non-object is much stronger than for object frontally as also found to be the case in the occipital region as revealed in Figure 3.17. Again, this is due to the fact that the non-object represented novel stimuli.

Summary of the PARAFAC analysis

From the PARAFAC analysis of the 11 subjects it was seen that the main activity was in the occipital region corresponding to an experiment having to do with visual stimuli. Furthermore, the difference between the two conditions was mostly present in the gamma band around 100 ms. In addition, this difference was primarily due to a delayed coherent signal at higher frequencies for non-object than the object condition. At the peak of each of the two conditions it seemed as if coherence was present in a larger region in the object situation than the non-object situation. Finally, imposing non-negativity to the ERP by addition of a constant made the PARAFAC model able to correctly split the ERP into an occipital and a frontal factor. This was not possible for an unconstrained model while imposing orthogonality as previously done didn't yield results that were as satisfying. From the ERP's it was found that non-object also had a stronger P200 frontal.

4 Discussion

“It is a good morning exercise for a research scientist to discard a pet hypothesis every day before breakfast. It keeps him young.”

Konrad Lorenz

PARAFAC and simulated data

The test of the PARAFAC algorithms on simulated EEG/ERP data revealed that both ALSPARAFAC and the ICAPARAFAC effectively found the right factors in the data. This was however not possible from the simple 2-D analysis using an ICA-algorithm on the raw data. Consequently, PARAFAC seemed effective in the analysis of the frequency changes in the EEG/ERP when using these two algorithms.

The Core Consistency Diagnostic and Bayesian Information Criterion both proved effective in accessing the correct number of factors in the data for the models that performed the best, i.e. ICAPARAFAC and ALSPARAFAC.

ICAPARAFAC performed better than ALSPARAFAC at estimating the true factors as the signal to noise ratio dropped. As the ICAPARAFAC algorithm also was faster than the ALSPARAFAC algorithm it has great potentials. In Appendix D: ICA- and ALSPARAFAC on Chemometric Data it was revealed that on chemometric data the ICAPARAFAC algorithm also performs well. Therefore, the new ICAPARAFAC algorithm seems promising in a wide range of fields where Combined Independence can be assumed in the data.

Herrmann and colleagues analysis vs. the PARAFAC analysis

Unfortunately, only 11 subjects were analyzed giving quite an uncertain picture to base a conclusive comparison of Herrmann et al.’s findings with the findings of PARAFAC. Furthermore, 2 of the 11 subjects lacked completely clear coherent gamma activity in the occipital region. However, the fact that Herrmann et al. by the basis of the peak of the object condition conclude that object is more coherent than the non-object seems wrong as it favors the object condition. Furthermore, they only examined channel 64 instead of taking the whole occipital region into considerations which is questionable – it might be that channel 64 was chosen for the only reason that it was the most significant. The PARAFAC analysis indicates that the difference found is mostly due to the fact that the non-object peak later and at higher frequencies than the object condition, rather than having to do with non-object in general being less coherent. Furthermore, the analysis based on the individually found peaks in the occipital region for object and non-object revealed that object seemed coherent in a larger region than the non-object.

The PARAFAC analysis was capable of integrating the information present in all channels into simple interpretable components. This made the analysis of Herrmann and colleagues’ paradigm much more complete by PARAFAC than by their own proposed analysis. In addition, it was seen that PARAFAC was able to work in several ways

analyzing both the averaged ERP and the wavelet transformed data over several different modalities taking into account subject and condition variability. Finally, when analyzing the ERP it was revealed that imposing non-negativity worked much better than forcing orthogonality.

The paradigm Herrmann et al. use isn't strong. Although the examples shown from the paradigm in Figure 3.12 are definitely recognizable as a pipe and chair whereas no such interpretation is present in the non-object situation, the non-object almost looks like a broken pipe. Furthermore, there are recognizable objects in the non-object condition such as a square, a cylinder and a bowl which in itself might have long term memory representations. This weakens the difference between object and non-object - weakening the whole paradigm.

5 Conclusion

“The important thing is not to stop questioning. “

Albert Einstein

It was shown that PARAFAC is an effective data exploratory tool in the analysis of the Inter Trial Phase Coherence of the ERP. Furthermore, it was revealed that imposing non-negativity when decomposing the ERP by PARAFAC seemed to solve the problem of degeneracy completely. Here PARAFAC was only used to analyze the ERP and ITPC. However, there is no reason why PARAFAC shouldn't also be an effective tool in the analysis of the ERSP and ERPCOH. These measures were not analyzed as the main interest was to access the Inter Trial Phase Coherence of the data. Furthermore, both measures are very susceptible to noise requiring an extensive preprocessing of the data. However, future work will focus on PARAFAC's ability to analyze these measures.

In the analysis of Herrmann et al.'s gamma band coherence it was found that the non-object was slightly more delayed and higher frequent than the object condition. New experiments are presently conducted to confirm these findings. Future work will also be done to use PARAFAC in analyzing ERP data from other experiments having multiple conditions including different forms of sensory stimuli. As PARAFAC is capable of analyzing complicated multi-modal data it might likely shed new light on these data.

The ICAPARAFAC algorithm also showed promising result, being a great alternative to the popular ALSPARAFAC algorithm when combined independence can be assumed. Work lies ahead in improving the underlying Independent Component Analysis algorithm both to deal with non-negativity as well as being optimized to find the correct components in the ERP.

The PARAFAC model analyzed the wavelet transformed data. As the wavelet transform corresponds to a convolution in which random noise ideally becomes a constant factor at all frequencies analyzing the data in the frequency domain is in itself an efficient way to handle noise full data. Furthermore, the PARAFAC algorithms were able to separate systematic oscillatory noise such as 50 Hz noise from electric devices into a designated component. As it is possible to reconstruct the signal from the wavelet coefficients [32], the EEG/ERP corresponding to each factor of the PARAFAC can also be reconstructed. Thus, the PARAFAC analysis might also work well in reducing systematic noise. Furthermore, knowing the signatures of the noise from a training set can be used to find the noise-signature in one of the modalities from a test set by keeping the components of all other modalities found from the training set fixed on the test set.

The PARAFAC analysis of the coherence relied on the wavelet analysis accessing the correct temporal frequency information of the data. In this thesis the wavelet suggested by Herrmann et al. [14],[15] was used. However, Figure 3.14 showed that the choice of wavelet had an impact on the coherence found. Consequently, the influence of the

wavelet's form used in the analysis of the ERP should also in future work be explored. Maybe even waveforms corresponding to more physiologic burst of neuronal activity could be developed and analyzed.

Wavelet analyzed data is known to be overcomplete, i.e. having more data than information present in the data. Methods such as various matching pursuit algorithms have tried to resolve this problem, reducing the wavelet data to an information level corresponding to the original signal [17]. Within this field, PARAFAC seems applicable as it is an effective tool for data reduction as proven for multiple image compression [33]. Consequently, work analyzing PARAFAC's ability to solve the overcomplete representation in the wavelet analysis could also prove an important application of the model.

Although the functional magnetic resonance imaging, fMRI, has taken much focus away from the EEG and ERP, EEG/ERP still has great potentials. First of all the EEG is much easier to use for experimentation as the subjects can stay in a natural environment rather than having to be put inside a scanner. Furthermore, the EEG offers a much higher temporal resolution below the micro second range whereas the fMRI still works in the ms range and is most probably limited to this range. Consequently, ways of integrating EEG and fMRI has lately gotten much attention. Hopefully, the PARAFAC analysis of the EEG can help in this work decomposing the EEG signal into atoms that can be related to the fMRI signal. This application of the PARAFAC model has already gotten some attention [24]. As PARAFAC shows promising results in the field of EEG its application to MEG will very likely also be good. Consequently, multi-way array analysis will probably in the future become an important tool in the analysis of brain-recordings from a variety of scanning techniques.

Although Harshman proposed the use of PARAFAC on EEG in 1970, the use has been very limited. Why this is so is hard to understand given its wide usability to explore the EEG/ERP. However, the limited use might have been the consequence of the fact that previous works didn't resolve the problem of degeneracy by imposing non-negativity. Furthermore, the PARAFAC analysis is very memory consuming and slow, putting great demands on computer power. In this analysis 2 GB of RAM was required in order to simultaneously analyze the ITPC of the 11 subjects in both conditions having 64 channels of 200 ms data sampled at 512 Hz at 60 different frequencies. These computer requirements might be the reason for the limited use of PARAFAC so far in the field of EEG/ERP research. Hopefully, however, PARAFAC will turn out to be an important tool that will be widely used in the future when analyzing brain-data.

References

- [1] Andersson, Claus A.; Bro, Rasmus “*The N-way Toolbox for MATLAB*” Chemometrics and Intelligent Laboratory Systems 52 2000 1–4, Toolbox can be downloaded from: <http://www.models.kvl.dl/source>
- [2] Arnfred, Sidse “*Proprioceptive Event Related Potentials: Gating and Task effects*” Article in press. *Clinical Neurophysiology*
- [3] Beal, J. Matthew “*VARIATIONAL ALGORITHMS FOR APPROXIMATE BAYESIAN INFERENCE*” Thesis submitted for the degree of Doctor of Philosophy of the University of London May 2003
- [4] Bishop, Christopher M. “*Neural Networks for Pattern Recognition*” Oxford University Press 1995
- [5] Bro, Rasmus “*Multi-way Analysis in the Food Industry Models, Algorithms, and Applications*” Ph.D. thesis
- [6] Busch, Niko A; Debener, Stefan; Kranczioch, Cornelia; Engel, Andreas K.; Herrmann et al., Christoph S. “*Size matters: effects of stimulus size, duration and eccentricity on the visual gamma-band response*” *Clinical Neurophysiology* 115 (2004) 1810–1820
- [7] Cole, H.W.; Ray, W.J. “*EEG correlates of emotional tasks related to attentional demands*” *Int. J. Psychophysiol* 3(1): 33-41, 1985
- [8] Delorme, Arnaud; Makeig, Scott “*EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis*” *Journal of Neuroscience Methods* 134 (2004), 9-21
- [9] EL-Bab, Mohamed Fath “*COGNITIVE EVENT RELATED POTENTIALS DURING A LEARNING TASK*”, A thesis presented for the degree of Doctor of Philosophy, Department of Clinical Neurological Sciences Faculty of Medicine, University of Southampton, United Kingdom, April 2001
- [10] Field, Aaron S.; Graupe, Daniel “*Topographic Component (Parallel Factor) analysis of Multichannel Evoked Potentials: Practical Issues in Trilinear Spatiotemporal Decomposition*” *Brain Topography*, Vol. 3, Nr. 4, 1991
- [11] Grasman, R. P. “*Sensor array signal processing and the neuro-electromagnetic inverse problem in functional connectivity analysis of the brain*” PhD-thesis, University of Amsterdam, Amsterdam (2004).
- [12] Griffiths, David J. “*Introduction to Electrodynamics*” Third edition Prentice Hall, 1999
- [13] Harshman, Richard A. “*FOUNDATIONS OF THE PARAFAC PROCEDURE: MODELS AND CONDITIONS FOR AN “EXPLANATORY” MULTIMODAL FACTOR ANALYSIS*” UCLA, December, 1970
- [14] Herrmann, Christoph S.; Munk, Matthias H.J.; Engel, Andreas K. “*Cognitive functions of gamma-band activity: memory match and utilization*” *TRENDS in Cognitive Sciences* Vol.8 No.8 August 2004
- [15] Herrmann, Christoph S; Lenz, Daniel; Junge, Stefanie ; Busch, Niko A; Maess, Burkhard “*Memory-matches evoke human gamma-responses*” *BMC Neuroscience* 2004, 5:13
- [16] Kandel, Eric R., Schwartz, James H., Jessell, Thomas M. “*Principles of Neural Science*” McGrawHill, fourth edition 2000, p. 106-108, 326
- [17] Koenig, T.; Marti-Lopez, F.; Valdes-Soas, P. “*Topographic Time-Frequency Decomposition of the EEG*” *NeuroImage* 14, p. 383-390, 2001

- [18] Kolendra, Thomas “*Adaptive tools in virtual environments*”, IMM-2002
- [19] Lathauwer, Lieven De.; Moor, Bart De; Vandewalle, Joos “*A MULTILINEAR SINGULAR VALUE DECOMPOSITION*” SIAM J. MATRIX ANAL. APPL. Vol. 21, No. 4, pp. 1253–1278
- [20] Lathauwer, Lieven De; Moor, Bart De; Vandewalle, Joos “*ON THE BEST RANK-1 AND RANK-(R_1, R_2, \dots, R_N) APPROXIMATION OF HIGHER-ORDER TENSORS.*” SIAM J. MATRIX ANAL. APPL. c Vol. 21, No. 4, pp. 1324–1342
- [21] Lathauwer, L. De ; Moor, B. De ; Vandewalle, J. “*Independent component analysis based on higher-order statistics only*”, Proceedings of the IEEE Signal Processing Workshop on Statistical Signal and Array Processing, Corfu, Greece, 1996, p. 356–359.
- [22] Lee, Daniel D.; Seung, H. Sebastian “*Algorithms for Non-negative Matrix Factorization*” Adv. Neural Info. Proc. Syst. 13, 556-562 (2001)
- [23] Martin, Carla D. Moravitz “*Tensor Decompositions Workshop Discussion Notes American Institute of Mathematics (AIM) Palo Alto, CA*” July 19-23, 2004
- [24] Martínez-Montes, E.; Valdés-Sosa, P.A.; Miwakeichi, F.; Goldman, R. I.; Cohen, MS. “*Concurrent EEG/fMRI analysis by multiway Partial Least Squares*” NeuroImage 22, p. 1023-1034, 2004.
- [25] Miwakeichi, Fumikazu; Martínez-Montes, Eduardo; Valdés-Sosa, Pedro A.; Nishiyama, Nobuaki; Mizuhara, Hiroaki; Yamaguchi, Yoko “*Decomposing EEG data into space-time-frequency components using Parallel Factor Analysis*” NeuroImage 22 (2004) p. 1035-1045.
- [26] Möcks, Joachim “*Decomposing Event-Related Potentials: A New Topographic Component Model*”, Biological Psychology 26 (1988), 199-215
- [27] Nielsen, Frederik Brink “*Variational Approach to Factor Analysis and Related Models*” Thesis, Technical University of Denmark, 2004.
- [28] Nunez, Paul L. “*Electric Fields Of the Brain*”, Oxford University Press 1981, p.50-300, p. 445-450
- [29] Nunez, Paul L. “*Neocortical Dynamics and Human EEG Rhythms*” Oxford university press 1995 chapter 1,4
- [30] Pfurtscheller, G., Lopes da Silva, F.H. “*Event-related EEG/MEG synchronization and desynchronization: basic principles*”, Clinical Neurophysiology 110 (1999) p. 1842-1857
- [31] Sidiropoulos, Nicholas D., Bro, Rasmus “*On the uniqueness of multilinear decomposition of N-way arrays*” J. Chemometrics 2000; 14 p. 229–239
- [32] Todorovska, M.I., “*Estimation of Instantaneous Signals Using the Continuous Wavelet Transform*”, University of Southern California, Department of Civil Engineering, Report CE-01-07, 2001.
- [33] Wang, Hongcheng; Ahuja, Narendra “*Compact Representation of Multidimensional Data Using Tensor Rank-One Decomposition*” Beckman Institute, University of Illinois at Urbana-Champaign, USA
- [34] Windhorst, U. and Johansson, H “*Modern Techniques in Neuroscience Research*”, Springer Verlag 1999
- [35] Wrightt, J.J., Kydd, R.R. “*The electroencephalogram and cortical neural networks*”, Network 3 1992 page 341-362
- [36] <http://www.univ.trieste.it/~brain/Segnali/Segnali2.html>

[37] <http://www.departments.bucknell.edu/linguistics/lectures/05lect14.html>

[38] <http://www.benbest.com/science/anatmind/FigV5.gif>

[39] http://bio.winona.msus.edu/bates/Human%20Bio/Images/hb10_05.jpg

[40] <http://psyphz.psych.wisc.edu/.../EEGintro.htm>

[41] http://www.models.kvl.dk/research/data/Amino_Acid_fluo/index.asp

Appendix A: Theorems with proofs

Theorem 1 (Jensens Inequality)

Let $f : R \rightarrow R$ be a convex function, $\lambda_i > 0 \forall i$ and $\sum_i \lambda_i = 1$ then

$$f\left(\sum_{j=1}^M \lambda_j x_j\right) \geq \sum_{j=1}^M \lambda_j f(x_j)$$

Proof:

Let $x_i \in [a; b]$ where $a, b \in \mathcal{R}$ and let $t \in [0; 1]$. As f is convex the following holds:

$$f((1-t)a + tb) \geq (1-t)f(a) + tf(b)$$

as $x_i \in [a; b]$ it follows that $x_i = (1-t_i)a + t_i b$. This gives:

$$\sum_i \lambda_i x_i = \sum_i \lambda_i ((1-t_i)a + t_i b) = a \sum_i \lambda_i (1-t_i) + b \sum_i \lambda_i t_i = (1-t)a + tb \in [a; b] \text{ as } t = \sum_i \lambda_i t_i$$

Proof now done by induction;

M=2:

$$f(\lambda_1 x_1 + \lambda_2 x_2) = f((1-\lambda_2)x_1 + \lambda_2 x_2) \geq (1-\lambda_2)f(x_1) + \lambda_2 f(x_2) = \lambda_1 f(x_1) + \lambda_2 f(x_2)$$

Assume $f\left(\sum_{j=1}^M \lambda_j x_j\right) \geq \sum_{j=1}^M \lambda_j f(x_j)$ valid for M-1. We now have:

$$\begin{aligned} f\left(\sum_{j=1}^M \lambda_j x_j\right) &= f\left(\sum_{j=1}^{M-1} \lambda_j x_j + \lambda_M x_M\right) = f\left((1-\lambda_M) \sum_{j=1}^{M-1} \frac{\lambda_j}{(1-\lambda_M)} x_j + \lambda_M x_M\right) \geq \\ (1-\lambda_M) f\left(\sum_{j=1}^{M-1} \frac{\lambda_j}{(1-\lambda_M)} x_j\right) &+ \lambda_M f(x_M) \geq (1-\lambda_M) \sum_{j=1}^{M-1} \frac{\lambda_j}{(1-\lambda_M)} f(x_j) + \lambda_M f(x_M) = \sum_{j=1}^M \lambda_j f(x_j) \end{aligned}$$

Theorem 2

Given the following model and assumptions:

Model :	$\mathbf{X}^{(i)} = \mathbf{AD}^{(i)}\mathbf{S} + \mathbf{E}^{(i)}$
Assumptions :	$\mathcal{N}(\mathbf{e}^{(i)} \mid \mathbf{0}, \mathbf{\Psi}^{(i)}) \quad (1)$
	$\mathcal{N}(\mathbf{s}_j \mid \mathbf{0}, \mathbf{I}) \quad (2)$
	$\mathbf{\Psi}^{(i)}, \mathbf{D}_i \text{ diagonal} \quad (3)$

The following holds:

$$p(\mathbf{s}_i | \mathbf{A}, \{\mathbf{x}_i^{(m)}, \mathbf{D}^{(m)}, \mathbf{\Psi}^{(m)}\}_{m=1}^M) \propto \exp\left(\mathbf{s}_i^T \mathbf{D}^{(m)} \mathbf{A}^T \mathbf{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} - \frac{1}{2} \mathbf{s}_i^T \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \mathbf{\Psi}^{(m)-1} \mathbf{AD}^{(m)} \right) \mathbf{s}_i\right)$$

Proof:

$$\begin{aligned} & p(\mathbf{s}_i | \mathbf{A}, \{\mathbf{x}_i^{(m)}, \mathbf{D}^{(m)}, \mathbf{\Psi}^{(m)}\}_{m=1}^M) = \\ & \frac{p(\mathbf{x}_i^{(m)} | \mathbf{s}_i, \mathbf{A}, \{\mathbf{D}^{(m)}, \mathbf{\Psi}^{(m)}\}_{m=1}^M) p(\mathbf{s}_i | \mathbf{A}, \{\mathbf{D}^{(m)}, \mathbf{\Psi}^{(m)}\}_{m=1}^M)}{p(\mathbf{x}_i^{(m)} | \mathbf{A}, \{\mathbf{D}^{(m)}, \mathbf{\Psi}^{(m)}\}_{m=1}^M)} \quad \infty_2 \\ & p(\mathbf{x}_i^{(m)} | \mathbf{s}_i, \mathbf{A}, \{\mathbf{D}^{(m)}, \mathbf{\Psi}^{(m)}\}_{m=1}^M) \cdot p(\mathbf{s}_i | \mathbf{A}, \{\mathbf{D}^{(m)}, \mathbf{\Psi}^{(m)}\}_{m=1}^M) \quad \infty_3 \\ & \exp\left(-\frac{1}{2} \sum_{m=1}^M (\mathbf{x}_i^{(m)} - \mathbf{AD}^{(m)} \mathbf{s}_i)^T \mathbf{\Psi}^{(m)-1} (\mathbf{x}_i^{(m)} - \mathbf{AD}^{(m)} \mathbf{s}_i)\right) \exp\left(-\frac{1}{2} (\mathbf{s}_i - \mathbf{0})^T \mathbf{I}^{-1} (\mathbf{s}_i - \mathbf{0})\right) = \\ & \exp\left(-\frac{1}{2} \sum_{m=1}^M (\mathbf{x}_i^{(m)T} \mathbf{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} + (\mathbf{AD}^{(m)} \mathbf{s}_i)^T \mathbf{\Psi}^{(m)-1} \mathbf{AD}^{(m)} \mathbf{s}_i - \mathbf{x}_i^{(m)T} \mathbf{\Psi}^{(m)-1} \mathbf{AD}^{(m)} \mathbf{s}_i - (\mathbf{AD}^{(m)} \mathbf{s}_i)^T \mathbf{\Psi}^{(m)-1} \mathbf{x}_i^{(m)}) - \frac{1}{2} \mathbf{s}_i^T \mathbf{I}^{-1} \mathbf{s}_i\right) \quad \infty_4 \\ & \exp\left(\mathbf{s}_i^T \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \mathbf{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} - \frac{1}{2} \mathbf{s}_i^T \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \mathbf{\Psi}^{(m)-1} \mathbf{AD}^{(m)} \right) \mathbf{s}_i\right) \end{aligned}$$

1. Follows from Bayes theorem
2. Denominator is a normalization constant
3. Follows from assumption 1 and 2
4. Result of the fact that $\mathbf{x}_i^{(m)T} \mathbf{\Psi}^{(m)-1} \mathbf{AD}^{(m)} \mathbf{s}_i = (\mathbf{x}_i^{(m)T} \mathbf{\Psi}^{(m)-1} \mathbf{AD}^{(m)} \mathbf{s}_i)^T = \mathbf{s}_i^T \mathbf{D}^{(m)} \mathbf{A}^T \mathbf{\Psi}^{(m)-1} \mathbf{x}_i^{(m)}$ and

$$\sum_{m=1}^M \mathbf{x}_i^{(m)T} \mathbf{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} = \text{Constant}$$

Theorem 3 (The E step of PARAFAC)

Given the model, assumptions and result of Theorem 2 the following holds:

$$\begin{aligned}\Sigma_{\mathbf{S}} &= \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \Psi^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \right)^{-1} \\ \langle \mathbf{S} \rangle &= \Sigma_{\mathbf{S}} \left(\sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \Psi^{(m)-1} \mathbf{X}_m \right) \\ \langle \mathbf{S} \mathbf{S}^T \rangle &= N \Sigma_{\mathbf{S}} + \langle \mathbf{S} \rangle \langle \mathbf{S} \rangle^T\end{aligned}$$

Proof:

As the posterior distribution of \mathbf{S} i.e. the distribution of \mathbf{S} conditioned on the model parameter is Gaussian distributed we have:

$$\begin{aligned}p(\mathbf{S} | \theta) &= \prod_{i=1}^N p(s_i | \theta) \propto \prod_{i=1}^N \exp \left(-\frac{1}{2} (s_i - \langle s_i \rangle)^T \Sigma_{\mathbf{S}}^{-1} (s_i - \langle s_i \rangle) \right) \\ &= \exp \left(\sum_{i=1}^N \left(-\frac{1}{2} s_i^T \Sigma_{\mathbf{S}}^{-1} s_i + s_i^T \Sigma_{\mathbf{S}}^{-1} \langle s_i \rangle - \frac{1}{2} \langle s_i \rangle^T \Sigma_{\mathbf{S}}^{-1} \langle s_i \rangle \right) \right) \\ &\propto \underbrace{\exp \left(\sum_{i=1}^N \left(-\frac{1}{2} s_i^T \Sigma_{\mathbf{S}}^{-1} s_i + s_i^T \Sigma_{\mathbf{S}}^{-1} \langle s_i \rangle \right) \right)}_{\text{result 1}}\end{aligned}$$

1. Follows as $\frac{1}{2} \langle \mathbf{S} \rangle^T \Sigma_{\mathbf{S}}^{-1} \langle \mathbf{S} \rangle$ is a constant

As:

$$p(s_i | \mathbf{A}, \{ \mathbf{x}_i^{(m)}, \mathbf{D}^{(m)}, \Psi^{(m)} \}_{m=1}^M) \propto \exp \left(s_i^T \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \Psi^{(m)-1} \mathbf{x}_i^{(m)} - \frac{1}{2} s_i^T \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \Psi^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \right) s_i \right)$$

It follows

$$\begin{aligned}p(\mathbf{S} | \theta) &= \prod_{i=1}^N p \left(s_i | \mathbf{A}, \{ \mathbf{x}_i^{(m)}, \mathbf{D}^{(m)}, \Psi^{(m)} \}_{m=1}^M \right) \propto \\ &\prod_{i=1}^N \exp \left(s_i^T \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \Psi^{(m)-1} \mathbf{x}_i^{(m)} - \frac{1}{2} s_i^T \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \Psi^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \right) s_i \right) = \\ &\exp \left(\underbrace{\sum_{i=1}^N \left(s_i^T \mathbf{D}^{(m)} \mathbf{A}^T \Psi^{(m)-1} \mathbf{x}_i^{(m)} - \frac{1}{2} s_i^T \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}^{(m)} \mathbf{A}^T \Psi^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \right) s_i \right)}_{\text{result 2}} \right)\end{aligned}$$

When comparing result 1 with result 2, it is immediately seen that:

$$\Sigma_{\mathbf{S}} = \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}_m \mathbf{A}^T \Psi_m^{-1} \mathbf{A} \mathbf{D}_m \right)^{-1}$$

$$\langle \mathbf{S} \rangle = \Sigma_{\mathbf{S}} \left(\sum_{m=1}^M \mathbf{D}_m \mathbf{A}^T \Psi_m^{-1} \mathbf{X}_m \right)$$

The last part of the theorem follows by noting that:

$$\begin{aligned} \text{Cov}(\mathbf{X}) &= \int (\mathbf{X} - \langle \mathbf{X} \rangle)(\mathbf{X} - \langle \mathbf{X} \rangle)^T p(\mathbf{X}) d\mathbf{X} = \\ &= \int \mathbf{X} \mathbf{X}^T p(\mathbf{X}) d\mathbf{X} - \int \mathbf{X} \langle \mathbf{X} \rangle^T p(\mathbf{X}) d\mathbf{X} - \int \langle \mathbf{X} \rangle \mathbf{X}^T p(\mathbf{X}) d\mathbf{X} + \int \langle \mathbf{X} \rangle \langle \mathbf{X} \rangle^T p(\mathbf{X}) d\mathbf{X} = \\ &= \langle \mathbf{X} \mathbf{X}^T \rangle - \langle \mathbf{X} \rangle \langle \mathbf{X} \rangle^T - \langle \mathbf{X} \rangle \langle \mathbf{X} \rangle^T + \langle \mathbf{X} \rangle \langle \mathbf{X} \rangle^T = \langle \mathbf{X} \mathbf{X}^T \rangle - \langle \mathbf{X} \rangle \langle \mathbf{X} \rangle^T \end{aligned}$$

As each \mathbf{s}_i has the covariance $\Sigma_{\mathbf{s}}$ it follows $\text{cov}(\mathbf{S}) = N \Sigma_{\mathbf{s}}$. Therefore:

$$\text{Cov}(\mathbf{S}) = \langle \mathbf{S} \mathbf{S}^T \rangle - \langle \mathbf{S} \rangle \langle \mathbf{S} \rangle^T \Leftrightarrow \langle \mathbf{S} \mathbf{S}^T \rangle = \text{Cov}(\mathbf{S}) + \langle \mathbf{S} \rangle \langle \mathbf{S} \rangle^T = N \Sigma_{\mathbf{s}} + \langle \mathbf{S} \rangle \langle \mathbf{S} \rangle^T$$

Theorem 4 : (The likelihood of EMPARAFAC)

Given the model with assumptions of Theorem 2 the following holds:

$$\mathcal{L}(\boldsymbol{\theta}) = -\frac{N}{2} \sum_{m=1}^M \ln |\boldsymbol{\Psi}^{(m)}| - \frac{1}{2} \sum_{i=1}^N \sum_{m=1}^M \mathbf{x}_i^{mT} \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} + \text{tr} \left(\mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \right) - 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle + \text{const}$$

Proof:

To prove the statement we make use the following result:

$$\langle \mathbf{s}_i^T \mathbf{W} \mathbf{s}_i \rangle = \text{tr}(\mathbf{W} \text{cov}(\mathbf{s}_i, \mathbf{s}_i)) + \langle \mathbf{s}_i^T \rangle \mathbf{W} \langle \mathbf{s}_i \rangle \quad (*)$$

Writing out the likelihood we get:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= \left\langle \ln \prod_{i=1}^N \prod_{m=1}^M (2\pi)^{-\frac{d}{2}} |\boldsymbol{\Psi}^{(m)}|^{-\frac{1}{2}} \exp \left[\left(\mathbf{x}_i^{(m)} - \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right)^T \boldsymbol{\Psi}^{(m)-1} \left(\mathbf{x}_i^{(m)} - \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right) \right] \right\rangle = \\ &= -\frac{dMN}{2} \ln(2\pi) - \frac{N}{2} \sum_{m=1}^M \ln |\boldsymbol{\Psi}^{(m)}| - \frac{1}{2} \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{mT} \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} + \right. \\ &\quad \left. \langle \mathbf{s}_i^T \mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \rangle - \langle 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \rangle \right) = \\ &= -\frac{dMN}{2} \ln(2\pi) - \frac{N}{2} \sum_{m=1}^M \ln |\boldsymbol{\Psi}^{(m)}| - \frac{1}{2} \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{mT} \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} + \right. \\ &\quad \left. \mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \text{cov}(\mathbf{s}_i, \mathbf{s}_i) + \langle \mathbf{s}_i^T \rangle \mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle - 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle \right) = \\ &= -\frac{N}{2} \sum_{m=1}^M \ln |\boldsymbol{\Psi}^{(m)}| - \frac{1}{2} \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{mT} \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} + \text{tr} \left(\mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \right) - 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle \right) + \text{const} \end{aligned}$$

1. Follows from assumption 1, the expectation is taken as the expectation of \mathbf{S} , making the likelihood only depend on $\boldsymbol{\theta}$.
2. Follow from assumption 3.
3. Follows from (*) and assumption 2.

Theorem 5 (The M step of PARAFAC)

Given the likelihood function:

$$\mathcal{L}(\boldsymbol{\theta}) = -\frac{N}{2} \sum_{m=1}^M \ln |\boldsymbol{\Psi}^{(m)}| - \frac{1}{2} \sum_{i=1}^N \sum_{m=1}^M \mathbf{x}_i^{mT} \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} + \text{tr} \left(\mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle \right) - 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle + \text{const}$$

The following holds:

$$\begin{aligned} \mathbf{a}_k &= \left(\sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \mathbf{D}^{(m)} \right)_k \left(\sum_{m=1}^M \left(\boldsymbol{\Psi}^{(m)-1} \right)_{kk} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle^T \mathbf{D}^{(m)} \right)^{-1} \\ \mathbf{d}_m &= \left(\langle \mathbf{S} \mathbf{S}^T \rangle \bullet \left(\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \right) \right)^{-1} \text{vec} \left[\text{diag} \left[\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \right] \right] \\ \boldsymbol{\Psi}^{(m)} &= \frac{1}{N} \text{diag} \left[\mathbf{X}^{(m)} \mathbf{X}^{(m)T} + \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle \mathbf{D}^{(m)} \mathbf{A}^T - 2 \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \rangle \mathbf{X}^{(m)T} \right] \end{aligned}$$

proof:

$$\mathbf{a}_k = \left(\sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \mathbf{D}^{(m)} \right)_k \left(\sum_{m=1}^M \left(\boldsymbol{\Psi}^{(m)-1} \right)_{kk} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle^T \mathbf{D}^{(m)} \right)^{-1}$$

$$\begin{aligned} \frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \mathbf{A}} &= \frac{\partial \left(\sum_{i=1}^N \sum_{m=1}^M \text{tr} \left(\mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \right) - 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle \right)}{\partial \mathbf{A}} \Big|_1 \\ &= \frac{\partial \left(\sum_{i=1}^N \sum_{m=1}^M \text{tr} \left(\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \mathbf{D}^{(m)} \right) - 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle \right)}{\partial \mathbf{A}} \Big|_2 \end{aligned}$$

$$\sum_{i=1}^N \sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \mathbf{D}^{(m)} + \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \left(\mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \mathbf{D}^{(m)} \right)^T - 2 \boldsymbol{\Psi}^{(m)-1} 2 \mathbf{x}_i^{(m)} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle^T \Big|_3$$

$$\sum_{i=1}^N \sum_{m=1}^M 2 \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \mathbf{D}^{(m)} - 2 \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} \langle \mathbf{s}_i \rangle^T \mathbf{D}^{(m)} \Rightarrow$$

$$\sum_{i=1}^N \sum_{m=1}^M 2 \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \mathbf{D}^{(m)} - 2 \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} \langle \mathbf{s}_i \rangle^T \mathbf{D}^{(m)} = 0 \Leftrightarrow \Big|_4$$

$$\sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle \mathbf{D}^{(m)} = \sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \mathbf{D}^{(m)} \Rightarrow$$

$$\sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{a}_k \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle \mathbf{D}^{(m)} = \left(\sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \mathbf{D}^{(m)} \right)_k \Leftrightarrow$$

$$\mathbf{a}_k = \left(\sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \mathbf{D}^{(m)} \right)_k \left(\sum_{m=1}^M \boldsymbol{\Psi}^{(m)-1} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle \mathbf{D}^{(m)} \right)^{-1}$$

1. follows as $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA}) = \text{tr}(\mathbf{CAB})$

2. result of $\frac{\partial \text{tr}(\mathbf{A}^T \mathbf{B} \mathbf{A} \mathbf{C})}{\partial \mathbf{A}} = \mathbf{B} \mathbf{A} \mathbf{C} + \mathbf{B}^T \mathbf{A} \mathbf{C}^T$ and $\frac{\partial \mathbf{b}^T \mathbf{A} \mathbf{c}}{\partial \mathbf{A}} = \mathbf{b} \mathbf{c}^T$
3. comes by $(\mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \mathbf{D}^{(m)})^T = (\mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \mathbf{D}^{(m)})$
4. consequence of $\sum_{i=1}^N \langle \mathbf{s}_i \mathbf{s}_i^T \rangle = N \langle \mathbf{S} \mathbf{S}^T \rangle$ and $\sum_{i=1}^N \mathbf{x}_i^{(m)} \langle \mathbf{s}_i \rangle^T = N \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T$

$$\mathbf{d}_m = \left(\langle \mathbf{S} \mathbf{S}^T \rangle \bullet (\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A}) \right)^{-1} \text{vec} \left[\text{diag} \left[\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \right] \right]$$

$$\begin{aligned} \frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \mathbf{D}^{(m)}} &= \frac{\partial \left(\sum_{i=1}^N \sum_{m=1}^M \text{tr} \left(\mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \right) - 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle \right)}{\partial \mathbf{D}^{(m)}} = \\ & \sum_{i=1}^N \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle + \left(\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \right)^T \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle^T - 2 \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} \langle \mathbf{s}_i \rangle^T = \\ & 2 \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle - 2 \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T = 0 \Rightarrow \\ & \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle = \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \Rightarrow \end{aligned}$$

$$\mathbf{d}_m = \left(\langle \mathbf{S} \mathbf{S}^T \rangle \bullet (\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A}) \right)^{-1} \text{vec} \left(\text{diag} \left[\mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{X}^{(m)} \langle \mathbf{S} \rangle^T \right] \right)$$

where $\mathbf{D}^{(m)} = \text{diag}(\mathbf{d}_m)$

$$\boldsymbol{\Psi}^{(m)} = \frac{1}{N} \text{diag} \left[\mathbf{X}^{(m)} \mathbf{X}^{(m)T} + \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle \mathbf{D}^{(m)} \mathbf{A}^T - 2 \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \rangle \mathbf{X}^{(m)T} \right]$$

$$\begin{aligned} \frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\Psi}^{(m)-1}} &= \frac{\partial -\frac{N}{2} \sum_{m=1}^M \ln |\boldsymbol{\Psi}^{(m)}| - \frac{1}{2} \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{x}_i^{(m)} + \text{tr} \left(\mathbf{D}^{(m)} \mathbf{A}^T \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \right) - 2 \mathbf{x}_i^{(m)T} \boldsymbol{\Psi}^{(m)-1} \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle \right)}{\partial \boldsymbol{\Psi}^{(m)-1}} = \\ & -\frac{N}{2} \boldsymbol{\Psi}^{(m)} - \frac{1}{2} \sum_{i=1}^N \sum_{m=1}^M \mathbf{x}_i^{(m)} \mathbf{x}_i^{(m)T} + \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \mathbf{D}^{(m)} \mathbf{A}^T - 2 \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{s}_i \rangle \mathbf{x}_i^{(m)T} = 0 \Rightarrow \\ & \boldsymbol{\Psi}^{(m)} = \frac{1}{N} \text{diag} \left(\sum_{m=1}^M \mathbf{X}^{(m)} \mathbf{X}^{(m)T} + \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \mathbf{S}^T \rangle \mathbf{D}^{(m)} \mathbf{A}^T - 2 \mathbf{A} \mathbf{D}^{(m)} \langle \mathbf{S} \rangle \mathbf{X}^{(m)T} \right) \end{aligned}$$

1. follows by: $\frac{\partial \ln |\mathbf{A}|}{\partial \mathbf{A}} = (\mathbf{A}^T)^{-1}$, $\frac{\partial \mathbf{b}^T \mathbf{A}^{-1} \mathbf{c}}{\partial \mathbf{A}} = -\mathbf{A}^{-T} \mathbf{b} \mathbf{c}^T \mathbf{A}^{-T}$ and

$$\frac{\partial \text{tr}(\mathbf{B} \mathbf{A}^{-1} \mathbf{C})}{\partial \mathbf{A}} = -(\mathbf{A}^{-1} \mathbf{C} \mathbf{B} \mathbf{A}^{-1})^T$$

Theorem 6 (BIC for EMPARAFAC)

The Bayesian Information Criterion (BIC) for the EMPARAFAC model is given by:

$$p(\mathbf{D} | F) \approx \left(\prod_{m=1}^M \left((2\pi)^{-d/2} |\boldsymbol{\Psi}^{(m)}|^{-\frac{1}{2}} \right)^N \right) \exp \left(-0.5 \cdot \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{(m)} - \mathbf{AD}^{(m)} \langle \mathbf{s}_i \rangle \right)^T \boldsymbol{\Psi}^{(m)-1} \left(\mathbf{x}_i^{(m)} - \mathbf{AD}^{(m)} \langle \mathbf{s}_i \rangle \right) \right) N^{-\frac{D}{2}}$$

F is the number of factors in the PARAFAC model, and $D = F(P + M)$ is the number of ‘effectively’ free parameters, where P is the number of rows of $\mathbf{X}^{(m)}$.

Proof:

From eq. 1.24 we have $p(\mathbf{D} | F) \approx p(\mathbf{D} | \boldsymbol{\theta}_{ML}, F) N^{-\frac{D}{2}}$. Combining this with eq. 1.53 we get:

$$\begin{aligned} p(\mathbf{D} | \boldsymbol{\theta}_{ML}, F) &= \int p \left(\left\{ \mathbf{X}^{(m)} \right\}_{m=1}^M | \mathbf{S}, \mathbf{A}, \left\{ \mathbf{D}^{(m)}, \boldsymbol{\Psi}^{(m)} \right\}_{m=1}^M \right) p(\mathbf{S}) d\mathbf{S} = \\ &= \int \prod_{i=1}^N p \left(\left\{ \mathbf{x}_i^{(m)} \right\}_{m=1}^M | \mathbf{s}_i, \mathbf{A}, \left\{ \mathbf{D}^{(m)}, \boldsymbol{\Psi}^{(m)} \right\}_{m=1}^M \right) p(\mathbf{s}_i) d\mathbf{s}_i = \\ &= \left(\prod_{m=1}^M \left((2\pi)^{-d/2} |\boldsymbol{\Psi}^{(m)}|^{-\frac{1}{2}} \right)^N \right) \exp \left(-0.5 \cdot \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{(m)} - \mathbf{AD}^{(m)} \langle \mathbf{s}_i \rangle \right)^T \boldsymbol{\Psi}^{(m)-1} \left(\mathbf{x}_i^{(m)} - \mathbf{AD}^{(m)} \langle \mathbf{s}_i \rangle \right) \right) \Rightarrow \\ &= p(\mathbf{D} | \boldsymbol{\theta}_{ML}, M) N^{-\frac{D}{2}} = \\ &= \left(\prod_{m=1}^M \left((2\pi)^{-d/2} |\boldsymbol{\Psi}^{(m)}|^{-\frac{1}{2}} \right)^N \right) \exp \left(-0.5 \cdot \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{(m)} - \mathbf{AD}^{(m)} \langle \mathbf{s}_i \rangle \right)^T \boldsymbol{\Psi}^{(m)-1} \left(\mathbf{x}_i^{(m)} - \mathbf{AD}^{(m)} \langle \mathbf{s}_i \rangle \right) \right) N^{-\frac{D}{2}} \end{aligned}$$

The true number of free parameters is $D = F(P + M) + PM$, however the free parameters of $\boldsymbol{\Psi}$ is not directly part of the factor model, and does consequently not help improving the fit corresponding to having PM parameters. Therefore, the free parameters are only considered to be $D = F(P + M)$. As N and D depends on which modalities M , N and P pertains to, the BIC measure is greatly dependent on how each of the 3-way-array’s ways are defined in terms of the model $\mathbf{X}^{(i)} = \mathbf{AD}^{(i)} \mathbf{S}$.

Theorem 7: (BIC for Least Square optimization of PARAFAC)

The Bayesian Information Criterion (BIC) for a least square optimization of PARAFAC as defined in eq. 1.30 is when assuming $\mathbf{X}^{(m)}$ and \mathbf{S} to be i.i.d., given by:

$$\log(p(\mathbf{D} | F) \approx \log(p(\mathbf{D} | \boldsymbol{\theta}, F) N^{-\frac{D}{2}}) = -\frac{1}{2} MN \log(\sigma) - \frac{D}{2} \log N + const$$

F is the number of factors in the PARAFAC model, and $D = F(P + M + N)$ is the number of free parameters, where P is the number of rows of $\mathbf{X}^{(m)}$.

Proof:

$$\begin{aligned} p(\mathbf{D} | \boldsymbol{\theta}, F) &= p\left(\left\{\mathbf{X}^{(m)}\right\}_{m=1}^M | \mathbf{S}, \mathbf{A}, \sigma, \left\{\mathbf{D}^{(m)}\right\}_{m=1}^M\right) = \\ &\prod_{i=1}^N p\left(\left\{\mathbf{x}_i^{(m)}\right\}_{m=1}^M | \mathbf{s}_i, \mathbf{A}, \sigma, \left\{\mathbf{D}^{(m)}\right\}_{m=1}^M\right) = \\ &\left(\prod_{m=1}^M \left((2\pi)^{-\frac{d}{2}} |\boldsymbol{\sigma} \cdot \mathbf{I}|^{-\frac{1}{2}}\right)^N\right) \exp\left(-0.5 \cdot \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i\right)^T (\boldsymbol{\sigma} \cdot \mathbf{I})^{-1} \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i\right)\right) \Rightarrow \\ &p(\mathbf{D} | \boldsymbol{\theta}, M) N^{-\frac{D}{2}} = \\ &\left(\prod_{m=1}^M \left((2\pi)^{-\frac{d}{2}} \boldsymbol{\sigma}^{-\frac{1}{2}}\right)^N\right) \exp\left(-0.5 \cdot \frac{1}{\sigma} \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i\right)^T \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i\right)\right) N_0^{-\frac{D}{2}} = \\ &\left(\prod_{m=1}^M \left((2\pi)^{-\frac{d}{2}} \boldsymbol{\sigma}^{-\frac{1}{2}}\right)^N\right) \exp(-0.5NM) N_0^{-\frac{D}{2}} \Rightarrow \\ &\log(p(\mathbf{D} | \boldsymbol{\theta}, M) N^{-\frac{D}{2}}) = -\frac{1}{2} MN \log(\sigma) - \frac{D}{2} \log N + const \end{aligned}$$

$$1. \text{ Consequence of } \sigma = \frac{1}{NM} \sum_{i=1}^N \sum_{m=1}^M \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i\right)^T \left(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i\right)$$

As N depends on which modality is considered observations, the BIC measure is greatly dependent on how each of the 3-way-array's ways are defined in terms of the model

$$\mathbf{X}^{(i)} = \mathbf{A}\mathbf{D}^{(i)}\mathbf{S}.$$

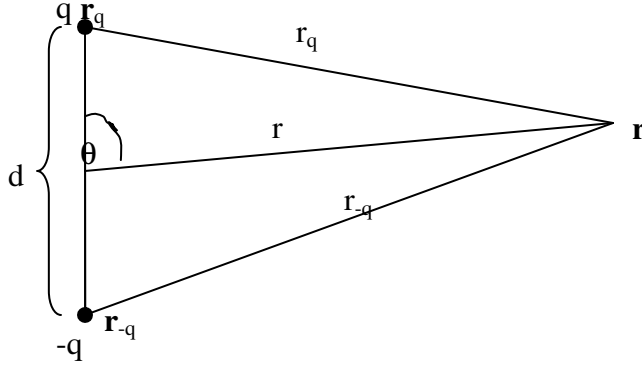
Theorem 8: (Potential of a dipole field)

The potential in a distance r of the center of a dipole, where the charges (q and $-q$) are separated by the distance d is given by:

$$\Phi_{dipole}(\mathbf{r}, \mathbf{d}) = \frac{1}{4\pi\epsilon_0} \frac{qd \cos \theta}{\|\mathbf{r}\|^2}$$

where θ is the angle between \mathbf{r} and the midpoint of the line connecting the two charges.

Proof:



$$r_q = \sqrt{r^2 + \frac{d^2}{2} - 2r\frac{d}{2}\cos\theta} = \sqrt{r^2 + \frac{d^2}{2} - rd\cos\theta}$$

$$r_{-q} = \sqrt{r^2 + \frac{d^2}{2} - 2r\frac{d}{2}\cos(\pi - \theta)} = \sqrt{r^2 + \frac{d^2}{2} + rd\cos(\theta)}$$

$$r = \|\mathbf{r}\|$$

$$d = \|\mathbf{d}\|$$

$$\|\vec{\mathbf{E}}(\mathbf{r}, \mathbf{d})\| = \frac{1}{4\pi\epsilon_0} \left(\frac{q}{r^2 + \frac{d^2}{2} - rd\cos\theta} - \frac{q}{r^2 + \frac{d^2}{2} + rd\cos(\theta)} \right) =$$

$$\frac{q}{4\pi\epsilon_0} \left(\frac{2rd\cos\theta}{\left(r^2 + \frac{d^2}{2} - rd\cos\theta\right)\left(r^2 + \frac{d^2}{2} + rd\cos(\theta)\right)} \right) \approx$$

$$\frac{q}{4\pi\epsilon_0} \left(\frac{2rd\cos\theta}{r^4} \right) = \frac{q}{4\pi\epsilon_0} \left(\frac{2d\cos\theta}{r^3} \right)$$

$$-\nabla\Phi(\mathbf{r}, \mathbf{d}) = \vec{\mathbf{E}}(\mathbf{r}, \mathbf{d}) \Rightarrow$$

$$\Phi(\mathbf{r}, \mathbf{d}) = \frac{1}{4\pi\epsilon_0} \frac{qd\cos\theta}{r^2} = \frac{1}{4\pi\epsilon_0} \frac{qd\cos\theta}{r^2}$$

1. Follows from the fact that $r \gg d$

Theorem 9 (Derivation of the VBEM algorithm for PARAFAC)

For the PARAFAC model of eq. 1.30, the following assumptions are made:

$$p(\mathbf{E}^{(m)}) = \prod_{i=1}^N \mathcal{N}(\mathbf{e}_i^{(m)} | 0, \text{diag}(\boldsymbol{\varphi}_m))$$

$$p(\mathbf{S}) = \prod_{i=1}^N \mathcal{N}(\mathbf{s}_i | 0, \mathbf{I})$$

$$p(\mathbf{A} | \boldsymbol{\alpha}) = \prod_{f=1}^F \mathcal{N}(\mathbf{a}_f | 0, \frac{1}{\alpha_f} \mathbf{I})$$

$$p(\boldsymbol{\alpha} | a^\alpha, b^\alpha) = \prod_{f=1}^F \mathcal{G}(\alpha_f | a_f^\alpha, b_f^\alpha)$$

$$p(\mathbf{d}_m | \boldsymbol{\mu}_m, \frac{1}{v_m}) = \mathcal{N}(\mathbf{d}_m | \boldsymbol{\mu}_m, \frac{1}{v_m} \mathbf{I})$$

$$p(\boldsymbol{\varphi}_m | \mathbf{a}_m^\varphi, \mathbf{b}_m^\varphi) = \prod_{p=1}^P \mathcal{G}(\varphi_{mp} | a_{mp}^\varphi, b_{mp}^\varphi)$$

where \mathbf{d}_m corresponds to the diagonal of $\mathbf{D}^{(m)}$

These yields the following update rules:

$$\boldsymbol{\Sigma}_S = \left(\mathbf{I} + \sum_{m=1}^M \langle \mathbf{d}_m^T \mathbf{d}_m \rangle \bullet \sum_{p=1}^P \langle \boldsymbol{\varphi}_{mp} \rangle \langle \mathbf{a}_p^T \mathbf{a}_p \rangle \right)^{-1}$$

$$\boldsymbol{\mu}_{s_i} = \boldsymbol{\Sigma}_S \sum_{m=1}^M \langle \mathbf{d}_m^T \rangle \langle \mathbf{A}^T \rangle \langle \boldsymbol{\varphi}_m \rangle \mathbf{d}_i$$

$$\boldsymbol{\Sigma}_a^p = \left\langle \left(\boldsymbol{\alpha} \mathbf{I} + \sum_{i=1}^N \sum_{m=1}^M \mathbf{s}_i^T \mathbf{D}^{(m)} \boldsymbol{\varphi}_j^{(m)} \mathbf{D}^{(m)} \mathbf{s}_i \right) \right\rangle^{-1} = \left\langle \langle \boldsymbol{\alpha} \rangle \mathbf{I} + \sum_{m=1}^M \langle \boldsymbol{\varphi}^{(m)} \rangle \langle \mathbf{d}^{(m)} \mathbf{d}^{(m)T} \rangle \bullet \sum_{i=1}^N \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \right\rangle^{-1}$$

$$\boldsymbol{\mu}_a^p = \boldsymbol{\Sigma}_a^p \left\langle \sum_{i=1}^N \sum_{m=1}^M \boldsymbol{\varphi}_p^{(m)} \mathbf{x}_i^{(m)T} \mathbf{D}^{(m)} \mathbf{s}_i \right\rangle = \boldsymbol{\Sigma}_a^p \left\langle \sum_{m=1}^M \langle \boldsymbol{\varphi}_p^{(m)} \rangle \langle \mathbf{D}^{(m)} \rangle \sum_{i=1}^N \mathbf{x}_i^{(m)T} \langle \mathbf{s}_i \rangle \right\rangle$$

$$\boldsymbol{\Sigma}_{\mathbf{D}^{(m)}} = \left(v_m \mathbf{I} + \sum_{j=1}^P \langle \boldsymbol{\varphi}_{mj} \rangle \langle \mathbf{a}_{:,j}^T \mathbf{a}_{:,j} \rangle \bullet \sum_{i=1}^N \langle \mathbf{s}_i^T \mathbf{s}_i \rangle \right)^{-1}$$

$$m_{\mathbf{D}^{(m)}} = \boldsymbol{\Sigma}_{\mathbf{D}^{(m)}} \left(v_m \boldsymbol{\mu}_m + \sum_{i=1}^N \sum_{j=1}^D \langle x_{ij}^{(m)} \rangle \boldsymbol{\varphi}_{mj} \langle \mathbf{s}_i \rangle \bullet \langle \mathbf{a}_{:,j} \rangle \right)$$

$$\hat{a}_f^\alpha = a_f^\alpha + \frac{p}{2}$$

$$\hat{b}_f^\alpha = b_f^\alpha + \frac{\langle \mathbf{a}_f^T \mathbf{a}_f \rangle}{2}$$

$$\hat{a}_{mp}^\varphi = a_{mp}^\varphi + \frac{N}{2}$$

$$\hat{b}_{mp}^\varphi = b_{mp}^\varphi + \frac{1}{2} \sum_{i=1}^N \left\langle \left(x_{ip}^{(m)} - \mathbf{a}_{:,p} \mathbf{D}^{(m)} \mathbf{s}_i \right)^2 \right\rangle$$

Proof:

$$\begin{aligned}
p(\mathbf{x}_i^{(m)}, \mathbf{s}_i, \mathbf{D}^{(m)}, \boldsymbol{\varphi}_m, \mathbf{A}, \alpha) &= \underbrace{k_1 |\boldsymbol{\varphi}_m|^{1/2} \exp\left(-\frac{1}{2}(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i)^T \boldsymbol{\varphi}_m (\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i)\right)}_{p(\mathbf{d}_i|\mathbf{s}_i, \mathbf{C}^{(m)}, \boldsymbol{\varphi}_m, \alpha)} \\
&\quad \underbrace{k_2 \exp\left(-\frac{1}{2}(\mathbf{s}_i)^T \mathbf{I}(\mathbf{s}_i)\right)}_{p(\mathbf{s}_i)} \\
&\quad \underbrace{k_3 \left|\frac{1}{v_m} \mathbf{I}\right|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{d}_m - \boldsymbol{\mu}_m)^T (v_m \mathbf{I})(\mathbf{d}_m - \boldsymbol{\mu}_m)\right)}_{p(\mathbf{D}_m)} \\
&\quad \underbrace{\prod_{p=1}^P \frac{1}{\Gamma(a_{mp}^\varphi)} \frac{1}{b_{mp}^\varphi} \boldsymbol{\varphi}_{mp}^{a_{mp}^\varphi - 1} \exp(-\boldsymbol{\varphi}_{mp} b_{mp}^\varphi)}_{p(\boldsymbol{\varphi}_m)} \\
&\quad \underbrace{\prod_{f=1}^F k_4 \left|\frac{1}{\alpha_f} \mathbf{I}\right|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{a}_f)^T \boldsymbol{\alpha}(\mathbf{a}_f)\right)}_{p(\mathbf{a}_f|\alpha)} \\
&\quad \underbrace{\prod_{f=1}^F \frac{1}{\Gamma(a_f^\alpha)} \frac{1}{b_f^\alpha} \alpha_f^{a_f^\alpha - 1} \exp(-\alpha_f b_f^\alpha)}_{p(\alpha)} \\
\ln p(\mathbf{X}, \mathbf{S}, \{\mathbf{D}^{(m)}, \boldsymbol{\varphi}_m\}_{m=1}^M, \mathbf{A}, \alpha) &= \sum_{i=1}^N \left\{ \underbrace{\sum_{m=1}^M \frac{1}{2} \ln |\boldsymbol{\varphi}_m| - \frac{1}{2}(\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i)^T \boldsymbol{\varphi}_m (\mathbf{x}_i^{(m)} - \mathbf{A}\mathbf{D}^{(m)}\mathbf{s}_i)}_{p(\mathbf{x}_i^{(m)}|\mathbf{s}_i, \mathbf{D}_m, \boldsymbol{\varphi}_m, \alpha)} - \frac{1}{2}(\mathbf{s}_i)^T \mathbf{I}(\mathbf{s}_i)} \right\} \\
&\quad \underbrace{\frac{1}{2} \ln \left|\frac{1}{v_m} \mathbf{I}\right| - \frac{1}{2}(\mathbf{d}_m - \boldsymbol{\mu}_m)^T v_m \mathbf{I}(\mathbf{d}_m - \boldsymbol{\mu}_m)}_{p(\mathbf{D}^{(m)})} + \underbrace{\sum_{p=1}^P -\ln \Gamma(a_{mp}^\varphi) - a_{mp}^\varphi \ln \left(\frac{1}{b_{mp}^\varphi}\right) + (a_{mp}^\varphi - 1) \ln \boldsymbol{\varphi}_{mp} - \boldsymbol{\varphi}_{mp} b_{mp}^\varphi}_{p(\boldsymbol{\varphi}_m)} \\
&\quad \underbrace{\sum_{f=1}^F -\frac{1}{2} \ln \left|\frac{1}{\alpha_f} \mathbf{I}\right| - \frac{1}{2}(\mathbf{a}_f)^T \boldsymbol{\alpha}(\mathbf{a}_f)}_{p(\mathbf{a}_f|\alpha)} + \underbrace{\sum_{f=1}^F -\ln \Gamma(a_f^\alpha) - a_f^\alpha \ln \left(\frac{1}{b_f^\alpha}\right) + (a_f^\alpha - 1) \ln \alpha_f - \alpha_f b_f^\alpha}_{p(\alpha)}
\end{aligned}$$

$$\begin{aligned}
q(\mathbf{s}_i) &: \sum_{m=1}^M \left(-\frac{1}{2} \left(\mathbf{x}_i^{(m)} - \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right)^T \varphi_m \left(\mathbf{x}_i^{(m)} - \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right) \right) - \frac{1}{2} (\mathbf{s}_i)^T \mathbf{I} (\mathbf{s}_i) \propto \\
& \sum_{m=1}^M \left(- \left(\mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right)^T \varphi_m \left(\mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right) - \mathbf{x}_i^{(m)} \varphi_m^{-1} \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right) - (\mathbf{s}_i)^T \mathbf{I} (\mathbf{s}_i) = \\
& \sum_{m=1}^M \left(-\mathbf{s}_i^T \mathbf{D}^{(m)T} \mathbf{A}^T \varphi_m \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i + \mathbf{x}_i^{(m)} \varphi_m^{-1} \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right) - (\mathbf{s}_i)^T \mathbf{I} (\mathbf{s}_i) = \\
& -\mathbf{s}_i^T \left(\mathbf{I} + \sum_{m=1}^M \mathbf{D}^{(m)T} \mathbf{A}^T \varphi_m \mathbf{A} \mathbf{D}^{(m)} \right) \mathbf{s}_i + \sum_{m=1}^M \mathbf{s}_i^T \mathbf{D}^{(m)T} \mathbf{A}^T \varphi_m \mathbf{x}_i^{(m)} \Rightarrow \\
\Sigma_{\mathbf{S}} &= \left(\mathbf{I} + \sum_{m=1}^M \left\langle \mathbf{D}^{(m)T} \mathbf{A}^T \varphi_m \mathbf{A} \mathbf{D}^{(m)} \right\rangle \right)^{-1} = \left(\mathbf{I} + \sum_{m=1}^M \left\langle \mathbf{d}_m^T \mathbf{d}_m \right\rangle \bullet \sum_{p=1}^P \left\langle \varphi_{mp} \right\rangle \left\langle \mathbf{a}_p^T \mathbf{a}_p \right\rangle \right)^{-1} \\
\boldsymbol{\mu}_{\mathbf{s}_i} &= \Sigma_{\mathbf{S}} \sum_{m=1}^M \left\langle \mathbf{D}^{(m)T} \mathbf{A}^T \varphi_m \mathbf{x}_i^{(m)} \right\rangle = \Sigma_{\mathbf{S}} \sum_{m=1}^M \left\langle \mathbf{D}^{(m)T} \right\rangle \left\langle \mathbf{A}^T \right\rangle \left\langle \varphi_m \right\rangle \mathbf{x}_i^{(m)} \\
1. \quad & \left\langle \mathbf{D}^{(m)T} \mathbf{A}^T \varphi_m \mathbf{A} \mathbf{D}^{(m)} \right\rangle = \left\langle \mathbf{D}^{(m)T} \sum_{p=1}^P \mathbf{a}_p^T \varphi_{mp} \mathbf{a}_p \mathbf{D}^{(m)} \right\rangle = \left\langle \mathbf{d}_m^T \mathbf{d}_m \right\rangle \bullet \sum_{p=1}^P \left\langle \varphi_{mp} \right\rangle \left\langle \mathbf{a}_p^T \mathbf{a}_p \right\rangle
\end{aligned}$$

$$\begin{aligned}
q(\mathbf{A}): & \sum_{i=1}^N \sum_{m=1}^M \left(- \left(\mathbf{x}_i^{(m)} - \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right)^T \varphi_m \left(\mathbf{x}_i^{(m)} - \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right) \right) - \sum_{f=1}^F \left(\mathbf{a}_f \right)^T \left(\alpha_f \mathbf{I} \right) \left(\mathbf{a}_f \right) = \\
& \sum_{i=1}^N \sum_{m=1}^M - \mathbf{s}_i^T \mathbf{D}^{(m)T} \mathbf{A}^T \varphi_m \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i + 2 \mathbf{x}_i^{(m)T} \varphi_m \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i - \sum_{f=1}^F \left(\mathbf{a}_f \right)^T \alpha_f \mathbf{I} \left(\mathbf{a}_f \right) = \\
& - \sum_{i=1}^N \sum_{m=1}^M \sum_{p=1}^P \left(\mathbf{a}_{p,:} \left(\mathbf{s}_i^T \mathbf{D}^{(m)T} \varphi_{mp} \mathbf{D}^{(m)} \mathbf{s}_i \right) \mathbf{a}_{p,:}^T - 2 \mathbf{a}_{p,:} \varphi_{mp} \mathbf{x}_i^{(m)T} \mathbf{D}^{(m)} \mathbf{s}_i \right) - \sum_{f=1}^F \left(\mathbf{a}_f \right)^T \alpha_f \mathbf{I} \left(\mathbf{a}_f \right) = \\
& \Rightarrow
\end{aligned}$$

$$\boldsymbol{\Sigma}_{\mathbf{a}}^p = \left\langle \left(\alpha_f \mathbf{I} + \sum_{i=1}^N \sum_{m=1}^M \mathbf{s}_i^T \mathbf{D}^{(m)} \varphi_{mp} \mathbf{D}^{(m)} \mathbf{s}_i \right) \right\rangle^{-1} = \left\langle \langle \alpha \rangle \mathbf{I} + \sum_{m=1}^M \langle \varphi_m \rangle \langle \mathbf{d}_m \mathbf{d}_m^T \rangle \bullet \sum_{i=1}^N \langle \mathbf{s}_i \mathbf{s}_i^T \rangle \right\rangle^{-1}$$

$$\boldsymbol{\mu}_{\mathbf{a}}^p = \boldsymbol{\Sigma}_{\mathbf{a}}^p \left\langle \sum_{i=1}^N \sum_{m=1}^M \varphi_{mj} \mathbf{x}_i^{(m)T} \mathbf{D}^{(m)} \mathbf{s}_i \right\rangle = \boldsymbol{\Sigma}_{\mathbf{a}}^p \left\langle \sum_{m=1}^M \langle \varphi_{mj} \rangle \langle \mathbf{D}^{(m)} \rangle \sum_{i=1}^N \mathbf{x}_i^{(m)T} \langle \mathbf{s}_i \rangle \right\rangle$$

$$q(\boldsymbol{\alpha}_f): \underbrace{- \frac{1}{2} \ln \left| \frac{1}{\alpha_f} \mathbf{I} \right| - \frac{1}{2} \left(\mathbf{a}_f \right)^T \left(\alpha_f \mathbf{I} \right) \left(\mathbf{a}_f \right) - \ln \Gamma(\alpha_f)}_{p(\boldsymbol{\alpha}_f | \boldsymbol{\alpha}_f)} - \underbrace{\alpha_f^\alpha \ln \left(\frac{1}{b_f^\alpha} \right) + (\alpha_f^\alpha - 1) \ln \alpha_f - \alpha_f b_f^\alpha}_{p(\boldsymbol{\alpha}_f)} =$$

$$\frac{p}{2} \ln \alpha_f - \frac{1}{2} \alpha_f \left(\mathbf{a}_f \right)^T \left(\mathbf{a}_f \right) - \ln \Gamma(\alpha_f) - \alpha_f^\alpha \ln \left(\frac{1}{b_f^\alpha} \right) + (\alpha_f^\alpha - 1) \ln \alpha_f - \alpha_f b_f^\alpha$$

$$\hat{a}_f^\alpha = \alpha_f^\alpha + \frac{p}{2}$$

$$\hat{b}_f^\alpha = b_f^\alpha + \left\langle \frac{\left(\mathbf{a}_f \right)^T \left(\mathbf{a}_f \right)}{2} \right\rangle$$

$$\langle \mathbf{a}_f \rangle = \frac{\hat{a}_f^\alpha}{\hat{b}_f^\alpha}$$

$$\begin{aligned}
q(\boldsymbol{\varphi}_{mp}) &: \sum_{i=1}^N \left(\frac{1}{2} \ln \varphi_{mj} - \frac{1}{2} \left(x_{ij}^{(m)} - \mathbf{a}_{:,j} \mathbf{D}^{(m)} \mathbf{s}_i \right)^T \varphi_{mj} \left(x_{ij}^{(m)} - \mathbf{a}_{:,j} \mathbf{D}^{(m)} \mathbf{s}_i \right) \right) + (a_{mj}^\varphi - 1) \ln \varphi_{mj} - \varphi_{mj} b_{mj}^\varphi = \\
& \left(a_{mp}^\varphi - 1 + \frac{N}{2} \right) \ln \varphi_{mj} - \varphi_{mj} \left(b_{mj}^\varphi + \sum_{i=1}^N \frac{1}{2} \left(x_{ij}^{(m)} - \mathbf{a}_{:,j} \mathbf{D}^{(m)} \mathbf{s}_i \right)^T \left(x_{ij}^{(m)} - \mathbf{a}_{:,j} \mathbf{D}^{(m)} \mathbf{s}_i \right) \right) \\
\hat{a}_{mp}^\varphi &= a_{mp}^\varphi + \frac{N}{2} \\
\hat{b}_{mp}^\varphi &= b_{mp}^\varphi + \frac{1}{2} \sum_{i=1}^N \left\langle \left(x_{ip}^{(m)} - \mathbf{a}_{:,p} \mathbf{D}^{(m)} \mathbf{s}_i \right)^2 \right\rangle \\
\langle \varphi_{mp} \rangle &= \frac{a_{mp}^\varphi}{b_{mp}^\varphi}
\end{aligned}$$

$$\begin{aligned}
q(\mathbf{D}^{(m)}) &: \sum_{i=1}^N \left(-\frac{1}{2} \left(\mathbf{x}_i^{(m)} - \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right)^T \varphi_{mj} \left(\mathbf{x}_i^{(m)} - \mathbf{A} \mathbf{D}^{(m)} \mathbf{s}_i \right) \right) - \frac{1}{2} (\mathbf{d}_m - \boldsymbol{\mu}_m)^T v_m \mathbf{I} (\mathbf{d}_m - \boldsymbol{\mu}_m) = \\
& \sum_{i=1}^N \sum_{p=1}^P \left(-\frac{1}{2} \left(x_{ij}^{(m)} - \mathbf{a}_{:,j} \mathbf{d}_m \mathbf{s}_i \right)^T \varphi_{mj} \left(x_{ij}^{(m)} - \mathbf{s}_i \mathbf{a}_{:,j} \mathbf{d}_m \right) \right) - \frac{1}{2} (\mathbf{d}_m - \boldsymbol{\mu}_m)^T v_m \mathbf{I} (\mathbf{d}_m - \boldsymbol{\mu}_m) = \\
& \sum_{i=1}^N \left(-\sum_{p=1}^P \left(\boldsymbol{\mu}_m^T v_m \mathbf{I} + x_{ip}^{(m)} \varphi_{mp} \mathbf{s}_i \mathbf{a}_{:,p} \right) \mathbf{d}_m - \frac{1}{2} \mathbf{d}_m^T \left(v_m \mathbf{I} + \varphi_{mp}^{-1} \mathbf{s}_i^T \mathbf{s}_i \mathbf{a}_{:,p}^T \mathbf{a}_{:,p} \right) \mathbf{d}_m \right) \\
\boldsymbol{\Sigma}_{\mathbf{D}^{(m)}} &= \left(v_m \mathbf{I} + \sum_{p=1}^P \langle \varphi_{mp} \rangle \langle \mathbf{a}_{:,p}^T \mathbf{a}_{:,p} \rangle \bullet \sum_{i=1}^N \langle \mathbf{s}_i^T \mathbf{s}_i \rangle \right)^{-1} \\
\mathbf{m}_{\mathbf{D}^{(m)}} &= \boldsymbol{\Sigma}_{\mathbf{D}^{(m)}} \left(v_m \boldsymbol{\mu}_m + \sum_{i=1}^N \sum_{p=1}^P \langle x_{ip}^{(m)} \rangle \varphi_{mp} \langle \mathbf{s}_i \rangle \bullet \langle \mathbf{a}_{:,p} \rangle \right)
\end{aligned}$$

Theorem 10: (Column-wise non-negativity constraints)

Consider the PARAFAC model as formulated by eq. 1.31, sought optimized in an unweighted least square sense. Let \mathbf{T} be defined as $\mathbf{T} = \mathbf{X}^{JK \times I} - (\mathbf{B}_{-F} | \otimes | \mathbf{C}_{-F}) \mathbf{A}_{-F}^T$ where \mathbf{M}_{-F} denotes that the F^{th} column has been removed from \mathbf{M} . Then finding the F^{th} column of \mathbf{A} subject to the constraint that the F^{th} column of \mathbf{A} has to be non-negative corresponds to finding the optimal value of the unconstrained problem and setting all negative values equal zero.

Proof (adapted from [3])

The optimal choice of \mathbf{A} unconstrained is the value that minimizes:

$$\min_{\mathbf{a}_F} \left\| \mathbf{T} - (\mathbf{b}_F | \otimes | \mathbf{c}_F) \mathbf{a}_F^T \right\|^2, \text{ let this value be denoted } \boldsymbol{\alpha}, \text{ i.e. } \boldsymbol{\alpha} = \frac{\mathbf{T}^T \mathbf{z}}{\mathbf{z}^T \mathbf{z}}, \text{ where } \mathbf{z} = (\mathbf{b}_F | \otimes | \mathbf{c}_F).$$

Furthermore, let $\mathbf{E} = \mathbf{T} - \mathbf{z} \boldsymbol{\alpha}^T$. It then follows that:

$$\begin{aligned} \min_{\mathbf{a}_F} \left\| \mathbf{E} + \mathbf{z} \boldsymbol{\alpha}^T - \mathbf{z} \mathbf{a}_F^T \right\|^2 &= \min_{\mathbf{a}_F} \left\| \mathbf{E} + \mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T \right\|^2 = \\ \min_{\mathbf{a}_F} \left(\text{tr}(\mathbf{E}^T \mathbf{E}) + 2 \text{tr} \mathbf{E}^T \left(\mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T \right) + \text{tr} \left[\mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T \mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T \right] \right) & \\ \min_{\mathbf{a}_F} \left(\text{tr} \left[\mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T \mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T \right] \right) &= \min_{\mathbf{a}_F} \left(\sum_{i=1}^p (\alpha_{pf} - a_{pi})^2 \right) \end{aligned}$$

1. Equality holds as $\mathbf{E}^T \mathbf{E}$ is a constant, furthermore,

$$\begin{aligned} \mathbf{E}^T \left(\mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T \right) &= \left(\mathbf{T}^T - \frac{\mathbf{T}^T \mathbf{z}}{\mathbf{z}^T \mathbf{z}} \mathbf{z}^T \right) \mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T = \left(\mathbf{T}^T \mathbf{z} - \frac{\mathbf{T}^T \mathbf{z}}{\mathbf{z}^T \mathbf{z}} \mathbf{z}^T \mathbf{z} \right) (\boldsymbol{\alpha} - \mathbf{a}_F)^T = \\ &= (\mathbf{T}^T \mathbf{z} - \mathbf{T}^T \mathbf{z}) (\boldsymbol{\alpha} - \mathbf{a}_F)^T = 0 \end{aligned}$$

2. Follows as $\mathbf{z}^T \mathbf{z}$ is a constant.

Notice: Proof also holds for weighted regression as $\boldsymbol{\alpha} = \frac{\mathbf{T}^T \boldsymbol{\Phi}^{-1} \mathbf{z}}{\mathbf{z}^T \boldsymbol{\Phi}^{-1} \mathbf{z}}$ and

$$\mathbf{E}^T \boldsymbol{\Phi}^{-1} \left(\mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T \right) = \left(\mathbf{T}^T - \frac{\mathbf{T}^T \boldsymbol{\Phi}^{-1} \mathbf{z}}{\mathbf{z}^T \boldsymbol{\Phi}^{-1} \mathbf{z}} \mathbf{z}^T \right) \boldsymbol{\Phi}^{-1} \mathbf{z} (\boldsymbol{\alpha} - \mathbf{a}_F)^T = 0$$

Theorem 11: (Regarding ICAPARAFAC)

Consider the 3-way array $\mathbf{X} \in R^{I \times J \times K}$ having the PARAFAC decomposition

$$x_{ijk} = \sum_{\lambda}^F b_{i\lambda} c_{j\lambda} s_{k\lambda} . \text{ Furthermore, let } \mathbf{X} \text{ be } CI_3. \text{ Finding the factors corresponding to the}$$

first and second dimension by SVD is then the same as finding the factors by ALSPARAFAC.

Proof:

As \mathbf{X} is CI_3 it can be written as $\mathbf{X}^{IJ \times K} = \mathbf{A}^{IJ \times F} \mathbf{S}^{F \times K}$ where \mathbf{A} and \mathbf{S} is found by ICA.

Finding \mathbf{B} and \mathbf{C} using ALSPARAFAC becomes the problem of minimizing:

$$\begin{aligned} \sum_{ijk} \left(\mathbf{X}_{ijk} - \sum_{\lambda=1}^F s_{i\lambda} b_{j\lambda} c_{k\lambda} \right)^2 &= \sum_{ijk} \left(\sum_{q=1}^F s_{kq} a_{(ij)q} - \sum_{\lambda=1}^F s_{k\lambda} b_{i\lambda} c_{j\lambda} \right)^2 = \\ \sum_k \left(\sum_{q,\lambda=1}^F s_{kq} s_{k\lambda} \sum_{ij} (a_{(ij)q} - b_{iq} c_{jq}) (a_{(ij)\lambda} - b_{i\lambda} c_{j\lambda}) \right) &= \text{tr} \left((\mathbf{A} - \mathbf{M}) \mathbf{\Psi} (\mathbf{A} - \mathbf{M})^T \right) \end{aligned} \quad \text{eq. 0.1}$$

1. Consequence of $m_{(ij)\lambda} = b_{i\lambda} c_{j\lambda}$, $\mathbf{\Psi} = \mathbf{S} \mathbf{S}^T$.

As the rows of \mathbf{S} are independent $\mathbf{\Psi} = \mathbf{S} \mathbf{S}^T$ is a diagonal matrix. Consequently, there is no interaction between the columns of $(\mathbf{A} - \mathbf{M})$. Therefore, the minimization problem can be split into minimizing the squared error of each column of $(\mathbf{A} - \mathbf{M})$. Let $q_{ij}^{(\lambda)} = a_{(ij)\lambda}$, i.e.

$\mathbf{Q} \in R^{I \times J}$ is the unmatricized version of the λ^{th} column of \mathbf{A} . The goal is now to find two vectors $\mathbf{b}_\lambda, \mathbf{c}_\lambda$ so $\left\| \mathbf{Q}^{(\lambda)} - \mathbf{b}_\lambda \mathbf{c}_\lambda^T \right\|_F$ is minimized. However, this minimization problem is solved by the singular valued decompositions, as the first vectors of the SVD decomposition explains the most of the variation of $\mathbf{Q}^{(\lambda)}$, i.e.

$$\left[\mathbf{U}, \mathbf{T}, \mathbf{V}^T \right] = \text{SVD}(\mathbf{Q}^{(\lambda)}) \Rightarrow \mathbf{b}_\lambda = \mathbf{u}_1, \mathbf{c}_\lambda = t_{11} \mathbf{v}_1 .$$

Appendix B: Multi-way array algebra

This compilation of definitions and manipulations is if not otherwise stated taken from the work of Lathauwer, Moor and Vandewalle [19].

Definition 1: The n-mode Multiplication (\times_n)

The n-mode Multiplication of the multi-way array $\mathbf{X} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_N}$ by a matrix $\mathbf{U} \in \mathcal{R}^{J_n \times I_n}$, denoted by $\mathbf{X} \times_n \mathbf{U}$, is an $(I_1 \times I_2 \times \dots \times I_{n-1} \times J_n \times I_{n+1} \times \dots \times I_N)$ -multi-way array of which the entries are given by:

$$(\mathbf{X} \times_n \mathbf{U})_{i_1 i_2 \dots i_{n-1} j_n i_{n+1} \dots i_N} \equiv \sum_{i_n} x_{i_1 i_2 \dots i_{n-1} i_n i_{n+1} \dots i_N} u_{j_n i_n}$$

Notice: let $\mathbf{V} \in \mathcal{R}^{J_m \times I_m}$ then:

$$\mathbf{X} \times_n \mathbf{U} \times_m \mathbf{V} = \mathbf{X} \times_m \mathbf{V} \times_n \mathbf{U}$$

Definition 2: The scalar product of two multi-way arrays

The scalar product $\langle \mathbf{A}, \mathbf{B} \rangle$ of two multi-way arrays $\mathbf{A}, \mathbf{B} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_N}$ is defined as:

$$\langle \mathbf{A}, \mathbf{B} \rangle \equiv \sum_{i_1} \sum_{i_2} \dots \sum_{i_n} b_{i_1 i_2 \dots i_n} a_{i_1 i_2 \dots i_n}$$

Definition 3: The norm of a multi-way array (Frobenius-norm)

The norm of a multi-way array is given by:

$$\|\mathbf{A}\| = \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle}$$

Definition 4: The rank of multi-way arrays

An N-way array \mathbf{A} has rank-1 when it equals the outer product of N vectors, i.e.

$$\mathbf{A} = \mathbf{u}^{(1)} \circ \mathbf{u}^{(2)} \circ \dots \circ \mathbf{u}^{(N)}$$

Furthermore, the rank of an arbitrary multi-way array \mathbf{A} denoted by $R = \text{rank}(\mathbf{A})$, is the minimal number of rank-1 multi-way arrays that yield \mathbf{A} in a linear combination.

Notice: The formulation for rank-1 can be restated as $a_{i_1 i_2 \dots i_N} = u_{i_1}^{(1)} u_{i_2}^{(2)} \dots u_{i_N}^{(N)}$.

Definition 5: The rank of a matrix

The normal rank of a matrix $r_{\mathbf{A}} := \text{rank}(\mathbf{A}) = r \Leftrightarrow \mathbf{A}$ contains at least a collection of r linearly independent columns, and this fails for r+1 columns. (Taken from [31])

Definition 6: The k-rank of a matrix (defined by J.B. Kruskal 1977)

The k-rank $k_{\mathbf{A}} = r \Leftrightarrow r$ columns of \mathbf{A} are linearly independent, but this fails for at least one set of $r+1$ columns (Taken from [31]). Mathematically this can be expressed as:

$k_{\mathbf{A}} = \arg \min_N (\mathbf{A} \mid \exists i; \mathbf{a}_i = \sum_{j \neq i}^N c_j \mathbf{a}_j)$, where N denotes the amount of columns of \mathbf{A} used to generate \mathbf{a}_i .

Definition 7: Diagonal multi-way arrays

A multi-way array is called diagonalizable if the core multi-way array \mathbf{S} of the HOSVD fulfills $s_{i_1 i_2 \dots i_N} = 0$ unless $i_1 = i_2 = \dots = i_N$.

Appendix C: MATLAB implementation of multi-way array manipulations

In the literature, I didn't find any fast implementations of the various multi-way array manipulations except for the first matricizing function given here. Therefore, I have suggested some very fast MATLAB implementations for these manipulations:

```
1 function X=matricizing(Y,n)
2
3 % Matricizes the multiway array Y around dimension n to give the matrix X
4 % Input:
5 % Y      Multiway array
6 % n      Dimension to use for matricizing
7 % Output:
8 % X      Matrix of matricized multiway array
9
10 N=ndims(Y);
11 X=reshape(permute(Y, [n 1:n-1 n+1:N]),size(Y,n),prod(size(Y))/size(Y,n));
12


---


1 function Y=unmatricizing(X,n,D)
2
3 % The inverse operation of matricizing, i.e. recreates the multiway
4 % array Y that was matricizes around n to give the matrix X.
5 %
6 % input:
7 % D      a vector containing the original dimensions of Y
8 % n      the dimension along which the matricizing was originally performed
9 % X      a matrix corresponding to the matricizing of Y around n
10 % output:
11 % Y      multi-way array
12
13 if n==1
14     perm=[1:length(D)];
15 else
16     perm=[2:n 1 n+1:length(D)];
17 end
18
19 Y=permute(reshape(X,D([n 1:n-1 n+1:length(D)])),perm);
20
```

```

1 function A=tmult(T,M,n)
2 % the n-mode multiplication or tensor multiplication:
3 % Let T be an I1xI2x...xIn-lxInxIn+lx...xIN multi-way array
4 % and M be an JxIn then
5 % T x_n M gives an I1xI2x...xIn-lxJxIn+lx...xIN multi-way array
6 % Input:
7 % T Multi-way array of dimensions I1xI2x...xIn-lxInxIn+lx...xIN
8 % M JxIn-matrix
9 % n the dimension to do the multiplication;
10 % Output:
11 % A Multi-way array of dimensions I1xI2x...xIn-lxJxIn+lx...xIN
12
13 - Dt=size(T);
14 - Dm=size(M);
15
16 - Tn=matricizing(T,n);
17 - Tnew=M*Tn;
18 - Dt(n)=Dm(1);
19 - A=unmatricizing(Tnew,n,Dt);
20
21


---


1 function T=outerprod(FACT)
2
3 % The outer products:
4 % T=sum_f (U1f o U2f o U3f o ...o Umf) where Uif is a vector corresponding
5 % to the f'th factor of the i'th dimension.
6 %
7 % Input:
8 % FACT Cell array containing the factor-vectors corresponding to
9 % each of the three dimensions, i.e. Ui=FACT{i}
10 % Output:
11 % T The multi-way array created from the outer product of each
12 % dimensions vector.
13
14 - T=0;
15 - for i=1:size(FACT{1},2)
16 - Y=1;
17 - for j=1:length(FACT)
18 - U=FACT{j};
19 - Y=tmult(Y,U(:,i),j);
20 - end
21 - T=T+Y;
22 - end
23
24

```

```

1  function [U,S]=HOSVD(X)
2  % The Higher Order Singular Value Decomposition as defined by
3  % Lathauwer, Lieven De., Moor, Bart De, Vandewalle, Joos
4  % "A MULTILINEAR SINGULAR VALUE DECOMPOSITION"
5  % SIAM J. MATRIX ANAL. APPL. Vol. 21, No. 4, pp. 1253 1278:
6  %
7  % Model:
8  % X=S x1 U1 x2 U2 ... xN UN, i.e. a TUCKER model restricted so
9  % Ui orthonormal being the orderen eigenvectors of the matrix
10 % found by matricizing(X,i). S the resulting core multi-way array.
11 %
12 % Input:
13 % X      Multi-way array
14 % Output:
15 % U      Cell of eigenvectors, i.e. U{i} corresponds to the eigenvectors of
16 %         the matrix M=matricizing(X,i).
17 % S      Multi-way array of same size as X
18
19 - S=X;
20 - for k=1:ndims(X)
21 -     Xk=matricizing(X,k);
22 -     [u,s]=eig(Xk*Xk'); %cheaper than the SVD since we only want U
23 -     u=fliplr(u);
24 -     U{k}=u;
25 -     S=tmult(S,U{k}',k)
26 - end
27

```

Appendix D: ICA- and ALSPARAFAC on Chemometric Data

The field in which the PARAFAC model has been the most applied is probably analyzing chemometric data. Here the dataset ‘Claus’ consisting of five samples containing different amounts of tyrosine, tryptophane and phenylalanine measured by fluorescence in the excitation spectra 240 to 300 nm and emission spectra 250 to 450 nm were analyzed. For a description of the dataset see [41]. The data set has the dimensions $sample \times emission \times excitation$. The ICAPARAFAC algorithm seems well justified as the goal is from the samples to find factors that are independent when emission and excitation is combined, i.e. $CI_{2,3}$.

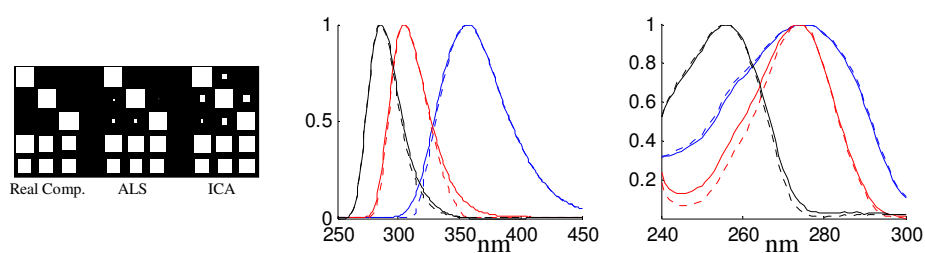


Figure 0.1: Left most panel; the true mixing of the five samples, the mixing found by ALSPARAFAC and mixing found by ICAPARAFAC. Center panel; the emission specters. Right panel; the excitation specters. ALSPARAFAC (solid) ICAPARAFAC (dashed).

From Figure 0.1 it is seen that the ICAPARAFAC algorithm is slightly worse at estimating the true mixing of the factors in the samples than the ALSPARAFAC algorithm. The ICAPARAFAC model seems however better at estimating the emission and excitation specters as they are slightly better defined. Consequently, the ICAPARAFAC method seems effective in analyzing this chemometric dataset.