

Kensuke Sekihara
Srikantan S. Nagarajan

Electromagnetic Brain Imaging

A Bayesian Perspective

Electromagnetic Brain Imaging

Kensuke Sekihara · Srikantan S. Nagarajan

Electromagnetic Brain Imaging

A Bayesian Perspective



Springer

Kensuke Sekihara
Department of Systems Design
and Engineering
Tokyo Metropolitan University
Tokyo
Japan

Srikantan S. Nagarajan
Department of Radiology and Biomedical
University of California
San Francisco, CA
USA

ISBN 978-3-319-14946-2
DOI 10.1007/978-3-319-14947-9

ISBN 978-3-319-14947-9 (eBook)

Library of Congress Control Number: 2014959853

Springer Cham Heidelberg New York Dordrecht London
© Springer International Publishing Switzerland 2015

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

Springer International Publishing AG Switzerland is part of Springer Science+Business Media
(www.springer.com)

Preface

Neuronal activities in a human brain generate coherent synaptic and intracellular currents in cortical columns, which generate electric potentials on the scalp surface and magnetic fields outside the head. Electroencephalography (EEG) measures these potentials and magnetoencephalography (MEG) measures these magnetic fields to obtain information on the state of the brain. The class of methodologies that reconstruct and visualize the neuronal activities based on MEG/EEG sensor measurements is referred to as the electromagnetic brain imaging.

In the past two decades there have been significant advances in signal processing and source reconstruction algorithms used in electromagnetic brain imaging. However, electromagnetic brain imaging can still be considered a young field. This is primarily due to the complexity associated with the electrophysiological brain activity underlying the signals. Also, it is true that applying the electromagnetic brain imaging with confidence requires an understanding of the relevant mathematics, physics, biology, and engineering as well as a broad perspective on human brain science. Due to its interdisciplinary nature, such a broad knowledge base takes years to acquire.

This book is intended to provide a coherent introduction to the body of mainstream algorithms used in electromagnetic brain imaging, with specific emphasis on novel Bayesian algorithms that we have developed. It is intended as a graduate level textbook with the goal of helping readers more easily understand the literature in biomedical engineering and in related fields, and be ready to pursue research in either the engineering or the neuroscientific aspects of electromagnetic brain imaging. We hope that this textbook will not only appeal to graduate students but all scientists and engineers engaged in research on electromagnetic brain imaging.

This book begins with an introductory overview of electromagnetic brain imaging in Chap. 1 and then discusses dominant algorithms that are used in electromagnetic brain imaging in Chaps. 2–4. Minimum-norm-based methods, which are classic algorithms and still widely used in this field, are described in Chap. 2, but with a Bayesian perspective in mind. Chapter 3 presents a concise review of adaptive beamformers, which have become a standard tool for analyzing brain

spontaneous activity such as resting-state MEG data, and we also include a Bayesian perspective on adaptive beamformers.

Chapters 4–6 review Bayesian algorithms for electromagnetic brain imaging that we have developed in the past decade. Chapter 4 presents a novel Bayesian algorithm, called Champagne, which has been developed by our group. Since we believe that the Champagne algorithm is a powerful next-generation imaging algorithm, the derivation of the algorithm is presented in detail. Chapter 5 presents Bayesian factor analysis, which is a group of versatile algorithms used for denoising, interference suppression, and source localization. Chapter 6 describes a unified theoretical Bayesian framework, which provides insightful perspective into various source imaging methods and reveals similarities and equivalences between methods that appear to be very different.

Chapters 7–9 deal with newer topics that are currently in vogue in electromagnetic brain imaging—functional connectivity, causality, and cross-frequency coupling analyses. Chapter 7 reviews functional connectivity analysis using imaginary coherence. Chapter 8 provides a review of several directional measures that can detect causal coupling of brain activities. Chapter 9 presents novel empirical results showing that the phase-amplitude coupling can be detected using MEG source-space analysis, and demonstrates that the electromagnetic brain imaging holds great potential in elucidating the mechanisms of brain information processing. This chapter was contributed by Eiichi Okumura and Takashi Asakawa. The first two chapters in the Appendix provide concise explanations of bioelectromagnetic forward modeling and basics of Bayesian inference, and the third chapter provides supplementary mathematical arguments. These chapters are included for the reader’s convenience.

Many people have made valuable contributions together with our own efforts in this area. Special mention must be made of Hagai Attias who is an invaluable collaborator. Hagai introduced us to probabilistic graphical models and Bayesian inference methods. He was an integral person in our fruitful collaboration, on which this book is based.

Many students, postdocs, fellows, and UCSF faculty members have collaborated with us over the years. These include Leighton Hinkley, David Wipf, Johanna Zumer, Julia Owen, Sarang Dalal, Isamu Kumihashi, Alex Herman, Naomi Kort, Ethan Brown, Rodney Gabriel, Maneesh Sahani, Adrian Guggisberg, Juan Martino, Kenneth Hild, Matthew Brookes, Carsten Stahlhut, Oleg Portniaguine, Erik Edwards, Ryan Canolty, David McGonigle, Tal Kenet, Theda Heinks-Maldonado, Aliu Sheye, Dameon Harrell, Josiah Ambrose, Sandeep Manyam, William McClay, Virginie van Wassenhove, Ilana Hairston, Maria Ventura, Zhao Zhu, Corby Dale, Tracy Luks, Kelly Westlake, Kamalini Ranasinghe, Karuna Subramanium, Carrie Niziollek, Zarinah Agnew, Carly Demopoulos, Megan Thomson, Peter Lin, Tulaya Limpiti, Lisa Dyson, Pew Puthividya, Jiucang Hao, Phiroz Tarapore, Dario Englot, Heidi Kirsch, Sophia Vinogradov, Michael Merzenich, Christoph Schreiner, Elizabeth Disbrow, Mitchel Berger, Edward Chang, Nick Barbaro, Roland Henry, Sarah Nelson, William Dillon, Jim Barkovich, Nancy Byl, David Blake, Keith Vossel, Elliot Sherr, Elysa Marco, Josh Wooley, Marilu Gorno-Tempini, Robert Knight, and

Pratik Mukherjee. Also, we are deeply indebted to the following fantastic and incredibly dedicated staff at the Biomagnetic Imaging Laboratory: Susanne Honma, Mary Mantle, Anne Findlay-Dowling, Danielle Mizuiri, and Garrett Coleman.

Furthermore, several people have been particularly influential, and have greatly helped us to enhance our understanding of this field. Such people include Thomas Ferree, Mike X. Cohen, Stephen E. Robinson, Jiri Vrba, Guido Nolte, and Barry Van Veen. We would also like to thank Daniel Palomo for his effort in editing the manuscript of this book. Finally, we thank the following collaborators for their friendship and invaluable support to our work over nearly 20 years: John Houde, Steven Cheung, David Poepel, Alec Marantz, and Tim Roberts.

We kindly ask readers to visit www.electromagneticbrainimaging.info. Supplementary information, as well as error corrections (if necessary), is uploaded to this website.

November 2014

Kensuke Sekihara
Srikantan S. Nagarajan

Contents

1	Introduction to Electromagnetic Brain Imaging	1
1.1	Functional Brain Imaging and Bioelectromagnetic Measurements	1
1.2	Sensing Magnetic Fields from the Brain	2
1.3	Electromagnetic Brain Imaging.	3
1.3.1	Forward Model.	3
1.3.2	Inverse Algorithms	4
1.4	From Source Imaging to Functional Connectivity Imaging	6
1.5	Examples of Clinical Applications	6
1.5.1	Functional Mapping for Preoperative Neurosurgical Planning	6
1.5.2	Functional Connectivity Imaging	7
References.		8
2	Minimum-Norm-Based Source Imaging Algorithms	9
2.1	Introduction	9
2.2	Definitions.	9
2.3	Sensor Lead Field.	10
2.4	Voxel Source Model and Tomographic Source Reconstruction	11
2.5	Maximum Likelihood Principle and the Least-Squares Method	13
2.6	Derivation of the Minimum-Norm Solution	14
2.7	Properties of the Minimum-Norm Solution.	15
2.8	L_2 -Regularized Minimum-Norm Solution.	17
2.9	L_1 -Regularized Minimum-Norm Solution.	19
2.9.1	L_1 -Norm Constraint.	19
2.9.2	Intuitive Explanation for Sparsity	20
2.9.3	Problem with Source Orientation Estimation.	23

2.10	Bayesian Derivation of the Minimum-Norm Method	24
2.10.1	Prior Probability Distribution and Cost Function	24
2.10.2	L_2 -Regularized Method	24
2.10.3	L_1 -Regularized Method	26
	References	28
3	Adaptive Beamformers	29
3.1	Introduction and Basic Formulation	29
3.2	Classical Derivation of Adaptive Beamformers	30
3.2.1	Minimum-Variance Beamformers with Unit-Gain Constraint	30
3.2.2	Minimum-Variance Beamformer with Array-Gain Constraint	31
3.2.3	Minimum-Variance Beamformer with Unit-Noise-Gain Constraint	32
3.3	Semi-Bayesian Derivation of Adaptive Beamformers	33
3.4	Diagonal-Loading and Bayesian Beamformers	35
3.5	Scalar Adaptive Beamformer with Unknown Source Orientation	36
3.5.1	Expressions for the Unit-Gain Constraint Beamformer	36
3.5.2	Expressions for the Array-Gain and Weight-Normalized Beamformers	37
3.6	Vector-Type Adaptive Beamformer	38
3.6.1	Vector Beamformer Formulation	38
3.6.2	Semi-Bayesian Formulation	40
3.7	Narrow-Band Beamformer	42
3.7.1	Background	42
3.7.2	Time-Domain Implementation	42
3.7.3	Frequency-Domain Implementation	43
3.7.4	Five-Dimensional Brain Imaging	44
3.8	Nonadaptive Spatial Filters	44
3.8.1	Minimum-Norm Filter	44
3.8.2	Weight-Normalized Minimum-Norm Filter	46
3.8.3	sLORETA Filter	46
3.9	Recursive Null-Steering (RENS) Beamformer	47
3.9.1	Beamformer Obtained Based on Beam-Response Optimization	47
3.9.2	Derivation of RENS Beamformer	48
	References	49
4	Sparse Bayesian (Champagne) Algorithm	51
4.1	Introduction	51
4.2	Probabilistic Model and Method Formulation	52

4.3	Cost Function for Marginal Likelihood Maximization	54
4.4	Update Equations for α	56
4.5	Modified Algorithm Integrating Interference Suppression	58
4.6	Convexity-Based Algorithm	59
4.6.1	Deriving an Alternative Cost Function	59
4.6.2	Update Equation for z	61
4.6.3	Update Equation for x_k	61
4.6.4	Update Equation for v	62
4.6.5	Summary of the Convexity-Based Algorithm	62
4.7	The Origin of the Sparsity	63
4.8	Extension to Include Source Vector Estimation	65
4.8.1	Update Equation for Z_j	66
4.8.2	Update Equation for $s(t_k)$	67
4.8.3	Update Equation for T_j	68
4.9	Source Vector Estimation Using Hyperparameter Tying	69
4.10	Appendix to This Chapter	71
4.10.1	Derivation of Eq. (4.21)	71
4.10.2	Derivation of Eq. (4.29)	72
4.10.3	Proof of Eq. (4.50)	73
	References	74
5	Bayesian Factor Analysis: A Versatile Framework for Denoising, Interference Suppression, and Source Localization	75
5.1	Introduction	75
5.2	Bayesian Factor Analysis	75
5.2.1	Factor Analysis Model	75
5.2.2	Probability Model	76
5.2.3	EM Algorithm	77
5.2.4	Computation of Marginal Likelihood	79
5.2.5	Summary of the BFA Algorithm	81
5.3	Variational Bayes Factor Analysis (VBFA)	82
5.3.1	Prior Distribution for Mixing Matrix	82
5.3.2	Variational Bayes EM Algorithm (VBEM)	84
5.3.3	Computation of Free Energy	92
5.3.4	Summary of the VBFA Algorithm	95
5.4	Partitioned Factor Analysis (PFA)	96
5.4.1	Factor Analysis Model	96
5.4.2	Probability Model	97
5.4.3	VBEM Algorithm for PFA	97
5.4.4	Summary of the PFA Algorithm	100
5.5	Saketini: Source Localization Algorithm Based on the VBFA Model	101
5.5.1	Data Model	101

5.5.2	Probability Model	102
5.5.3	VBEM Algorithm	103
5.5.4	Summary of the Saketini Algorithm	107
5.6	Numerical Examples	107
5.7	Appendix to This Chapter	112
5.7.1	Proof of Eq. (5.84)	112
5.7.2	Proof of Eq. (5.94)	114
5.7.3	Proof of Eq. (5.103)	115
5.7.4	Proof of Eq. (5.166)	115
	References	117
6	A Unified Bayesian Framework for MEG/EEG Source Imaging	119
6.1	Introduction	119
6.2	Bayesian Modeling Framework	121
6.3	Bayesian Modeling Using General Gaussian Scale Mixtures and Arbitrary Covariance Components	122
6.3.1	The Generative Model	122
6.3.2	Estimation and Inference	123
6.3.3	Source MAP or Penalized Likelihood Methods	127
6.3.4	Variational Bayesian Approximation	131
6.4	Selection of Covariance Components C	133
6.5	Discussion	134
	References	137
7	Source-Space Connectivity Analysis Using Imaginary Coherence	139
7.1	Introduction	139
7.2	Source-Space Coherence Imaging	140
7.3	Real and Imaginary Parts of Coherence	141
7.4	Effects of the Leakage	143
7.4.1	Leakage Effects on the Magnitude Coherence	143
7.4.2	Leakage Effects on the Imaginary Coherence	144
7.5	Corrected Imaginary Coherence	145
7.5.1	Modification of Imaginary Coherence	145
7.5.2	Factorization of Mutual Information	146
7.5.3	Residual Coherence	148
7.5.4	Phase Dependence of the Corrected Imaginary Coherences	150
7.6	Canonical Coherence	151
7.6.1	Canonical Magnitude Coherence	151
7.6.2	Canonical Imaginary Coherence	154
7.6.3	Canonical Residual Coherence	157

7.6.4	Computing Coherence When Each Voxel has Multiple Time Courses	158
7.7	Envelope Correlation and Related Connectivity Metrics	159
7.7.1	Envelope Correlation	159
7.7.2	Residual Envelope Correlation	160
7.7.3	Envelope Coherence	160
7.8	Statistical Thresholding of Coherence Images	161
7.9	Mean Imaginary Coherence (MIC) Mapping	162
7.10	Numerical Examples	163
	References	168
8	Estimation of Causal Networks: Source-Space Causality Analysis	171
8.1	Introduction	171
8.2	Multivariate Vector Autoregressive (MVAR) Process	171
8.2.1	MVAR Modeling of Time Series	171
8.2.2	Coherence and Partial Coherence of the MVAR Process	173
8.3	Time-Domain Granger Causality	174
8.3.1	Granger Causality for a Bivariate Process	174
8.3.2	Multivariate Granger Causality	175
8.3.3	Total Interdependence	177
8.4	Spectral Granger Causality: Geweke Measures	178
8.4.1	Basic Relationships in the Frequency Domain	178
8.4.2	Total Interdependence and Coherence	179
8.4.3	Deriving Causal Relationships in the Frequency Domain	180
8.5	Other MVAR-Modeling-Based Measures	182
8.5.1	Directed Transfer Function (DTF)	182
8.5.2	Relationship Between DTF and Coherence	183
8.5.3	Partial Directed Coherence (PDC)	184
8.6	Transfer Entropy	185
8.6.1	Definition	185
8.6.2	Transfer Entropy Under Gaussianity Assumption	186
8.6.3	Equivalence Between Transfer Entropy and Granger Causality	187
8.6.4	Computation of Transfer Entropy	188
8.7	Estimation of MVAR Coefficients	190
8.7.1	Least-Squares Algorithm	190
8.7.2	Sparse Bayesian (Champagne) Algorithm	191
8.8	Numerical Examples	192
8.8.1	Experiments Using Bivariate Causal Time Series	192
8.8.2	Experiments Using Trivariate Causal Time Series	194
	References	198

9 Detection of Phase–Amplitude Coupling in MEG Source Space:	
An Empirical Study	199
9.1 Introduction	199
9.2 Types of Cross-Frequency Coupling	200
9.3 Local PAC and Cross-Location PAC	201
9.4 Quantification of Phase–Amplitude Coupling	202
9.4.1 Instantaneous Amplitude and Phase.	202
9.4.2 Amplitude–Phase Diagram	202
9.4.3 Modulation Index (MI)	203
9.4.4 Phase-Informed Time-Frequency Map	204
9.5 Source Space PAC Analysis: An Example Study	
Using Hand-Motor MEG Data	204
9.5.1 Experimental Design and Recordings	204
9.5.2 Data Analysis.	204
9.5.3 Results of Local PAC Analysis.	206
9.5.4 Results of Analyzing Cross-Location PACs	209
9.6 Summary	212
References	212
Appendix A: Bioelectromagnetic Forward Modeling	215
Appendix B: Basics of Bayesian Inference	231
Appendix C: Supplementary Mathematical Arguments	247
Index	267

Chapter 1

Introduction to Electromagnetic Brain Imaging

1.1 Functional Brain Imaging and Bioelectromagnetic Measurements

Noninvasive functional brain imaging has made a tremendous impact in improving our understanding of the human brain. Functional magnetic resonance imaging (fMRI) has been the predominant modality for imaging the functioning brain since the middle of the 1990s. fMRI measures changes in blood oxygenation-level-dependent (BOLD) signals caused by neuronal activation. It is a noninvasive method that allows for whole-brain measurement and the examination of activity in deep brain structures.

However, fMRI lacks the temporal resolution required to image the dynamic and oscillatory spatiotemporal patterns associated with activities in a brain. This is because the BOLD signal, which is only an indirect measure of neural activity, is fundamentally limited by the rate of oxygen consumption and subsequent blood flow. Furthermore, since the BOLD signal is only an approximate and indirect measure of neuronal activity, it might not accurately reflect true neuronal processes. Hence, to observe neurophysiological processes more directly within relevant timescales, imaging techniques that have both high temporal and adequate spatial resolution are needed.

Neurons in the brain function electrically, as well as chemically. Therefore, their activity generates associated electric and magnetic fields that can be detected outside the head. Electromagnetic brain imaging is a term intended to encompass noninvasive techniques which probe brain electromagnetic activity and properties. Two techniques currently exist for detecting electrical brain activities: electroencephalography (EEG) and magnetoencephalography (MEG). EEG measures the electric potential by means of electrodes placed on the scalp, and MEG measures magnetic fields by means of sensors placed near the head.

In contrast to fMRI, both MEG and EEG directly measure electromagnetic fields emanating from the brain with excellent temporal resolution of less than 1 ms, and allow the study of neural oscillatory processes over a wide frequency range (1–600 Hz). Because of such high temporal resolution, MEG and EEG can provide

information not obtainable with other functional brain imaging techniques. Most notably, MEG and EEG track neural population activity on millisecond timescales, revealing large-scale dynamics that are crucial for understanding brain function. Furthermore, by applying modern inverse algorithms, it is possible to obtain three-dimensional images, which are reasonable estimates of neural activity. Such images are extremely useful for answering many questions on brain science.

While MEG is mainly sensitive to tangential currents in the brain closer to the surface and relatively insensitive to the conductive properties of the skull, EEG is primarily sensitive to radial sources while being highly sensitive to the conductive properties of the brain, skull, and scalp. Therefore, MEG and EEG can be viewed as being complementary in terms of the sensitivity to underlying neural activity. However, since magnetic fields generated from neurons are not distorted by the heterogeneous electrical properties of a brain, these magnetic fields can be considered an undistorted signature of underlying cortical activity.

In addition, there are several practical or physiological reasons why neuroscientists prefer MEG to EEG. First, MEG setup time is very short, because MEG measurement does not require much preparation in attaching and checking electrodes, as is needed in performing EEG measurement. This simplifies matters both for experimenters and subjects. Second, the anatomical location of primary sensory cortices in sulci makes MEG ideally suited for electrophysiological studies. Furthermore, with whole-head sensor arrays, MEG is also well-suited to investigate hemispheric lateralization effects. Therefore, this chapter is primarily dedicated to giving a review of the methodologies associated with MEG.

1.2 Sensing Magnetic Fields from the Brain

The long apical dendrites of cortical pyramidal cells are arranged perpendicularly to the cortical surface and parallel to each other. This fortuitous anatomical arrangement of these cells allows the magnetic fields to sum up to magnitudes large enough to detect at the scalp. Synchronously fluctuating dendritic currents result in equivalent current dipoles that produce such magnetic fields. However, biomagnetic fields from a brain are extremely small, (in range of tens-to-hundreds of femto-Tesla (fT)) which is about seven orders of magnitude smaller than the earth's magnetic field. As a result, appropriate data collection necessitates a magnetically shielded room and highly sensitive detectors known as superconducting quantum interference devices (SQUIDs). Biomagnetic fields from a brain are typically sensed using detection coils called flux transformers or magnetometers, which are positioned close to the scalp and connected to SQUIDs. A SQUID acts as a magnetic-field-to-voltage converter, and its nonlinear response is linearized by flux-locked loop electronics. SQUIDs have a sensitivity of up to 5 femto-Tesla per square root of Hz, which is adequate for the detection of brain-generated magnetic fields.

MEG sensors are often configured for measuring differential magnetic fields so as to reduce ambient noise in measurements. Such sensors are referred to as

gradiometers, although some MEG systems do not use gradiometers relying on clever noise cancellation methods. The two commonly used gradiometer configurations are axial and planar gradiometers. Axial (first order) gradiometers consist of two coils that share the same axis, whereas planar (first order) gradiometers consist of two coils that share the same plane. The sensitivity profile of a planar gradiometer is somewhat similar to EEG, whereby a sensor is maximally sensitive to a source closest to it. In contrast, the sensitivity profile of an axial gradiometer can be somewhat counterintuitive because it is not maximally sensitive to sources closest to the sensors.

Modern MEG systems consist of simultaneous recordings from many sensors that provide whole head coverage. The total number of sensors varies from 100 to 300. The advent of such large sensor-array systems has significantly advanced MEG studies. Although the maximum sampling rate for many MEG systems reaches more than 10 kHz, most MEG data is usually recorded at a sampling rate of around 1,000 Hz, which still provides excellent temporal resolution for measuring the dynamics of cortical neuronal activity at millisecond order.

1.3 Electromagnetic Brain Imaging

MEG sensor data only provides qualitative information about underlying brain activities. The analysis of the sensor data is typically performed based on the intuitions of experienced users regarding the sensitivity profile of the sensors. To extract more precise information from the observed sensor data, it is essential to apply imaging-type analysis involving the reconstruction of brain activities from the sensor data. Major components for the electromagnetic brain imaging are the forward model and the inverse algorithms.

1.3.1 Forward Model

The forward model consists of three subcomponents: the source model, the volume conductor model, and the measurement model. Typical source models assume that the brain magnetic fields are generated by equivalent current dipoles in the brain. This model is consistent with available measurements of coherent synaptic and intracellular currents in cortical columns that are thought to be major contributors to MEG and EEG signals. Although several more complex source models have been proposed, the equivalent current dipole is the dominant source model. This is because, given the distance between the sources and sensors, the dipole is a reasonable approximation of brain sources.

The volume conductor model refers to the equations that govern the relation between the source model and the sensor measurements, namely the electric potentials or the magnetic fields. These surface integral equations, obtained by solving Maxwell's equations under quasi-static conditions, can be solved analytically for

special geometries of the volume conductor, such as the sphere or an ellipsoid. A concise review on the spherical homogeneous conductor model is found in Appendix A. For realistic volume conductors, various numerical techniques such as finite-element and boundary-element methods may be employed, although these methods are generally time-consuming.

Measurement models refer to the specific measurement systems used in EEG and MEG including the position of the sensors relative to the head. For instance, different MEG systems measure axial versus planar gradients of the magnetic fields with respect to different locations of sensors. The measurement model incorporates such information about the type of measurement and the geometry of the sensors. Measurement of the position of the head relative to the sensor array is accomplished by attaching head-localization coils to fiducial landmarks on the scalp. Modern MEG systems are sometimes equipped with continuous head-localization procedures that enable constant updating of sensor locations relative to the head to compensate for subjects head movements. The source, volume conductor, and measurement models are typically combined into a concept called the lead field that describes a linear relationship between sources and the measurements. When discussing inverse methods, we assume that the lead field matrix is known.

1.3.2 Inverse Algorithms

Inverse algorithms are used for solving the bioelectromagnetic inverse problem, i.e., for estimating the parameters of neural sources from MEG and EEG sensor data. When implementing electromagnetic brain imaging, this estimation of spatial locations and timing of brain sources is a challenging problem because it involves solving for unknown brain activity across thousands of voxels from the recordings of just a few hundred sensors. In general, there are no unique solutions to the inverse problem because there are many source configurations that could produce sensor data equal to the sensor observations, even in the absence of noise and (if given) infinite spatial or temporal sampling. This nonuniqueness is referred to as the ill-posed nature of the inverse problem. Nevertheless, to get around this nonuniqueness, various estimation procedures incorporate prior knowledge and constraints about source characteristics.

Inverse algorithms can be classified into two categories: model-based dipole fitting and (non-model-based) imaging methods. Dipole fitting methods assume that a small set of current dipoles can adequately represent an unknown source distribution. In this case, the dipole locations and its moments form a set of unknown parameters, which are typically estimated using the least-squares fit. The dipole fitting method—particularly the single-dipole fitting method—has clinically been used for localization of early sensory responses in somatosensory and auditory cortices.

However, two major problems exist in a dipole fitting procedure. First, the nonlinear optimization causes a problem of local minima when more than two dipole parameters are estimated. A second, more difficult problem is that the dipole fitting methods require *a priori* knowledge of the number of dipoles. Often, such information

about the model order is not available a priori, especially for complex brain mapping conditions, and the resulting localization of higher order cortical functions can sometimes be unreliable.

An alternative approach is whole-brain source imaging methods. These methods apply voxel discretization over a whole brain volume, and assume a fixed source at each voxel. These methods estimate the amplitudes (and directions) of the sources by minimizing a cost function. One classic method of this kind is the minimum-norm method [1]. The minimum-norm and related methods are the topic of Chap. 2.

Since the number of voxels is usually much larger than the number of sensors, the cost function should contain a constraint term that is derived from various prior information about the nature of the sources. We have developed a powerful imaging algorithm that incorporates a sparsity constraint that facilitates the sparse source configurations. This algorithm, called the Champagne algorithm, is described in detail in Chap. 4.

The other class of imaging algorithm is the spatial filter. The spatial filter estimates source amplitude at each location independently. It is often called a virtual sensor method, because it forms a virtual sensor, which scans an entire source space to produce a source image. The most popular spatial filter algorithms are adaptive spatial filtering techniques, more commonly referred to as adaptive beamformers. Adaptive beamformers are the topic of Chap. 3 in which a concise review on adaptive beamformers is presented. A comprehensive review of these algorithms is found in [2]. The implementation of adaptive beamformers is quite simple, and they have proven to be powerful algorithms for characterizing cortical oscillations. Therefore, they are popular for the source localization from spontaneous brain activity, particularly resting-state MEG. We have recently proposed novel scanning algorithms [3, 4]. These methods have shown performance improvements over the adaptive beamformers. One of these methods, called the Saketini algorithm, is described in Chap. 5.

In Chap. 6, we discuss a hierarchical Bayesian framework which encompasses various source imaging algorithms in a unified framework and reveals a close relationship between algorithms that are considered quite different [5].

An enduring problem in MEG imaging is that the brain evoked responses to sensory or cognitive events are small compared to the interfering magnetic field. Typical sources of such interference include the background room interference from power lines and electronic equipment, the interference with biological origins such as heartbeat, eye blink, or other muscle artifacts. Ongoing brain activity itself, including the drowsy-state alpha rhythm, is also a major source of interference. All existing methods for brain source imaging are hampered by such interferences present in MEG data. Several signal-processing methods have been developed to reduce interferences by preprocessing the sensor data before submitting it to source localization algorithms. One such algorithm, called partitioned factor analysis, is described in Chap. 5.

1.4 From Source Imaging to Functional Connectivity Imaging

It is now well recognized in neuroscience that it is necessary to examine how the brain integrates information across multiple regions. The term functional connectivity essentially defines the complex functional interaction between separate brain areas. Although functional connectivity analysis using fMRI is common, MEG is better suited for modeling and detecting such interactions, because of its high temporal resolution.

Functional connectivity analysis can be applied either in sensor-space or in source-space. In sensor-space analysis, the field spread across many sensors arriving from a single brain region leads to uncertainties in interpreting the estimation results of brain interactions [6]. Therefore, a number of studies have begun to use source-space analysis, in which voxel time courses are first estimated by an inverse algorithm and brain interactions are then analyzed using those estimated voxel time courses. However, a serious problem arises from the leakage of an inverse algorithm, and such leakages are more or less inevitable in any inverse algorithm. Chapter 7 describes the imaginary coherence analysis in source space. The imaginary coherence analysis is known to be robust to the leakage of an inverse algorithm, and has become a popular method in analyzing functional connectivity using MEG.

Connectivity measures commonly used these days—such as coherence—are bidirectional measures, which cannot detect the directionality of brain interactions. There has been growing interest in analyzing causal networks in the brain, and directional measures are needed to detect such causal networks. Measures for detecting causal networks are the topic of Chap. 8.

Neural mechanisms of brain information processing are the subject of intense investigation. A prevailing hypothesis is that the brain uses temporal encoding based on firing phase, and that there may exist a coupling between oscillators of different frequencies. Chapter 9 presents empirical results that such phase-amplitude coupling can be detected using MEG source-space analysis. The results in Chap. 9 demonstrate that MEG functional connectivity imaging holds great potential in revealing mechanisms of brain information processing.

1.5 Examples of Clinical Applications

1.5.1 *Functional Mapping for Preoperative Neurosurgical Planning*

The surgical management of brain tumors or epilepsy often requires detailed functional mapping of cortical regions around a tumor or epileptogenic zone. Preoperative mapping techniques aim at accurately estimating the location of functional

areas in relation to a tumor or an epileptic focus to minimize surgical risk. Mass lesions can frequently distort normal neuroanatomy, which makes the identification of eloquent cortices inaccurate with normal neuroanatomical landmarks. MEG is increasingly being used for preoperative functional brain imaging. By mapping relevant somatosensory, auditory, and motor cortices preoperatively, retained areas of function can be delineated. Preoperative localization of functionally intact brain tissue helps guide neurosurgical planning and limits the region of resection, allowing for improved long-term patient morbidity and neurological function.

Examples that show how MEG imaging is used to map the motor cortex in presurgical patients are given in [7]. This study performed localization of β -band desynchronization preceding the index finger flexion for a subject with a frontal tumor, and showed the location of hand motor cortex relative to a single dipole localization of hand somatosensory cortex. These results were confirmed with electrical cortical stimulation. This investigation demonstrated that MEG source images obtained using beta-band event-related desynchronization reliably localize the hand motor cortex.

1.5.2 Functional Connectivity Imaging

Functional connectivity analysis has been shown to have profound clinical significance, because disturbances in networks are manifested as abnormalities in functional connectivity and such abnormalities can be detected using resting state MEG. Recent studies have shown this to be the case, and abnormalities in functional connectivity during resting state are observed in many clinical conditions such as brain tumors [8], strokes [9], traumatic brain injury [10], schizophrenia [11, 12], and Alzheimer disease [13].

Utilizing MEG imaging, changes in the alpha-band functional connectivity in patients with traumatic brain injury were measured and compared to healthy controls in [10]. In this investigation, mean imaginary coherence (MIC) (described in Sect. 7.9) at each brain voxel was calculated as an index for functional connectivity. In one male patient's case, his initial MEG scan was obtained 9 months after his injury and the results of MIC mapping demonstrated several regions of decreased resting state functional connectivity compared to the control group. The follow-up second MEG scan was obtained 23 months after the initial scan and the results of MIC mapping demonstrated a decrease in the volume of cortex with reduced connectivity, although some of the cortical regions with the greatest reductions in functional connectivity seen in the initial MEG remained abnormal even in the second MEG scan.

The investigation described next identifies brain regions that exhibited abnormal resting-state connectivity in patients with schizophrenia [11, 12]. Associations between functional connectivity and clinical symptoms were found in stable outpatient participants. Resting-state MEG was measured from thirty schizophrenia patients and fifteen healthy control subjects. The MIC mapping was computed in the alpha

frequency band, and the results showed that the functional connectivity of the left inferior parietal cortex was negatively related to positive symptoms and the left prefrontal cortical connectivity was associated with negative symptoms. This study demonstrates direct functional disconnection in schizophrenia patients between specific cortical fields within low-frequency resting state oscillations. Such findings indicate that this level of functional disconnection between cortical regions is an important treatment target in schizophrenia patients.

References

1. M.S. Hämäläinen, R.J. Ilmoniemi, Interpreting measured magnetic fields of the brain: estimates of current distributions. Technical Report TKK-F-A559, Helsinki University of Technology (1984)
2. K. Sekihara, S.S. Nagarajan, *Adaptive Spatial Filters for Electromagnetic Brain Imaging* (Springer, Berlin, 2008)
3. J.M. Zumer, H.T. Attias, K. Sekihara, S.S. Nagarajan, A probabilistic algorithm integrating source localization and noise suppression for MEG and EEG data. *NeuroImage* **37**, 102–115 (2007)
4. J.M. Zumer, H.T. Attias, K. Sekihara, S.S. Nagarajan, Probabilistic algorithms for MEG/EEG source reconstruction using temporal basis functions learned from data. *NeuroImage* **41**(3), 924–940 (2008)
5. D. Wipf, S.S. Nagarajan, A unified Bayesian framework for MEG/EEG source imaging. *NeuroImage* **44**, 947–966 (2009)
6. J.-M. Schoffelen, J. Gross, Source connectivity analysis with MEG and EEG. *Hum. Brain Mapp.* **30**, 1857–1865 (2009)
7. P.E. Tarapore, M.C. Tate, A.M. Findlay, S.M. Honma, D. Mizuiri, M.S. Berger, S.S. Nagarajan, Preoperative multimodal motor mapping: a comparison of magnetoencephalography imaging, navigated transcranial magnetic stimulation, and direct cortical stimulation: clinical article. *J. Neurosurg.* **117**(2), 354–362 (2012)
8. A.G. Guggisberg, S.M. Honma, A.M. Findlay, S.S. Dalal, H.E. Kirsch, M.S. Berger, S.S. Nagarajan, Mapping functional connectivity in patients with brain lesions. *Ann. Neurol.* **63**(2), 193–203 (2008)
9. K.P. Westlake, L.B. Hinkley, M. Bucci, A.G. Guggisberg, A.M. Findlay, R.G. Henry, S.S. Nagarajan, N. Byl, Resting state alpha-band functional connectivity and recovery after stroke. *Exp. Neurol.* **237**(1), 160–169 (2012)
10. P.E. Tarapore, A.M. Findlay, S.C. LaHue, H. Lee, S.M. Honma, D. Mizuiri, T.L. Luks, G.T. Manley, S.S. Nagarajan, P. Mukherjee, Resting state magnetoencephalography functional connectivity in traumatic brain injury: clinical article. *J. Neurosurg.* **118**(6), 1306–1316 (2013)
11. L.B. Hinkley, J.P. Owen, M. Fisher, A.M. Findlay, S. Vinogradov, S.S. Nagarajan, Cognitive impairments in schizophrenia as assessed through activation and connectivity measures of magnetoencephalography (MEG) data. *Front. Hum. Neurosci.* **3**, 73 (2009)
12. L.B. Hinkley, S. Vinogradov, A.G. Guggisberg, M. Fisher, A.M. Findlay, S.S. Nagarajan, Clinical symptoms and alpha band resting-state functional connectivity imaging in patients with schizophrenia: implications for novel approaches to treatment. *Biol. Psychiatry* **70**(12), 1134–1142 (2011)
13. K.G. Ranasinghe, L.B. Hinkley, A.J. Beagle, D. Mizuiri, A.F. Dowling, S.M. Honma, M.M. Finucane, C. Scherling, B.L. Miller, S.S. Nagarajan et al., Regional functional connectivity predicts distinct cognitive impairments in Alzheimers disease spectrum. *NeuroImage: Clin.* **5**, 385–395 (2014)

Chapter 2

Minimum-Norm-Based Source Imaging Algorithms

2.1 Introduction

In this chapter, we describe the minimum-norm and related methods, which are classic algorithms for electromagnetic brain imaging [1, 2]. In this chapter, the minimum-norm method is first formulated based on the maximum-likelihood principle, and the properties of the minimum-norm solution are discussed. This discussion leads to the necessity of regularization when implementing the minimum-norm method. We discuss two different representative regularization methods: the L_2 -norm regularization and the L_1 -norm regularization. The minimum-norm method is, then, formulated based on Bayesian inference—Bayesian formulation providing a form of the minimum-norm method where the regularization is already embedded.

2.2 Definitions

In electromagnetic brain imaging, we use an array of sensors to obtain bioelectromagnetic measurements. We define the output of the m th sensor at time t as $y_m(t)$, and the column vector containing outputs from all sensors, such that

$$\mathbf{y}(t) = \begin{bmatrix} y_1(t) \\ y_2(t) \\ \vdots \\ y_M(t) \end{bmatrix}, \quad (2.1)$$

where M is the total number of sensors. This column vector $\mathbf{y}(t)$ expresses the outputs of the sensor array, and it may be called the data vector or array measurement.

A spatial location is represented by a three-dimensional vector \mathbf{r} : $\mathbf{r} = (x, y, z)$. A source vector is defined as a three-dimensional column vector $\mathbf{s}(\mathbf{r}, t)$:

$$\mathbf{s}(\mathbf{r}, t) = \begin{bmatrix} s_x(\mathbf{r}, t) \\ s_y(\mathbf{r}, t) \\ s_z(\mathbf{r}, t) \end{bmatrix}, \quad (2.2)$$

where $s_x(\mathbf{r}, t)$, $s_y(\mathbf{r}, t)$, and $s_z(\mathbf{r}, t)$ are the x , y , and z components. The physical nature of the source vector is the electromotive force generated by neuronal activities in the brain. Additional discussion regarding the nature of the sources is presented in Sect. A.1 in the Appendix. The magnitude of the source vector is denoted as a scalar $s(\mathbf{r}, t)$, and the orientation of the source is denoted as a three-dimensional unit vector $\boldsymbol{\eta}(\mathbf{r}) = [\eta_x(\mathbf{r}), \eta_y(\mathbf{r}), \eta_z(\mathbf{r})]^T$, where the superscript T indicates the matrix transpose. Then, the relationship

$$\mathbf{s}(\mathbf{r}, t) = s(\mathbf{r}, t)\boldsymbol{\eta}(\mathbf{r}) = s(\mathbf{r}, t) \begin{bmatrix} \eta_x(\mathbf{r}) \\ \eta_y(\mathbf{r}) \\ \eta_z(\mathbf{r}) \end{bmatrix} \quad (2.3)$$

holds.

2.3 Sensor Lead Field

We assume that a unit-magnitude source exists at \mathbf{r} . We denote the output of the m th sensor due to this unit-magnitude source as $l_m^x(\mathbf{r})$, $l_m^y(\mathbf{r})$, and $l_m^z(\mathbf{r})$ when the unit-magnitude source is directed in the x , y , and z directions, respectively. The column vectors $\mathbf{l}_x(\mathbf{r})$, $\mathbf{l}_y(\mathbf{r})$, and $\mathbf{l}_z(\mathbf{r})$ are defined as

$$\begin{aligned} \mathbf{l}_x(\mathbf{r}) &= [l_1^x(\mathbf{r}), l_2^x(\mathbf{r}), \dots, l_M^x(\mathbf{r})]^T, \\ \mathbf{l}_y(\mathbf{r}) &= [l_1^y(\mathbf{r}), l_2^y(\mathbf{r}), \dots, l_M^y(\mathbf{r})]^T, \\ \mathbf{l}_z(\mathbf{r}) &= [l_1^z(\mathbf{r}), l_2^z(\mathbf{r}), \dots, l_M^z(\mathbf{r})]^T. \end{aligned}$$

These vectors express the sensor array sensitivity for a source located at \mathbf{r} and directed in the x , y , and z directions. Using these column vectors, the sensitivity of the whole sensor array for a source at \mathbf{r} is expressed using an $M \times 3$ matrix:

$$\mathbf{L}(\mathbf{r}) = [\mathbf{l}_x(\mathbf{r}), \mathbf{l}_y(\mathbf{r}), \mathbf{l}_z(\mathbf{r})]. \quad (2.4)$$

This matrix $\mathbf{L}(\mathbf{r})$ is called the lead-field matrix. We also define the lead-field vector, $\mathbf{l}(\mathbf{r})$, that expresses the sensitivity of the sensor array in a particular source direction $\boldsymbol{\eta}(\mathbf{r})$, such that

$$\mathbf{l}(\mathbf{r}) = \mathbf{L}(\mathbf{r})\boldsymbol{\eta}(\mathbf{r}). \quad (2.5)$$

The problem of estimating the sensor lead field is referred to as the bioelectromagnetic forward problem. Arguments on how to compute the sensor lead field are presented in Appendix A.

2.4 Voxel Source Model and Tomographic Source Reconstruction

Using the lead-field matrix in Eq. (2.4), the relationship between the sensor data, $\mathbf{y}(t)$, and the source vector, $\mathbf{s}(\mathbf{r}, t)$, is expressed as

$$\mathbf{y}(t) = \int_{\Omega} \mathbf{L}(\mathbf{r}) \mathbf{s}(\mathbf{r}, t) d\mathbf{r}. \quad (2.6)$$

Here, $d\mathbf{r}$ indicates the volume element, and the integral is performed over a volume where sources are assumed to exist. This volume is called the source space, which is denoted Ω . Equation (2.6) expresses the relationship between the sensor outputs $\mathbf{y}(t)$ and the source distribution $\mathbf{s}(\mathbf{r}, t)$.

The bioelectromagnetic inverse problem is the problem of estimating the source-vector spatial distribution, $\mathbf{s}(\mathbf{r}, t)$, from the measurements, $\mathbf{y}(t)$. Here, we assume that we know the sensor lead field $\mathbf{L}(\mathbf{r})$, although our knowledge of the sensor lead field is to some degree imperfect because it must be estimated using an analytical model or numerical computations.

When estimating $\mathbf{s}(\mathbf{r}, t)$ from $\mathbf{y}(t)$, $\mathbf{s}(\mathbf{r}, t)$ is continuous in space, while $\mathbf{y}(t)$ is discrete in space. A common strategy here is to introduce voxel discretization over the source space. Let us define the number of voxels as N , and the locations of the voxels are denoted as $\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N$. Then, the discrete form of Eq. (2.6) is expressed as:

$$\mathbf{y}(t) = \sum_{j=1}^N \mathbf{L}(\mathbf{r}_j) \mathbf{s}(\mathbf{r}_j, t) = \sum_{j=1}^N \mathbf{L}(\mathbf{r}_j) \mathbf{s}_j(t). \quad (2.7)$$

where the source vector at the j th voxel, $\mathbf{s}(\mathbf{r}_j, t)$, is denoted $\mathbf{s}_j(t)$ for simplicity. We introduce the augmented lead-field matrix over all voxel locations as

$$\mathbf{F} = [\mathbf{L}(\mathbf{r}_1), \mathbf{L}(\mathbf{r}_2), \dots, \mathbf{L}(\mathbf{r}_N)], \quad (2.8)$$

which is an $M \times 3N$ matrix. We define a $3N \times 1$ column vector containing the source vectors at all voxel locations, $\mathbf{x}(t)$, such that

$$\mathbf{x}(t) = \begin{bmatrix} \mathbf{s}_1(t) \\ \mathbf{s}_2(t) \\ \vdots \\ \mathbf{s}_N(t) \end{bmatrix}. \quad (2.9)$$

Equation (2.7) is then rewritten as

$$\mathbf{y}(t) = [\mathbf{L}(\mathbf{r}_1), \mathbf{L}(\mathbf{r}_2), \dots, \mathbf{L}(\mathbf{r}_N)] \begin{bmatrix} s_1(t) \\ s_2(t) \\ \vdots \\ s_N(t) \end{bmatrix} = \mathbf{F}\mathbf{x}(t). \quad (2.10)$$

Here, since the augmented lead-field matrix \mathbf{F} is a known quantity, the only unknown quantity is the $3N \times 1$ column vector, $\mathbf{x}(t)$. This vector $\mathbf{x}(t)$ is called the voxel source vector.

The spatial distribution of the source orientation, $\boldsymbol{\eta}(\mathbf{r})$, may be a known quantity if accurate subject anatomical information (such as high-precision subject MRI) can be obtained with accurate co-registration between the MRI coordinate and the sensor coordinate. In this case, the inverse problem is the problem of estimating the source magnitude, $s(\mathbf{r}, t)$, instead of the source vector, $\mathbf{s}(\mathbf{r}, t)$. Let us consider a situation in which the source orientations at all voxel locations are predetermined. Defining the orientation of a source at the j th voxel as $\boldsymbol{\eta}_j$, the lead field at the j th voxel is expressed as the column vector \mathbf{l}_j , which is obtained as $\mathbf{l}_j = \mathbf{L}(\mathbf{r}_j)\boldsymbol{\eta}_j$, according to Eq. (2.5). Thus, the augmented lead field is expressed as an $M \times N$ matrix \mathbf{H} defined such that

$$\mathbf{H} = [\mathbf{L}(\mathbf{r}_1)\boldsymbol{\eta}_1, \mathbf{L}(\mathbf{r}_2)\boldsymbol{\eta}_2, \dots, \mathbf{L}(\mathbf{r}_N)\boldsymbol{\eta}_N] = [\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_N], \quad (2.11)$$

whereby Eq. (2.10) can be reduced as follows:

$$\begin{aligned} \mathbf{y}(t) &= [\mathbf{L}(\mathbf{r}_1), \mathbf{L}(\mathbf{r}_2), \dots, \mathbf{L}(\mathbf{r}_N)] \begin{bmatrix} s_1(t) \\ s_2(t) \\ \vdots \\ s_N(t) \end{bmatrix} \\ &= [\mathbf{L}(\mathbf{r}_1), \mathbf{L}(\mathbf{r}_2), \dots, \mathbf{L}(\mathbf{r}_N)] \begin{bmatrix} \boldsymbol{\eta}_1 & 0 & \cdots & 0 \\ 0 & \boldsymbol{\eta}_2 & \cdot & \vdots \\ \vdots & \cdot & \ddots & 0 \\ 0 & \cdots & 0 & \boldsymbol{\eta}_N \end{bmatrix} \begin{bmatrix} s_1(t) \\ s_2(t) \\ \vdots \\ s_N(t) \end{bmatrix} \\ &= [\mathbf{L}(\mathbf{r}_1)\boldsymbol{\eta}_1, \mathbf{L}(\mathbf{r}_2)\boldsymbol{\eta}_2, \dots, \mathbf{L}(\mathbf{r}_N)\boldsymbol{\eta}_N] \begin{bmatrix} s_1(t) \\ s_2(t) \\ \vdots \\ s_N(t) \end{bmatrix} = \mathbf{H}\mathbf{x}(t). \end{aligned} \quad (2.12)$$

Thus, the voxel source vector $\mathbf{x}(t)$, in this case, is an $N \times 1$ column vector,

$$\mathbf{x}(t) = \begin{bmatrix} s_1(t) \\ s_2(t) \\ \vdots \\ s_N(t) \end{bmatrix}, \quad (2.13)$$

in which the j th component of $\mathbf{x}(t)$ is $s_j(t)$, which is the scalar intensity at the j th voxel. In this book, the same notation $\mathbf{x}(t)$ is used to indicate either the $3N \times 1$ vector in Eq. (2.9) or the $N \times 1$ vector in Eq. (2.13), unless any confusion arises.

In summary, denoting the additive noise in the sensor data $\boldsymbol{\varepsilon}$, the relationship between the sensor data $\mathbf{y}(t)$ and the voxel source vector $\mathbf{x}(t)$ is expressed as

$$\mathbf{y}(t) = \mathbf{F}\mathbf{x}(t) + \boldsymbol{\varepsilon}, \quad (2.14)$$

where $\mathbf{x}(t)$ is a $3N \times 1$ column vector in Eq. (2.9). When voxels have predetermined orientations, using the augmented lead field matrix \mathbf{H} in Eq. (2.11), the relationship between $\mathbf{y}(t)$ and $\mathbf{x}(t)$ is expressed as

$$\mathbf{y}(t) = \mathbf{H}\mathbf{x}(t) + \boldsymbol{\varepsilon}, \quad (2.15)$$

where $\mathbf{x}(t)$ is an $N \times 1$ column vector in Eq. (2.13).

2.5 Maximum Likelihood Principle and the Least-Squares Method

When estimating the unknown quantity \mathbf{x} from the sensor data \mathbf{y} , the basic principle is to interpret the data \mathbf{y} as a realization of most probable events. That is, the sensor data \mathbf{y} is considered the result of the most likely events. We call this the maximum likelihood principle. In this chapter, we first derive the maximum likelihood solution of the unknown source vector \mathbf{x} .

We assume that the noise distribution is Gaussian, i.e.,

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(\boldsymbol{\varepsilon} | \mathbf{0}, \sigma^2 \mathbf{I}).$$

Namely, the noise in the sensor data is the identically and independently distributed Gaussian noise with a mean of zero, and the same variance σ^2 . According to (C.1) in the Appendix, the explicit form of the noise probability distribution is given by

$$p(\boldsymbol{\varepsilon}) = \frac{1}{(2\pi\sigma^2)^{M/2}} \exp\left[-\frac{1}{2\sigma^2} \|\boldsymbol{\varepsilon}\|^2\right]. \quad (2.16)$$

Since the linear relationship in Eq. (2.14) holds, the probability distribution of the sensor data $\mathbf{y}(t)$ is expressed as

$$p(\mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{M/2}} \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2\right], \quad (2.17)$$

where the explicit time notation (t) is omitted from the vector notations $\mathbf{x}(t)$ and $\mathbf{y}(t)$ for simplicity.¹

This $p(\mathbf{y})$ as a function of the unknown parameter \mathbf{x} is called the likelihood function, and the maximum likelihood estimate $\hat{\mathbf{x}}$ is obtained such that²

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmax}} \log p(\mathbf{y}), \quad (2.18)$$

where $\log p(\mathbf{y})$ is called the log-likelihood function. Using the probability distribution in Eq. (2.17), the log-likelihood function $\log p(\mathbf{y})$ is expressed as

$$\log p(\mathbf{y}) = -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \mathcal{C}, \quad (2.19)$$

where \mathcal{C} expresses terms that do not contain \mathbf{x} . Therefore, the \mathbf{x} that maximizes $\log p(\mathbf{y})$ is equal to the one that minimizes $\mathcal{F}(\mathbf{x})$ defined such that

$$\mathcal{F}(\mathbf{x}) = \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2. \quad (2.20)$$

That is, the maximum likelihood solution $\hat{\mathbf{x}}$ is obtained using

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \mathcal{F}(\mathbf{x}) : \quad \text{where } \mathcal{F} = \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2. \quad (2.21)$$

This $\mathcal{F}(\mathbf{x})$ in Eq. (2.20) is referred to as the least-squares cost function, and the method that estimates \mathbf{x} through the minimization of the least-squares cost function is the method of least-squares.

2.6 Derivation of the Minimum-Norm Solution

In the bioelectromagnetic inverse problem, the number of voxels N , in general, is much greater than the number of sensors M . Thus, the estimation of the source vector \mathbf{x} is an ill-posed problem. When applying the least-squares method to such an ill-posed problem, the problem arises that an infinite number of \mathbf{x} could make the cost

¹ For the rest of this chapter, the explicit time notation is omitted from these vector notations, unless otherwise noted.

² The notation argmax indicates the value of \mathbf{x} that maximizes $\log p(\mathbf{y})$ which is an implicit function of \mathbf{x} .

function equal to zero. Therefore, we cannot obtain an optimum solution of \mathbf{x} based only on the least-squares method.

A general strategy for overcoming this problem is to integrate a “desired property” of the unknown parameter \mathbf{x} into the estimation problem. That is, we choose \mathbf{x} so as to maximize this “desired property,” and also satisfy $\mathbf{y} = \mathbf{F}\mathbf{x}$. Quite often, a small norm of the solution vector is used as this “desired property,” and in this case, the optimum estimate $\hat{\mathbf{x}}$ is obtained using

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x}\|^2 \quad \text{subject to } \mathbf{y} = \mathbf{F}\mathbf{x}. \quad (2.22)$$

In the optimization above, the notation of “subject to” indicates a constraint, (i.e., the above optimization requires that the estimate $\hat{\mathbf{x}}$ be chosen such that \mathbf{x} minimizes $\|\mathbf{x}\|^2$ as well as satisfies $\mathbf{y} = \mathbf{F}\mathbf{x}$.) To solve the constraint optimization problem in Eq. (2.22), we use the method of Lagrange multipliers that can convert a constrained optimization problem to an unconstrained optimization problem. In this method, using an $M \times 1$ column vector \mathbf{c} as the Lagrange multipliers, we define a function called the Lagrangian $\mathbb{L}(\mathbf{x}, \mathbf{c})$ such that

$$\mathbb{L}(\mathbf{x}, \mathbf{c}) = \|\mathbf{x}\|^2 + \mathbf{c}^T (\mathbf{y} - \mathbf{F}\mathbf{x}). \quad (2.23)$$

The solution $\hat{\mathbf{x}}$ is obtained by minimizing $\mathbb{L}(\mathbf{x}, \mathbf{c})$ above with respect to \mathbf{x} and \mathbf{c} —the solution $\hat{\mathbf{x}}$ being equal to $\hat{\mathbf{x}}$ obtained by solving the constrained optimization in Eq. (2.22).

To derive an \mathbf{x} that minimizes Eq. (2.23), we compute the derivatives of $\mathbb{L}(\mathbf{x}, \mathbf{c})$ with respect to \mathbf{x} and \mathbf{c} , and set them to be zero, giving

$$\frac{\partial \mathbb{L}(\mathbf{x}, \mathbf{c})}{\partial \mathbf{x}} = 2\mathbf{x} - \mathbf{F}^T \mathbf{c} = 0, \quad (2.24)$$

$$\frac{\partial \mathbb{L}(\mathbf{x}, \mathbf{c})}{\partial \mathbf{c}} = \mathbf{y} - \mathbf{F}\mathbf{x} = 0. \quad (2.25)$$

Using the equations above, we can derive

$$\hat{\mathbf{x}} = \mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} \mathbf{y}. \quad (2.26)$$

The solution in Eq. (2.26) is called the minimum-norm solution, which is well known as a solution for the ill-posed linear inverse problem.

2.7 Properties of the Minimum-Norm Solution

The minimum-norm solution is expressed as

$$\hat{\mathbf{x}} = \mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} (\mathbf{F}\mathbf{x} + \boldsymbol{\varepsilon}) = \mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} \mathbf{F}\mathbf{x} + \mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} \boldsymbol{\varepsilon}. \quad (2.27)$$

The first term on the right-hand side is expressed as $E(\hat{\mathbf{x}})$, which indicates the expectation of $\hat{\mathbf{x}}$. This term represents how the solution deviates from its true value even in the noiseless cases. The second term indicates the influence of the noise ε . The first term is rewritten as

$$E(\hat{\mathbf{x}}) = \mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} \mathbf{F}\mathbf{x} = \mathbf{Q}\mathbf{x}, \quad (2.28)$$

where

$$\mathbf{Q} = \mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} \mathbf{F}. \quad (2.29)$$

Apparently, the first term is not equal to the true value \mathbf{x} , and the matrix \mathbf{Q} in Eq. (2.29) expresses the relationship between the true value \mathbf{x} and the estimated value $E(\hat{\mathbf{x}})$.

Denoting the (i, j) th element of \mathbf{Q} as $Q_{i,j}$, the j th element of $E(\hat{\mathbf{x}})$, $E(\hat{x}_j)$, is expressed as

$$E(\hat{x}_j) = \sum_{k=1}^N Q_{j,k} x_k. \quad (2.30)$$

The above equation shows how each element of the true vector \mathbf{x} affects the value of $E(\hat{x}_j)$. That is, $Q_{j,k}$ expresses the amount of leakage of x_k into \hat{x}_j when $j \neq k$. If the weight $Q_{j,1}, \dots, Q_{j,N}$ has a sharp peak at j , \hat{x}_j may be close to the true value x_j . If the weight has no clear peak or if the weight has a peak at j' that is different from j , \hat{x}_j may be very different from x_j . Because of such properties, the matrix \mathbf{Q} is called the resolution matrix.

We next examine the second term, which expresses the noise influence. The noise influence is related to the singular values of \mathbf{F} . The singular value decomposition of \mathbf{F} is defined as

$$\mathbf{F} = \sum_{j=1}^M \gamma_j \mathbf{u}_j \mathbf{v}_j^T, \quad (2.31)$$

where we assume that $M < N$, and the singular values are numbered in decreasing order. Using

$$\mathbf{F}^T (\mathbf{F}\mathbf{F}^T)^{-1} = \sum_{j=1}^N \frac{1}{\gamma_j} \mathbf{v}_j \mathbf{u}_j^T, \quad (2.32)$$

we can express the second term in Eq. (2.27) as

$$\sum_{j=1}^N \frac{(\mathbf{u}_j^T \varepsilon)}{\gamma_j} \mathbf{v}_j. \quad (2.33)$$

The equation above shows that the denominator contains the singular values. Thus, if higher order singular values are very small and close to zero, the terms containing such small singular values amplify the noise influence, resulting in a situation where

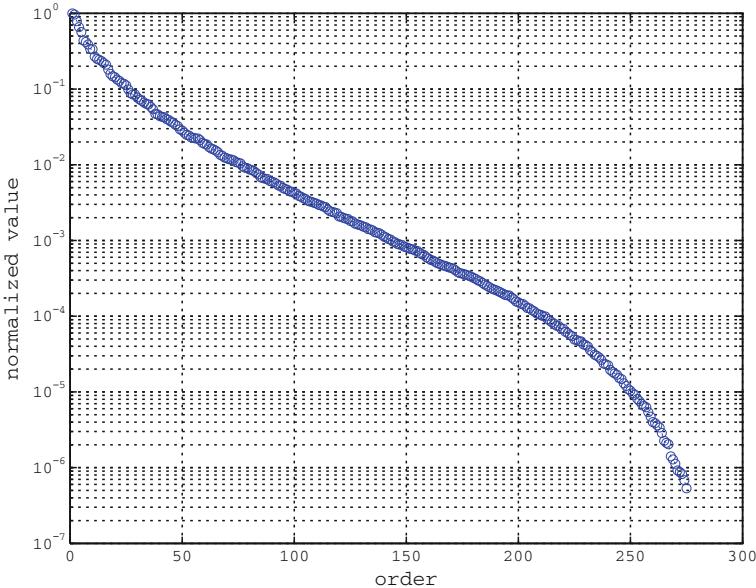


Fig. 2.1 Typical plot of the singular values of the lead field matrix \mathbf{F} . We assume the 275-channel CTF Omega whole-head MEG sensor system (VMS MedTech Ltd., BC, Canada). A typical location of the subject head relative to the whole head sensor array is assumed. An $8 \times 8 \times 10$ cm region is also assumed as the source space within the subject's head. The spherical homogeneous conductor model is used for computing the sensor lead field. The singular values are normalized with the maximum (i.e., the first) singular value

the second term is dominated in Eq. (2.27), and the minimum norm solution would contain large errors due to the noise.

A plot of a typical singular-value spectrum of the lead field matrix \mathbf{F} is shown in Fig. 2.1. To obtain the plot, we used the sensor array of the 275-channel CTF Omega whole-head MEG sensor system (VMS MedTech Ltd., BC, Canada) and spherical homogeneous conductor model to compute the sensor lead field [3].³ The plot shows that higher order singular values of the lead field matrix are very small. In Fig. 2.1, the ratio of the maximum and minimum singular values reaches the order of 10^{-7} . Therefore, the minimum-norm method in Eq. (2.26) generally produces results highly susceptible to the noise in the sensor data.

2.8 L_2 -Regularized Minimum-Norm Solution

When a large amount of noise is overlapped onto the sensor data \mathbf{y} , if we seek a solution that satisfies $\mathbf{y} = \mathbf{F}\mathbf{x}$, the resultant solution \mathbf{x} would be severely affected by the noise. In other words, when noise exists in the sensor data, it is more or less

³ Computing the lead field using the spherical homogeneous conductor model is explained in Sect. A.2.4 in the Appendix.

meaningless to impose the constraint $\mathbf{y} = \mathbf{F}\mathbf{x}$, so, instead of using the optimization in Eq. (2.22), we should use

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x}\|^2 \text{ subject to } \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 \leq d, \quad (2.34)$$

where d is a positive constant. In Eq. (2.34), the condition $\|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 \leq d$ does not require $\mathbf{F}\mathbf{x}$ to be exactly equal to \mathbf{y} , but allow $\mathbf{F}\mathbf{x}$ to be different from \mathbf{y} within a certain range specified by d . Therefore, the solution $\hat{\mathbf{x}}$ is expected to be less affected by the noise in the sensor data \mathbf{y} .

Unfortunately, there is no closed-form solution for the optimization problem in Eq. (2.34), because of the inequality constraint. Although we can solve Eq. (2.34) numerically, we proceed in solving it by replacing the inequality constraint with the equality constraint. This is possible because the solution of Eq. (2.34) generally exists on the border of the constraint. Thus, we can change the optimization problem in Eq. (2.34) to

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x}\|^2 \text{ subject to } \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 = d. \quad (2.35)$$

Since this is an equality-constraint problem, we can use the method of Lagrange multipliers. Using the Lagrange multiplier λ , the Lagrangian is defined as

$$\mathbb{L}(\mathbf{x}, \mathbf{c}) = \|\mathbf{x}\|^2 + \lambda (\|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 - d). \quad (2.36)$$

Thus, the solution $\hat{\mathbf{x}}$ is given as

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \mathbb{L}(\mathbf{x}, \mathbf{c}) = \underset{\mathbf{x}}{\operatorname{argmin}} \left[\|\mathbf{x}\|^2 + \lambda \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 \right]. \quad (2.37)$$

In the above expression, we disregard the term $-\lambda d$, which does not affect the results of the minimization. Also, we can see that the multiplier λ works as a balancer between the L_2 -norm⁴ term $\|\mathbf{x}\|^2$ and the squared error term $\|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2$.

To derive the solution of \mathbf{x} that minimizes $\mathbb{L}(\mathbf{x}, \mathbf{c})$, we compute the derivative of $\mathbb{L}(\mathbf{x}, \mathbf{c})$ with respect to \mathbf{x} and set it to zero, i.e.,

$$\begin{aligned} \frac{\mathbb{L}(\mathbf{x}, \mathbf{c})}{\partial \mathbf{x}} &= \frac{1}{\partial \mathbf{x}} \left(\mathbf{y}^T \mathbf{y} - \mathbf{x}^T \mathbf{F}^T \mathbf{y} - \mathbf{y}^T \mathbf{F} \mathbf{x} + \mathbf{x}^T \mathbf{F}^T \mathbf{F} \mathbf{x} + \xi \mathbf{x}^T \mathbf{x} \right) \\ &= -2 \mathbf{F}^T \mathbf{y} + 2 \left(\mathbf{F}^T \mathbf{F} + \xi \mathbf{I} \right) \mathbf{x} = \mathbf{0}, \end{aligned} \quad (2.38)$$

where we use $1/\lambda = \xi$. We can then derive

$$\hat{\mathbf{x}} = \left(\mathbf{F}^T \mathbf{F} + \xi \mathbf{I} \right)^{-1} \mathbf{F}^T \mathbf{y}. \quad (2.39)$$

⁴ A brief summary of the norm of vectors is presented in Sect. C.4 in the Appendix.

Using the matrix inversion lemma in Eq. (C.92), we obtain

$$\hat{\mathbf{x}} = \mathbf{F}^T \left(\mathbf{F}\mathbf{F}^T + \xi \mathbf{I} \right)^{-1} \mathbf{y}. \quad (2.40)$$

The solution in Eq. (2.40) is called the L_2 -norm-regularized minimum-norm solution, or simply L_2 -regularized minimum-norm solution.

Let us compute the noise influence term for the L_2 -regularized minimum-norm solution. Using Eq. (2.31), we have

$$\mathbf{F}^T \left(\mathbf{F}\mathbf{F}^T + \xi \mathbf{I} \right)^{-1} = \sum_{j=1}^N \frac{\gamma_j}{\gamma_j^2 + \xi} \mathbf{v}_j \mathbf{u}_j^T, \quad (2.41)$$

and the L_2 -regularized minimum-norm solution is expressed as

$$\begin{aligned} \hat{\mathbf{x}} &= \mathbf{F}^T \left(\mathbf{F}\mathbf{F}^T + \xi \mathbf{I} \right)^{-1} (\mathbf{H}\mathbf{x} + \boldsymbol{\varepsilon}) \\ &= \mathbf{F}^T \left(\mathbf{F}\mathbf{F}^T + \xi \mathbf{I} \right)^{-1} \mathbf{H}\mathbf{x} + \sum_{j=1}^N \frac{\gamma_j}{\gamma_j^2 + \xi} \mathbf{v}_j \mathbf{u}_j^T \boldsymbol{\varepsilon}. \end{aligned} \quad (2.42)$$

The second term, expressing the influence of noise, is

$$\sum_{j=1}^N \frac{\gamma_j (\mathbf{u}_j^T \boldsymbol{\varepsilon})}{\gamma_j^2 + \xi} \mathbf{v}_j. \quad (2.43)$$

In the expression above, the denominator contains the positive constant ξ , and it is easy to see that this ξ prevents the terms with smaller singular values from being amplified.

One problem here is how to choose an appropriate value for ξ . Our argument above only suggests that if the noise is large, we need a large ξ , but if small, a smaller ξ can be used. However, the arguments above do not lead to the derivation of an appropriate ξ . We will return to this problem in Sect. 2.10.2 where L_2 -regularized minimum-norm solution is re-derived based on a Bayesian formulation, in which deriving the optimum ξ is embedded.

2.9 L_1 -Regularized Minimum-Norm Solution

2.9.1 L_1 -Norm Constraint

In the preceding section, we derived a solution that minimizes the L_2 -norm of the solution vector \mathbf{x} . In this section, we argue for a solution that minimizes the L_1 -norm of \mathbf{x} , which is defined in Eq. (C.64). The L_1 -norm-regularized solution is obtained

using [4–6]

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \sum_{j=1}^N |x_j| \quad \text{subject to} \quad \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 = d. \quad (2.44)$$

The only difference between the equation above and Eq. (2.35) is to minimize either L_2 norm $\|\mathbf{x}\|^2$ in Eq. (2.35) or L_1 norm, $\|\mathbf{x}\|_1 = \sum_j |x_j|$ in Eq. (2.44). Although it may look as if there is no significant difference between the two methods, the results of source estimation are significantly different. The L_1 -norm regularization gives a “so-called” sparse solution, in which only few x_j have nonzero values and a majority of other x_j have values close to zero.

Using the method of Lagrange multipliers and following exactly the same arguments as in Sect. 2.8, the L_1 -norm solution can be obtained by minimizing the cost function \mathcal{F} , i.e.,

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \mathcal{F}: \quad \mathcal{F} = \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \xi \sum_{j=1}^N |x_j|, \quad (2.45)$$

where again ξ is a positive constant that controls the balance between the first and the second terms in the cost function above. Unfortunately, the minimization problem in Eq. (2.45) does not have a closed-form solution, so numerical methods are used here to obtain the solution $\hat{\mathbf{x}}$.

2.9.2 Intuitive Explanation for Sparsity

Actually, it is not easy to provide an intuitive explanation regarding why the optimization in Eq. (2.44) or (2.45) causes a sparse solution. The straightforward (and intuitively clear) formulation to obtain a sparse solution should use the L_0 -norm minimization, such that

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \sum_{j=1}^N \mathcal{T}(x_j) \quad \text{subject to} \quad \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 = d, \quad (2.46)$$

where the function $\mathcal{T}(x)$ is defined in Eq. (C.65). In the above formulation, since $\sum_{j=1}^N \mathcal{T}(x_j)$ indicates the number of nonzero x_j , $\hat{\mathbf{x}}$ is the solution that has the smallest number of nonzero x_j and still satisfies $\|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 = d$. The optimization problem in Eq. (2.46) is known to require impractically long computational time. The optimization for the L_1 -norm cost function in Eq. (2.44) approximates this L_0 -norm optimization in Eq. (2.46) so as to obtain a sparse solution within a reasonable range of computational time [7].

The regularization methods mentioned above can be summarized to have a form of the cost functions expressed as

$$\mathcal{F} = \|\mathbf{y} - \mathbf{Hx}\|^2 + \xi\phi(\mathbf{x}). \quad (2.47)$$

The first term is the data-fitting term, and the second term $\phi(\mathbf{x})$ expresses the constraint, which has the following form for general L_p -norm cases ($0 \leq p \leq 1$):

$$\phi(\mathbf{x}) = \sum_{j=1}^N \mathcal{T}(x_j) \quad \text{for } L_0\text{-norm,} \quad (2.48)$$

$$\phi(\mathbf{x}) = \sum_{j=1}^N |x_j| \quad \text{for } L_1\text{-norm,} \quad (2.49)$$

$$\phi(\mathbf{x}) = \left[\sum_{j=1}^N x_j^p \right]^{1/p} \quad \text{for } L_p\text{-norm.} \quad (2.50)$$

The plots of $\phi(\mathbf{x})$ with respect to the x_j axis are shown in Fig. 2.2. In this figure, the four kinds of plots of $\phi(\mathbf{x}) = \|\mathbf{x}\|_p$ when $p = 0$, $p = 0.3$, $p = 0.7$, and $p = 1$

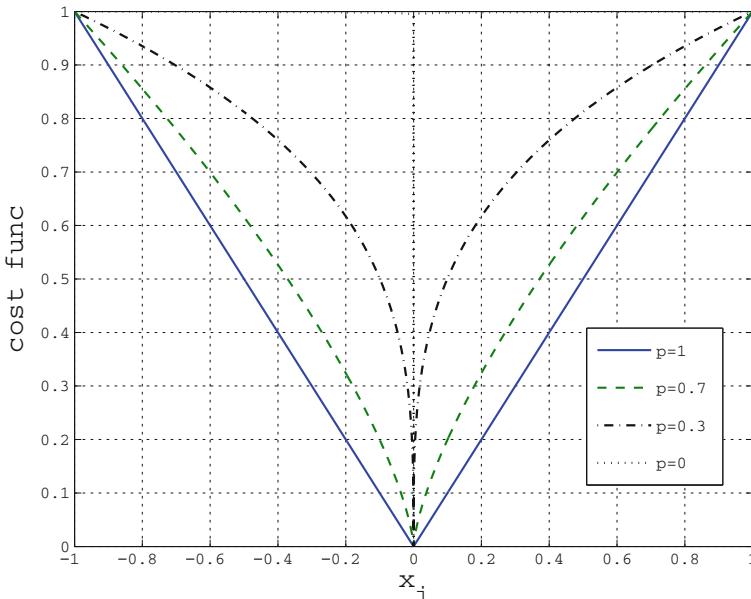


Fig. 2.2 Plots of objective function $\phi(\mathbf{x})$ defined in Eqs. (2.48)–(2.50) with respect to the x_j axis. The four cases of $p = 0$, $p = 0.3$, $p = 0.7$, and $p = 1$ are shown. The cases of $p = 0$, and $p = 1$ correspond to the L_0 and L_1 norm constraints (A brief summary of the norm of vectors is presented in Sect. C.4 in the Appendix.)

are shown. It can be seen in this figure that the L_0 -norm constraint is approximated by the L_p -norm constraint, and as p becomes closer to 0, the L_p -norm provides a better approximation.

Let us see how the L_p -norm regularization causes sparse solutions when $0 \leq p \leq 1$. To do so, we consider a simplest estimation problem in which only two voxels exist and the voxels have source intensity of x_1 and x_2 . We assume a noiseless measurement using a single-sensor; the sensor data being represented by a scalar y . The optimization for the L_1 -norm solution is expressed in this case as

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} (|x_1| + |x_2|) \quad \text{subject to} \quad y = h_1x_1 + h_2x_2, \quad (2.51)$$

where $\hat{\mathbf{x}} = (\hat{x}_1, \hat{x}_2)^T$, and h_1 and h_2 are the sensor lead field. For the sake of comparison, we also argue the L_2 -norm regularization whose optimization is given as follows:

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} (x_1^2 + x_2^2)^{1/2} \quad \text{subject to} \quad y = h_1x_1 + h_2x_2. \quad (2.52)$$

The optimization process is depicted in Fig. 2.3. In Fig. 2.3a, the tetragon at the center represents the L_1 -norm objective function, $|x_1| + |x_2| = \text{constant}$. The broken line represents the x_1 and x_2 that satisfy the measurement equation $y = h_1x_1 + h_2x_2$. Thus, as a result of the optimization in Eq. (2.51), the x_1 and x_2 on the broken line that minimize $|x_1| + |x_2|$ should be chosen as the solution, i.e., the point (x_1, x_2) at which the tetragon touches the broken line is chosen as the solution for the optimization. Such solution is indicated by the small filled circle in Fig. 2.3a. In this solution, x_2 has a nonzero value but x_1 is zero, i.e., a sparse solution is obtained. It can be seen in this figure that in most cases, the point at which the tetragon touches the broken

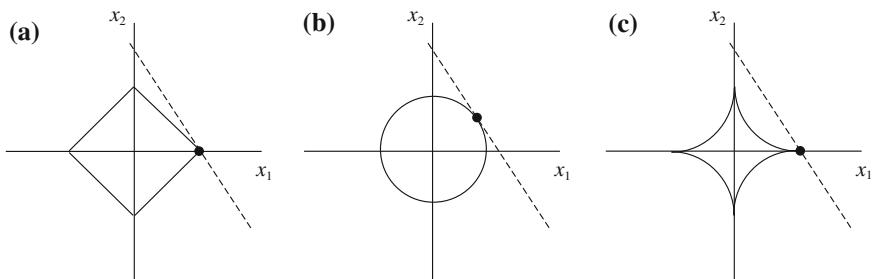


Fig. 2.3 The optimization process is depicted for the simple case in which a single sensor and two voxels exist. Source magnitudes at the voxels are represented by x_1 and x_2 . The broken lines represent the x_1 and x_2 that satisfy the measurement equation, $y = h_1x_1 + h_2x_2$. The filled black circles indicate an example of the solution for each case. **a** L_1 -norm regularization in Eq. (2.51). The tetragon at the center represents the L_1 -norm objective function $|x_1| + |x_2| = \text{constant}$. **b** L_2 -norm regularization in Eq. (2.52). The circle at the center represents the L_2 -norm objective function $x_1^2 + x_2^2 = \text{constant}$. **c** L_p -norm regularization where $0 < p < 1$

line is likely to be located at one of its vertices, so a sparse solution is likely to be obtained.

Figure 2.3b shows the case of the L_2 -norm minimization in Eq. (2.52). In this figure, the broken line again represents the x_1 and x_2 that satisfy the measurement equation $y = h_1x_1 + h_2x_2$, and the circle represents the L_2 -norm objective function $x_1^2 + x_2^2 = \text{constant}$. In this case, the x_1 and x_2 on the broken line that minimizes $x_1^2 + x_2^2$ should be chosen, and the resultant solution is (x_1, x_2) at which the circle touches the broken line. An example of such solution is indicated by the small filled circle. In this case, both x_1 and x_2 have nonzero values, and a non-sparse solution is likely to be obtained using L_2 -norm regularization.

Finally, Fig. 2.3c shows a case of the general L_p norm minimization ($0 < p < 1$). An example of such solution is indicated by the small, filled circle. Using the general L_p norm regularization, the solution is more likely to be sparse than the case of the L_1 -norm minimization. However, the computational burden for the general L_p norm minimization is so high that it is seldom used in practical applications.

2.9.3 Problem with Source Orientation Estimation

When applying the L_1 -norm regularization to the bioelectromagnetic source localization, it has been known that the method fails in estimating correct source orientations. The reason for this is described as follows: The components of the solution vector \mathbf{x} is denoted explicitly as

$$\mathbf{x} = [s_1^x, s_1^y, s_1^z, \dots, s_j^x, s_j^y, s_j^z, \dots, s_N^x, s_N^y, s_N^z]^T,$$

where s_j^x, s_j^y, s_j^z are the x , y , and z components of the source at the j th voxel. When the j th voxel has a source activity, it is generally true that s_j^x, s_j^y, s_j^z have nonzero values. However, when using the L_1 regularization, only one of s_j^x, s_j^y, s_j^z tends to have nonzero value, and others tend to be close to zero because of the nature of a sparse solution. As a result, the source orientation may be erroneously estimated.

To avoid this problem, the source orientation is estimated in advance using some other method [4] such as the L_2 -norm minimum-norm method. Then, the L_1 -norm method is formulated using the orientation-embedded data model in Eq. (2.15). That is, we use

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \mathcal{F} : \quad \mathcal{F} = \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + \xi \sum_{j=1}^N |x_j|. \quad (2.53)$$

In this case, the sparsity is imposed on the source vector magnitude, s_1, s_2, \dots, s_N , and only a few of s_1, s_2, \dots, s_N have nonzero values, allowing for the reconstruction of a sparse source distribution.

2.10 Bayesian Derivation of the Minimum-Norm Method

2.10.1 Prior Probability Distribution and Cost Function

In this section, we derive the minimum-norm method based on Bayesian inference. As in Eq. (2.16), we assume that the noise ε is independently and identically distributed Gaussian, i.e.,

$$\varepsilon \sim \mathcal{N}(\varepsilon | \mathbf{0}, \beta^{-1} \mathbf{I}), \quad (2.54)$$

where the precision β is used, which is the inverse of the noise variance, $\beta^{-1} = \sigma^2$. Thus, using Eq. (2.14), the conditional probability distribution of the sensor data for a given \mathbf{x} , $p(\mathbf{y}|\mathbf{x})$ is

$$p(\mathbf{y}|\mathbf{x}) = \left(\frac{\beta}{2\pi} \right)^{M/2} \exp \left[-\frac{\beta}{2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 \right]. \quad (2.55)$$

This conditional probability $p(\mathbf{y}|\mathbf{x})$ is equal to the likelihood $p(\mathbf{y})$ in the arguments in Sect. 2.5. Since \mathbf{x} is a random variable in the Bayesian arguments, we use the conditional probability $p(\mathbf{y}|\mathbf{x})$, instead of $p(\mathbf{y})$.

Let us derive a cost function for estimating \mathbf{x} . Taking a logarithm of the Bayes's rule in Eq. (B.3) in the Appendix, we have

$$\log p(\mathbf{x}|\mathbf{y}) = \log p(\mathbf{y}|\mathbf{x}) + \log p(\mathbf{x}) + \mathcal{C}, \quad (2.56)$$

where \mathcal{C} represents the constant terms. Neglecting \mathcal{C} , the cost function $\mathcal{F}(\mathbf{x})$ in general form is obtained as

$$\mathcal{F}(\mathbf{x}) = -2 \log p(\mathbf{x}|\mathbf{y}) = \beta \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 - 2 \log p(\mathbf{x}). \quad (2.57)$$

The first term on the right-hand side is a squared error term, which expresses how well the solution \mathbf{x} fits the sensor data \mathbf{y} . The second term $-2 \log p(\mathbf{x})$ is a constraint imposed on the solution. The above equation indicates that the constraint term in the cost function is given from the prior probability distribution in the Bayesian formulation. The optimum estimate of \mathbf{x} is obtained by minimizing the cost function $\mathcal{F}(\mathbf{x})$.

2.10.2 L_2 -Regularized Method

Let us assume the following Gaussian distribution for the prior probability distribution of \mathbf{x} ,

$$p(\mathbf{x}) = \left(\frac{\alpha}{2\pi} \right)^{N/2} \exp \left[-\frac{\alpha}{2} \|\mathbf{x}\|^2 \right]. \quad (2.58)$$

Substituting Eq. (2.58) into (2.57), we get the cost function

$$\mathcal{F}(\mathbf{x}) = \beta \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \alpha \|\mathbf{x}\|^2. \quad (2.59)$$

The cost function in Eq. (2.59) is the same as the cost function in Eq. (2.37), assuming $\lambda = \beta/\alpha$. Thus, the solution obtained by minimizing this cost function is equal to the solution of the L_2 -norm regularized minimum-norm method introduced in Sect. 2.8.

To obtain the optimum estimate of \mathbf{x} , we should compute the posterior distribution. In this case, the posterior is known to have a Gaussian distribution because $p(\mathbf{y}|\mathbf{x})$ and $p(\mathbf{x})$ are both Gaussian, and the mean and the precision matrix of this posterior distribution is derived as in Eqs. (B.24) and (B.25). Substituting $\Phi = \alpha\mathbf{I}$ and $\Lambda = \beta\mathbf{I}$ into these equations, we have

$$\boldsymbol{\Gamma} = \alpha\mathbf{I} + \beta\mathbf{F}^T\mathbf{F}, \quad (2.60)$$

$$\bar{\mathbf{x}}(t) = \left(\mathbf{F}^T\mathbf{F} + \frac{\alpha}{\beta}\mathbf{I} \right)^{-1} \mathbf{F}^T\mathbf{y}(t). \quad (2.61)$$

The Bayesian solution which minimizes the cost function in Eq. (2.59) is given in Eq. (2.61). This solution is the same as Eq. (2.39). Comparison between Eqs. (2.61) and (2.39) shows that the regularization constant is equal to α/β , which is the inverse of the signal-to-noise ratio of the sensor data. This is in accordance with the arguments in Sect. 2.8 that when the sensor data contains larger amounts of noise, a larger regularization constant must be used.

The optimum values of the hyperparameters α and β can be obtained using the EM algorithm, as described in Sect. B.5.6. The update equations for the hyperparameters are:

$$\hat{\alpha}^{-1} = \frac{1}{3N} \left[\frac{1}{K} \sum_{k=1}^K \bar{\mathbf{x}}^T(t_k) \bar{\mathbf{x}}(t_k) + \text{tr}(\boldsymbol{\Gamma}^{-1}) \right], \quad (2.62)$$

$$\hat{\beta}^{-1} = \frac{1}{M} \left[\frac{1}{K} \sum_{k=1}^K \|\mathbf{y}(t_k) - \mathbf{F}\bar{\mathbf{x}}(t_k)\|^2 + \text{tr}(\mathbf{F}^T\mathbf{F}\boldsymbol{\Gamma}^{-1}) \right]. \quad (2.63)$$

Here, we assume that multiple K time-point data is available to determine α and β .

The Bayesian minimum-norm method is summarized as follows. First, $\boldsymbol{\Gamma}$ and $\bar{\mathbf{x}}(t_k)$ are computed using Eqs. (2.60) and (2.61) with initial values set to α and β . Then, the values of α and β are updated using (2.62) and (2.63). Using the updated α and β , the values of $\boldsymbol{\Gamma}$ and $\bar{\mathbf{x}}(t_k)$ are updated using Eqs. (2.60) and (2.61). These procedures are repeated and the resultant $\bar{\mathbf{x}}(t_k)$ is the optimum estimate of $\mathbf{x}(t_k)$.

The EM iteration may be stopped by monitoring the marginal likelihood, which is obtained using Eq. (B.29) as

$$\log p(\mathbf{y}(t_1), \dots, \mathbf{y}(t_K) | \alpha, \beta) = -\frac{1}{2} K \log |\boldsymbol{\Sigma}_y| - \frac{1}{2} \sum_{k=1}^K \mathbf{y}^T(t_k) \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t_k), \quad (2.64)$$

where according to Eq. (B.30), Σ_y is expressed as

$$\Sigma_y = \beta^{-1} \mathbf{I} + \alpha^{-1} \mathbf{F} \mathbf{F}^T. \quad (2.65)$$

If the increase of the likelihood in Eq. (2.64) with respect to the iteration count becomes very small, the iteration may be stopped.

2.10.3 L_1 -Regularized Method

The method of L_1 -norm regularization can also be derived based on the Bayesian formulation. To derive the L_1 -regularization, we use the Laplace distribution as the prior distribution

$$p(\mathbf{x}) = \prod_{j=1}^N \frac{1}{2b} \exp\left[-\frac{1}{b}|x_j|\right]. \quad (2.66)$$

Then, using Eq. (2.57), (and replacing \mathbf{F} with \mathbf{H}), the cost function is derived as

$$\mathcal{F}(\mathbf{x}) = \beta \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + 2b \sum_{j=1}^N |x_j|, \quad (2.67)$$

which is exactly equal to Eq. (2.53), if we set $\xi = 2b/\beta$.

Another formulation for deriving the L_1 -regularized method is known. It uses the framework of the sparse Bayesian learning described in Chap. 4. In Chap. 4, assuming the Gaussian prior,

$$p(\mathbf{x}|\boldsymbol{\alpha}) = \prod_{j=1}^N \mathcal{N}(x_j|0, \alpha_j^{-1}) = \prod_{j=1}^N \left(\frac{\alpha_j}{2\pi}\right)^{1/2} \exp\left[-\frac{\alpha_j}{2}x_j^2\right], \quad (2.68)$$

we derive the marginal likelihood for the hyperparameter $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]$, $p(\mathbf{y}|\boldsymbol{\alpha})$, using,

$$p(\mathbf{y}|\boldsymbol{\alpha}) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\boldsymbol{\alpha})d\mathbf{x}, \quad (2.69)$$

and eventually derive the Champagne algorithm. However, instead of implementing Eq. (2.69), there is another option in which we compute the posterior distribution $p(\mathbf{x}|\mathbf{y})$ using

$$p(\mathbf{x}|\mathbf{y}) \propto \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\boldsymbol{\alpha})p(\boldsymbol{\alpha})d\boldsymbol{\alpha} = p(\mathbf{y}|\mathbf{x})p(\mathbf{x}), \quad (2.70)$$

where

$$p(\mathbf{x}) = \int p(\mathbf{x}|\boldsymbol{\alpha})p(\boldsymbol{\alpha})d\boldsymbol{\alpha}. \quad (2.71)$$

The estimate $\hat{\mathbf{x}}$ is, then, obtained by

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{y}|\mathbf{x})p(\mathbf{x}).$$

To compute $p(\mathbf{x})$ using Eq. (2.71), we need to specify the hyperprior $p(\boldsymbol{\alpha})$. However, we usually have no such information and may use noninformed prior $p(\boldsymbol{\alpha}) = \text{const}$. Substituting this flat prior into Eq. (2.71), we have

$$p(\mathbf{x}) \propto \int p(\mathbf{x}|\boldsymbol{\alpha})d\boldsymbol{\alpha}.$$

However, the integral in the above equation is difficult to compute. The formal procedure to compute $p(\mathbf{x})$ in this case is to first assume the Gamma distribution for the hyperprior $p(\boldsymbol{\alpha})$, such that

$$p(\boldsymbol{\alpha}) = \prod_{j=1}^N p(\alpha_j) = \prod_{j=1}^N \Gamma(a)^{-1} b^a (\alpha_j)^{a-1} e^{-b\alpha_j}. \quad (2.72)$$

Then, $p(\mathbf{x})$ in Eq. (2.71) is known to be obtained as Student t -distribution, such that [8]

$$\begin{aligned} p(x_j) &= \int p(x_j|\alpha_j)p(\alpha_j)d\alpha_j \\ &= \int \left(\frac{\alpha_j}{2\pi}\right)^{1/2} \exp\left(-\frac{\alpha_j}{2}x_j^2\right) \frac{b^a}{\Gamma(a)} (\alpha_j)^{a-1} e^{-b\alpha_j} d\alpha_j \\ &= \frac{b^a \Gamma(a + \frac{1}{2})}{\sqrt{2\pi} \Gamma(a)} \left(b + \frac{x_j^2}{2}\right)^{-(a+\frac{1}{2})}. \end{aligned} \quad (2.73)$$

We then assume that $a \rightarrow 0$ and $b \rightarrow 0$, (which is equivalent to making $p(\boldsymbol{\alpha})$ a noninformed prior,) $p(x_j)$ then becomes

$$p(x_j) \rightarrow \frac{1}{|x_j|} \quad \text{i.e.} \quad p(\mathbf{x}) \rightarrow \prod_{j=1}^N \frac{1}{|x_j|}. \quad (2.74)$$

Using Eq. (2.57), the cost function, in this case, is derived as

$$\mathcal{F}(\mathbf{x}) = \beta \|\mathbf{y} - \mathbf{Fx}\|^2 + \sum_{j=1}^N \log |x_j|. \quad (2.75)$$

This Eq. (2.75) is not exactly equal to the L_1 -norm cost function, since the constraint term is not equal to $\sum_{j=1}^N |x_j|$ but has form of $\sum_{j=1}^N \log |x_j|$. Since these constraint terms have similar properties, the solution obtained by minimizing this cost function has a property very similar to the L_1 -norm-regularized minimum-norm solution. Related arguments are found in Chap. 6.

References

1. M.S. Hämäläinen, R.J. Ilmoniemi, Interpreting measured magnetic fields of the brain: estimates of current distributions. Technical Report TKK-F-A559, Helsinki University of Technology (1984)
2. M.S. Hämäläinen, R.J. Ilmoniemi, Interpreting magnetic fields of the brain: minimum norm estimates. *Med. Biol. Eng. Comput.* **32**, 35–42 (1994)
3. J. Sarvas, Basic mathematical and electromagnetic concepts of the biomagnetic inverse problem. *Phys. Med. Biol.* **32**, 11–22 (1987)
4. K. Uutela, M. Hämäläinen, E. Somersalo, Visualization of magnetoencephalographic data using minimum current estimate. *NeuroImage* **10**, 173–180 (1999)
5. B.D. Jeffs, Maximally sparse constrained optimization for signal processing applications. Ph.D. thesis, University of Southern California (1989)
6. K. Matsuura, Y. Okabe, Multiple current-dipole distribution reconstructed by modified selective minimum-norm method, in *Biomag 96*, (Springer, Heidelberg, 2000), pp. 290–293
7. R. Tibshirani, Regression shrinkage and selection via the Lasso. *J. R. Stat. Soc. Ser. B (Methodological)* **58**(1), 267–288 (1996)
8. M.E. Tipping, Sparse Bayesian learning and relevance vector machine. *J. Mach. Learn. Res.* **1**, 211–244 (2001)

Chapter 3

Adaptive Beamformers

3.1 Introduction and Basic Formulation

The beamformer is a location-dependent spatial filter applied to the sensor data. It is used for estimating the strength of brain activity at a particular spatial location, which is referred to as the beamformer pointing location. Thus, the beamformer may be interpreted as a technique that forms a virtual sensor whose sensitivity pattern is localized at its pointing location. By postprocessing, the pointing location can be scanned over the source space to obtain the source three-dimensional reconstruction.

Let us assume momentarily that the source orientation is predetermined and avoid the orientation estimation. In this case, using the data vector, $\mathbf{y}(t)$, the beamformer reconstructs the source magnitude, $s(\mathbf{r}, t)$, using

$$\hat{s}(\mathbf{r}, t) = \mathbf{w}^T(\mathbf{r})\mathbf{y}(t), \quad (3.1)$$

where $\hat{s}(\mathbf{r}, t)$ is the estimated source magnitude at location \mathbf{r} and time t . In Eq. (3.1), a column vector $\mathbf{w}(\mathbf{r})$ expresses the beamformer's weight, which characterizes the properties of the beamformer. There are two types of beamformers. One is the non-adaptive beamformer in which the weight vector depends solely on the lead field of the sensor array, and the other is the adaptive beamformer in which the weight depends on the measured data as well as the lead field of the sensor array.

Adaptive beamformers were originally developed in the field of seismic exploration [1] and introduced later into the field of electromagnetic brain imaging [2–5]. In recent years, brain imaging with adaptive beamformers has increasingly been used for clinical and basic human neuroscience studies, and it is a very popular method for analyzing rhythmic brain activity. You may find detailed arguments on various aspects of adaptive beamformers in [6].

This chapter presents a concise review on the adaptive beamformers. It presents Bayesian-flavored formulations of scalar and vector adaptive beamformers, as well as the conventional derivations. It also describes the narrow-band beamformer and its application to five-dimensional (space–time–frequency) brain imaging [7], which can be used for the source-space connectivity analysis.

3.2 Classical Derivation of Adaptive Beamformers

3.2.1 Minimum-Variance Beamformers with Unit-Gain Constraint

The weight vector of the adaptive beamformer is derived using the optimization:

$$\mathbf{w}(\mathbf{r}) = \underset{\mathbf{w}(\mathbf{r})}{\operatorname{argmin}} \mathbf{w}^T(\mathbf{r}) \mathbf{R} \mathbf{w}(\mathbf{r}), \quad \text{subject to } \mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = 1. \quad (3.2)$$

Here, \mathbf{R} is the data covariance matrix obtained using $\langle \mathbf{y}(t) \mathbf{y}^T(t) \rangle$ where $\langle \cdot \rangle$ indicates the ensemble average. In Eq. (3.2), the inner product $\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r})$ represents the beamformer output from a unit-magnitude source located at \mathbf{r} . Therefore, setting $\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = 1$ guarantees that the beamformer passes the signal from \mathbf{r} with the gain equal to one. The constraint $\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = 1$ is called the unit-gain constraint.

The output power of the beamformer $\mathbf{w}^T(\mathbf{r}) \mathbf{R} \mathbf{w}(\mathbf{r})$ generally contains not only the noise contributions but also unwanted contributions such as the influence of sources at locations other than \mathbf{r} . Accordingly, by minimizing the output power with this unit-gain constraint, we can derive a weight that minimizes such unwanted influence without affecting the signal coming from \mathbf{r} , the pointing location of the beamformer.

This constrained minimization problem can be solved using a method of the Lagrange multiplier. We define the Lagrange multiplier as a scalar ζ , and the Lagrangian as $\mathbb{L}(\mathbf{w}, \zeta)$, such that

$$\mathbb{L}(\mathbf{w}, \zeta) = \mathbf{w}^T \mathbf{R} \mathbf{w} + \zeta (\mathbf{w}^T \mathbf{l}(\mathbf{r}) - 1), \quad (3.3)$$

where the explicit notation of (\mathbf{r}) is omitted from $\mathbf{w}(\mathbf{r})$ for simplicity. The weight vector satisfying Eq. (3.2) can be obtained by minimizing the Lagrangian $\mathbb{L}(\mathbf{w}, \zeta)$ in Eq. (3.3) with no constraints.

The derivative of $\mathbb{L}(\mathbf{w}, \zeta)$ with respect to \mathbf{w} is given by:

$$\frac{\partial \mathbb{L}(\mathbf{w}, \zeta)}{\partial \mathbf{w}} = 2 \mathbf{R} \mathbf{w} + \zeta \mathbf{l}(\mathbf{r}). \quad (3.4)$$

By setting the right-hand side of the above equation to zero, we obtain

$$\mathbf{w} = -\zeta \mathbf{R}^{-1} \mathbf{l}(\mathbf{r}) / 2. \quad (3.5)$$

Substituting this relationship back into the constraint equation $\mathbf{w}^T \mathbf{l}(\mathbf{r}) = 1$, we get $\zeta = -2 / [\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{l}(\mathbf{r})]$. Substituting this ζ into Eq. (3.5), the weight vector satisfying Eq. (3.2) is obtained as

$$\mathbf{w}(\mathbf{r}) = \frac{\mathbf{R}^{-1} \mathbf{l}(\mathbf{r})}{[\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{l}(\mathbf{r})]}. \quad (3.6)$$

Using Eq. (3.1), the estimated source magnitude (the beamformer output) is expressed as

$$\hat{s}(\mathbf{r}, t) = \frac{\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{y}(t)}{[\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{l}(\mathbf{r})]}, \quad (3.7)$$

and the output power is given by

$$\langle \hat{s}(\mathbf{r}, t)^2 \rangle = \frac{1}{[\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{l}(\mathbf{r})]}. \quad (3.8)$$

3.2.2 Minimum-Variance Beamformer with Array-Gain Constraint

We have derived the minimum-variance beamformer with the unit-gain constraint, $\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = 1$. However, imposing the unit-gain constraint is somewhat ad hoc, and there may be other possibilities. In electromagnetic brain imaging, the norm of the lead field $\|\mathbf{l}(\mathbf{r})\|$ represents the gain of the sensor array, and is spatially dependent. Particularly, when the spherical homogeneous conductor model¹ is used, $\|\mathbf{l}(\mathbf{r})\|$ is zero at the center of the sphere and a false intensity increase around the center of the sphere arises, because the weight vector (in Eq. (3.6)) becomes infinite at the center of the sphere. Such false results occur because the unit-gain constraint forces a nonzero gain at a location where the sensor array has zero gain.

When $\|\mathbf{l}(\mathbf{r})\|$ has a spatial dependence, it is more reasonable to use the constraint $\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = \|\mathbf{l}(\mathbf{r})\|$. We can derive, by using $\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = \|\mathbf{l}(\mathbf{r})\|$, a beamformer whose gain exactly matches the gain of the sensor array. The weight, in this case, is obtained using

$$\mathbf{w}(\mathbf{r}) = \underset{\mathbf{w}(\mathbf{r})}{\operatorname{argmin}} \mathbf{w}^T(\mathbf{r}) \mathbf{R} \mathbf{w}(\mathbf{r}), \quad \text{subject to } \mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = \|\mathbf{l}(\mathbf{r})\|. \quad (3.9)$$

Exactly, the same derivation used for deriving Eq. (3.6) leads to the weight vector:

$$\mathbf{w}(\mathbf{r}) = \frac{\mathbf{R}^{-1} \tilde{\mathbf{l}}(\mathbf{r})}{[\tilde{\mathbf{l}}^T(\mathbf{r}) \mathbf{R}^{-1} \tilde{\mathbf{l}}(\mathbf{r})]}, \quad (3.10)$$

where $\tilde{\mathbf{l}}(\mathbf{r})$ is the normalized lead field vector defined as $\tilde{\mathbf{l}}(\mathbf{r}) = \mathbf{l}(\mathbf{r}) / \|\mathbf{l}(\mathbf{r})\|$. In Eq. (3.10), the weight is independent of the norm of the lead field, and we can avoid the artifacts around the center of the sphere. This type of beamformer is referred to as the array-gain minimum-variance beamformer [6], and the constraint $\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = \|\mathbf{l}(\mathbf{r})\|$ is referred to as the array-gain constraint. The estimated source magnitude for the array-gain minimum-variance beamformer is given by

¹ An explanation of the spherical homogeneous conductor model is given in the Appendix A.

$$\hat{s}(\mathbf{r}, t) = \frac{\tilde{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{y}(t)}{[\tilde{l}^T(\mathbf{r}) \mathbf{R}^{-1} \tilde{l}(\mathbf{r})]}, \quad (3.11)$$

and the output power is expressed as

$$\langle \hat{s}(\mathbf{r}, t)^2 \rangle = \frac{\mathbf{l}^T(\mathbf{r}) \mathbf{l}(\mathbf{r})}{[\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{l}(\mathbf{r})]}. \quad (3.12)$$

3.2.3 Minimum-Variance Beamformer with Unit-Noise-Gain Constraint

Another possible constraint is the unit-noise-gain constraint, which is expressed as $\mathbf{w}^T(\mathbf{r})\mathbf{w}(\mathbf{r}) = 1$. That is, the filter weight is obtained using

$$\begin{aligned} \mathbf{w}(\mathbf{r}) = \operatorname{argmin}_{\mathbf{w}(\mathbf{r})} \mathbf{w}^T(\mathbf{r}) \mathbf{R} \mathbf{w}(\mathbf{r}), \text{ subject to } & \mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = \tau, \\ & \text{and } \mathbf{w}^T(\mathbf{r}) \mathbf{w}(\mathbf{r}) = 1, \end{aligned} \quad (3.13)$$

where the minimization problem is solved with the first constraint, $\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = \tau$. The scalar constant τ is determined by the second constraint, $\mathbf{w}^T(\mathbf{r}) \mathbf{w}(\mathbf{r}) = 1$. To obtain the weight vector derived from the above minimization, we first calculate the weight using

$$\mathbf{w}(\mathbf{r}) = \operatorname{argmin}_{\mathbf{w}(\mathbf{r})} \mathbf{w}^T(\mathbf{r}) \mathbf{R} \mathbf{w}(\mathbf{r}) \text{ subject to } \mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = \tau. \quad (3.14)$$

Following the same steps from Eqs. (3.3) to (3.6), the weight satisfying Eq. (3.14) is obtained as

$$\mathbf{w}(\mathbf{r}) = \tau \frac{\mathbf{R}^{-1} \mathbf{l}(\mathbf{r})}{[\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{l}(\mathbf{r})]}. \quad (3.15)$$

Substituting this expression to $\mathbf{w}^T(\mathbf{r}) \mathbf{w}(\mathbf{r}) = 1$ leads to

$$\tau = \frac{\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{l}(\mathbf{r})}{\sqrt{\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-2} \mathbf{l}(\mathbf{r})}}, \quad (3.16)$$

and the weight is given by:

$$\mathbf{w}(\mathbf{r}) = \frac{\mathbf{R}^{-1} \mathbf{l}(\mathbf{r})}{\sqrt{\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-2} \mathbf{l}(\mathbf{r})}}. \quad (3.17)$$

This weight vector again does not depend on the norm of the lead field $\|\mathbf{l}(\mathbf{r})\|$. The output power of this beamformer is given by

$$\langle \hat{s}(\mathbf{r}, t)^2 \rangle = \frac{\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{l}(\mathbf{r})}{[\mathbf{l}^T(\mathbf{r}) \mathbf{R}^{-2} \mathbf{l}(\mathbf{r})]}. \quad (3.18)$$

This beamformer was first proposed by Borgiotti and Kaplan [8] and it is referred to as the unit-noise-gain (constraint) minimum-variance beamformer,² or the weight-normalized minimum-variance beamformer.

3.3 Semi-Bayesian Derivation of Adaptive Beamformers

The adaptive beamformer can be derived based on a Bayesian formulation [9]. Let the relationship between the sensor data $\mathbf{y}(t)$ and the voxel source distribution $\mathbf{x}(t)$ be expressed in Eq. (2.15), rewritten here as:

$$\mathbf{y}(t) = \mathbf{H}\mathbf{x}(t) + \boldsymbol{\varepsilon},$$

where \mathbf{H} is the lead field matrix defined in Eq. (2.11). We also assume that the noise in the sensor data is assumed such that

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(\boldsymbol{\varepsilon} | \mathbf{0}, \sigma^2 \mathbf{I}). \quad (3.19)$$

We assume that the prior distribution for the source vector $\mathbf{x}(t)$ is the zero mean Gaussian with a diagonal precision matrix, i.e.,

$$p(\mathbf{x}(t)) = \mathcal{N}(\mathbf{x}(t) | \mathbf{0}, \boldsymbol{\Phi}^{-1}), \quad (3.20)$$

where the precision matrix is expressed as

$$\boldsymbol{\Phi} = \begin{bmatrix} \alpha_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \alpha_N \end{bmatrix}.$$

An entirely rigorous Bayesian treatment using this prior distribution leads to the Champagne algorithm, which is the topic of Chap. 4. In the following, we show that adaptive beamformer algorithm can also be derived using this prior distribution.

² This name comes from the fact that spatial filter's noise gain is equal to the squared weight norm $\|\mathbf{w}(\mathbf{r})\|^2$.

The prior distribution in Eq. (3.20) leads to Bayes' estimate of \mathbf{x} expressed in Eq. (B.26) in Appendix. Setting the noise precision Λ equal to $\sigma^{-2}\mathbf{I}$ in Eq. (B.26), we get

$$\bar{\mathbf{x}}(t) = \boldsymbol{\Phi}^{-1} \mathbf{H}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t), \quad (3.21)$$

where

$$\boldsymbol{\Sigma}_y = \mathbf{H} \boldsymbol{\Phi}^{-1} \mathbf{H}^T + \sigma^2 \mathbf{I}$$

is the model data covariance matrix, which is the covariance of the marginal distribution $p(\mathbf{y}(t))$. Taking a look at the j th component of $\bar{\mathbf{x}}(t)$ in Eq. (3.21), we have

$$\bar{x}_j(t) = \hat{s}(\mathbf{r}_j, t) = \alpha_j^{-1} \mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t), \quad (3.22)$$

where \mathbf{l}_j is the lead field vector: $\mathbf{l}_j = \mathbf{l}(\mathbf{r}_j)$. In the equation above, the voxel precision α_j is an unknown quantity. Thus, computing $\bar{x}_j(t)$ first requires estimating α_j . A rigorous Bayesian treatment for estimating the voxel precision leads to the Champagne algorithm. Here, we perform a non-Bayesian treatment for estimating α_j .

Let us treat $\boldsymbol{\Phi}$ as a deterministic (i.e., nonrandom variable) unknown matrix. Using Eq. (3.22), we compute the power of $\bar{x}_j(t)$ as

$$\langle \bar{x}_j(t)^2 \rangle = \alpha_j^{-2} \mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \langle \mathbf{y}(t) \mathbf{y}^T(t) \rangle \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j. \quad (3.23)$$

If we assume that $\langle \mathbf{y}(t) \mathbf{y}^T(t) \rangle$ is equal to the model data covariance $\boldsymbol{\Sigma}_y$, Eq. (3.23) may be reexpressed as

$$\langle \hat{s}(\mathbf{r}_j, t)^2 \rangle = \alpha_j^{-2} \mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \boldsymbol{\Sigma}_y \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j = \alpha_j^{-2} \mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j. \quad (3.24)$$

Next we assume that the relationship,

$$\langle \hat{s}(\mathbf{r}_j, t)^2 \rangle = \alpha_j^{-1} \quad (3.25)$$

holds. This relationship indicates that the reconstructed power of a brain source is equal to its true power, and this is equal to the unit-gain constraint.

Assuming that the unit-gain constraint holds, we have

$$\alpha_j^{-1} = \alpha_j^{-2} \mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j,$$

and thus,

$$\alpha_j^{-1} = \frac{1}{\mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j}. \quad (3.26)$$

Substituting this equation into Eq. (3.22), we can derive the beamformer expression,

$$\hat{s}(\mathbf{r}_j, t) = \mathbf{w}^T(\mathbf{r}_j) \mathbf{y}(t), \quad (3.27)$$

where the weight vector is expressed as

$$\mathbf{w}(\mathbf{r}_j) = \frac{\boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j^T}{\mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j}. \quad (3.28)$$

If the model data covariance matrix $\boldsymbol{\Sigma}_y$ is replaced with the sample data covariance matrix \mathbf{R} , the weight equation (3.28) becomes

$$\mathbf{w}(\mathbf{r}_j) = \frac{\mathbf{R}^{-1} \mathbf{l}^T(\mathbf{r}_j)}{\mathbf{l}^T(\mathbf{r}_j) \mathbf{R}^{-1} \mathbf{l}(\mathbf{r}_j)}, \quad (3.29)$$

which is exactly equal to the weight expression of the minimum-variance beamformer. To derive the array-gain constraint beamformer, we use the relationship,

$$\langle \hat{s}(\mathbf{r}_j, t)^2 \rangle = \alpha_j^{-1} \|\mathbf{l}(\mathbf{r}_j)\|^2, \quad (3.30)$$

The weight vector in this case is obtained as

$$\mathbf{w}(\mathbf{r}_j) = \frac{\mathbf{R}^{-1} \tilde{\mathbf{l}}(\mathbf{r}_j)}{[\tilde{\mathbf{l}}^T(\mathbf{r}_j) \mathbf{R}^{-1} \tilde{\mathbf{l}}(\mathbf{r}_j)]}, \quad (3.31)$$

which is equal to the weight expression of the array-gain constraint minimum-variance beamformer.

3.4 Diagonal-Loading and Bayesian Beamformers

The sample covariance matrix sometimes has a large condition number. This situation happens, for example, when the number of time samples is considerably fewer than the size of the matrix or when the SNR of the sensor data is very high. In such cases, since the direct inversion of the sample covariance \mathbf{R}^{-1} may cause numerical instability, the regularized inverse, $(\mathbf{R} + \kappa \mathbf{I})^{-1}$, may be a better approximation of the inverse of the model data covariance $\boldsymbol{\Sigma}_y^{-1}$, where κ is a small positive real number called the regularization constant. This gives the weight expression,

$$\mathbf{w}(\mathbf{r}) = \frac{(\mathbf{R} + \kappa \mathbf{I})^{-1} \mathbf{l}(\mathbf{r})}{[\mathbf{l}^T(\mathbf{r}) (\mathbf{R} + \kappa \mathbf{I})^{-1} \mathbf{l}(\mathbf{r})]}. \quad (3.32)$$

The beamformer that uses the above weight is called the diagonal-loading minimum-variance beamformer [10]. The diagonal loading is a form of regularization used when inverting the sample covariance matrix, and is nearly equivalent to adding noise to the sensor data. Contrary to the L_2 -regularization in the minimum-norm method, the diagonal loading is not needed when the SNR of the sensor data is low, because the condition number of the sample covariance is generally low in such cases.

The diagonal loading is needed when the sensor data has high SNR, because, in such cases, a significant SNR degradation is caused due to the problem of array mismatch [6], which indicates a situation where the lead field used for computing the beamformer's weight vector is different from the true lead field. When computing a weight vector, the exact lead field is generally unknown and the lead field is usually estimated from some kind of forward models such as the homogeneous conductor model. Thus, we cannot completely avoid this array mismatch problem. The diagonal loading can reduce the SNR degradation due to the array mismatch. On the other hand, however, it degrades the spatial resolution of reconstructed source distribution, so it provides a trade-off between the SNR degradation and spatial resolution [6].

The sample data covariance \mathbf{R} , which is the maximum-likelihood estimate of the model data covariance, is not necessarily the best estimate of Σ_y . We can obtain a better estimate by applying the Bayesian factor analysis described in Chap. 5. Applying the VBFA algorithm, we can get an estimate of the data covariance $\bar{\mathbf{R}}$, as shown in Eq. (5.104). The beamformer with a weight expression computed using a Bayesian-inferred data covariance matrix, (such as $\bar{\mathbf{R}}$ in Eq. (5.104)), is called the Bayesian beamformer [11]. Note that a Bayesian-inferred data covariance matrix has an intrinsic regularization term, and thus the regularization is embedded in the Bayesian beamformer.

3.5 Scalar Adaptive Beamformer with Unknown Source Orientation

3.5.1 Expressions for the Unit-Gain Constraint Beamformer

So far, we have derived the weight of the adaptive beamformer by assuming that the source orientation is predetermined. The source orientation may be predetermined using an accurate three-dimensional anatomical image of the subject, if available. However, in general, the source orientation $\eta(\mathbf{r})$ is an unknown quantity, and should be estimated from the data. There are two types of beamformers that can handle the estimation of the source vector. One is the scalar-type and the other is the vector-type adaptive beamformer. In the following, we first describe the scalar adaptive beamformer [5].

In the scalar-type adaptive beamformer, the weight is first formulated using the unknown source orientation η , such that [5]

$$\mathbf{w}(\mathbf{r}, \boldsymbol{\eta}) = \frac{\mathbf{R}^{-1}\mathbf{L}(\mathbf{r})\boldsymbol{\eta}}{[\boldsymbol{\eta}^T \mathbf{L}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r}) \boldsymbol{\eta}]}, \quad (3.33)$$

and the output power obtained using this weight is given by

$$\langle \hat{s}(\mathbf{r}, \boldsymbol{\eta})^2 \rangle = \frac{1}{\boldsymbol{\eta}^T [\mathbf{L}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r})] \boldsymbol{\eta}}. \quad (3.34)$$

The source orientation is then determined by maximizing this output power. That is, the optimum orientation $\boldsymbol{\eta}_{\text{opt}}(\mathbf{r})$ is derived as [5]

$$\boldsymbol{\eta}_{\text{opt}}(\mathbf{r}) = \underset{\boldsymbol{\eta}(\mathbf{r})}{\operatorname{argmax}} \left[\frac{1}{\boldsymbol{\eta}^T(\mathbf{r}) \mathbf{L}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r}) \boldsymbol{\eta}(\mathbf{r})} \right]. \quad (3.35)$$

According to the Rayleigh–Ritz formula in Sect.C.9, the orientation $\boldsymbol{\eta}_{\text{opt}}(\mathbf{r})$ is obtained as

$$\boldsymbol{\eta}_{\text{opt}}(\mathbf{r}) = \underset{\boldsymbol{\eta}(\mathbf{r})}{\operatorname{argmin}} \left[\boldsymbol{\eta}^T(\mathbf{r}) [\mathbf{L}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r})] \boldsymbol{\eta}(\mathbf{r}) \right] = \vartheta_{\min}\{\mathbf{L}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r})\}, \quad (3.36)$$

where $\vartheta_{\min}\{\cdot\}$ indicates the eigenvector corresponding to the minimum eigenvalue of the matrix in the curly braces.³ Namely, the optimum orientation $\boldsymbol{\eta}_{\text{opt}}(\mathbf{r})$ is given by the eigenvector corresponding to the minimum eigenvalue of $\mathbf{L}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r})$.

Once $\boldsymbol{\eta}_{\text{opt}}(\mathbf{r})$ is obtained, the explicit form of the weight vector for the scalar minimum-variance beamformer is expressed as

$$\mathbf{w}(\mathbf{r}) = \frac{\mathbf{R}^{-1}\mathbf{L}(\mathbf{r})\boldsymbol{\eta}_{\text{opt}}(\mathbf{r})}{[\boldsymbol{\eta}_{\text{opt}}^T(\mathbf{r}) \mathbf{L}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r}) \boldsymbol{\eta}_{\text{opt}}(\mathbf{r})]}. \quad (3.37)$$

The output power is given by

$$\langle \hat{s}(\mathbf{r}, t)^2 \rangle = \frac{1}{[\boldsymbol{\eta}_{\text{opt}}^T(\mathbf{r}) \mathbf{L}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r}) \boldsymbol{\eta}_{\text{opt}}(\mathbf{r})]} = \frac{1}{\mathcal{S}_{\min}\{\mathbf{L}^T(\mathbf{r}) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r})\}}, \quad (3.38)$$

where $\mathcal{S}_{\min}\{\cdot\}$ is the minimum eigenvalue of the matrix in the curly braces.

3.5.2 Expressions for the Array-Gain and Weight-Normalized Beamformers

For the scalar-type array-gain minimum-variance beamformer, the optimum orientation is derived using

³ The notations such as $\vartheta_{\min}\{\cdot\}$ and $\mathcal{S}_{\min}\{\cdot\}$ are defined in Sect.C.9 in the Appendix.

$$\eta_{\text{opt}}(\mathbf{r}) = \underset{\eta(\mathbf{r})}{\operatorname{argmax}} \left[\frac{\eta^T(\mathbf{r})[\mathbf{L}^T(\mathbf{r})\mathbf{L}(\mathbf{r})]\eta(\mathbf{r})}{\eta^T(\mathbf{r})[\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-1}\mathbf{L}(\mathbf{r})]\eta(\mathbf{r})} \right], \quad (3.39)$$

which can be computed using

$$\eta_{\text{opt}}(\mathbf{r}) = \vartheta_{\min}\{\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-1}\mathbf{L}(\mathbf{r}), \mathbf{L}^T(\mathbf{r})\mathbf{L}(\mathbf{r})\}. \quad (3.40)$$

Once $\eta_{\text{opt}}(\mathbf{r})$ is obtained, the weight vector is obtained using Eq. (3.10) with $\mathbf{l}(\mathbf{r}) = \mathbf{L}(\mathbf{r})\eta_{\text{opt}}(\mathbf{r})$ and $\tilde{l}(\mathbf{r}) = l(\mathbf{r})/\|l(\mathbf{r})\|$. The output power of this scalar beamformer is given by

$$\langle \hat{s}(\mathbf{r}, t)^2 \rangle = \frac{1}{S_{\min}\{\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-1}\mathbf{L}(\mathbf{r}), \mathbf{L}^T(\mathbf{r})\mathbf{L}(\mathbf{r})\}}. \quad (3.41)$$

For the scalar-type weight-normalized minimum-variance beamformer, the optimum orientation is derived using

$$\eta_{\text{opt}}(\mathbf{r}) = \underset{\eta(\mathbf{r})}{\operatorname{argmax}} \left[\frac{\eta^T(\mathbf{r})[\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-1}\mathbf{L}(\mathbf{r})]\eta(\mathbf{r})}{\eta^T(\mathbf{r})[\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-2}\mathbf{L}(\mathbf{r})]\eta(\mathbf{r})} \right], \quad (3.42)$$

which can be computed using

$$\eta_{\text{opt}}(\mathbf{r}) = \vartheta_{\min}\{\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-2}\mathbf{L}(\mathbf{r}), \mathbf{L}^T(\mathbf{r})\mathbf{R}^{-1}\mathbf{L}(\mathbf{r})\}. \quad (3.43)$$

Once $\eta_{\text{opt}}(\mathbf{r})$ is obtained, the weight vector is obtained using Eq. (3.17) with $\mathbf{l}(\mathbf{r}) = \mathbf{L}(\mathbf{r})\eta_{\text{opt}}(\mathbf{r})$. The output power of this scalar beamformer is given by

$$\langle \hat{s}(\mathbf{r}, t)^2 \rangle = \frac{1}{S_{\min}\{\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-2}\mathbf{L}(\mathbf{r}), \mathbf{L}^T(\mathbf{r})\mathbf{R}^{-1}\mathbf{L}(\mathbf{r})\}}. \quad (3.44)$$

3.6 Vector-Type Adaptive Beamformer

3.6.1 Vector Beamformer Formulation

The vector adaptive beamformer is another type of adaptive beamformers that can reconstruct the source orientation as well as the source magnitude. It uses a set of three weight vectors $\mathbf{w}_x(\mathbf{r})$, $\mathbf{w}_y(\mathbf{r})$, and $\mathbf{w}_z(\mathbf{r})$, which, respectively, detect the x , y , and z components of the source vector $\mathbf{s}(\mathbf{r}, t)$. That is, the weight matrix $\mathbf{W}(\mathbf{r})$ is defined as

$$\mathbf{W}(\mathbf{r}) = [\mathbf{w}_x(\mathbf{r}), \mathbf{w}_y(\mathbf{r}), \mathbf{w}_z(\mathbf{r})], \quad (3.45)$$

and the source vector can be estimated using

$$\hat{\mathbf{s}}(\mathbf{r}, t) = [\hat{s}_x(\mathbf{r}, t), \hat{s}_y(\mathbf{r}, t), \hat{s}_z(\mathbf{r}, t)]^T = \mathbf{W}^T(\mathbf{r})\mathbf{y}(t). \quad (3.46)$$

To derive the weight matrix of a vector-type minimum-variance beamformer, we use the optimization

$$\mathbf{W}(\mathbf{r}) = \underset{\mathbf{W}(\mathbf{r})}{\operatorname{argmin}} \operatorname{tr}[\mathbf{W}^T(\mathbf{r})\mathbf{R}\mathbf{W}(\mathbf{r})], \text{ subject to } \mathbf{W}^T(\mathbf{r})\mathbf{L}(\mathbf{r}) = \mathbf{I}. \quad (3.47)$$

A derivation similar to that for Eq. (3.6) leads to the solution for the weight matrix $\mathbf{W}(\mathbf{r})$, which is expressed as [4]

$$\mathbf{W}(\mathbf{r}) = \mathbf{R}^{-1}\mathbf{L}(\mathbf{r})[\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-1}\mathbf{L}(\mathbf{r})]^{-1}. \quad (3.48)$$

The beamformer that uses the above weight is called the vector-type adaptive beamformer. Using the weight matrix above, the source-vector covariance matrix is estimated as

$$\langle \hat{\mathbf{s}}(\mathbf{r}, t)\hat{\mathbf{s}}^T(\mathbf{r}, t) \rangle = [\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-1}\mathbf{L}(\mathbf{r})]^{-1}, \quad (3.49)$$

and the source power estimate is obtained using

$$\langle \|\hat{\mathbf{s}}(\mathbf{r}, t)\|^2 \rangle = \langle \hat{s}(\mathbf{r}, t)^2 \rangle = \operatorname{tr} \left[[\mathbf{L}^T(\mathbf{r})\mathbf{R}^{-1}\mathbf{L}(\mathbf{r})]^{-1} \right]. \quad (3.50)$$

This vector-type beamformer is sometimes referred to as the linearly constrained minimum-variance (LCMV) beamformer [4].

The weight matrix of the array-gain vector beamformer is derived, such that

$$\begin{aligned} \mathbf{W}(\mathbf{r}) &= \underset{\mathbf{W}(\mathbf{r})}{\operatorname{argmin}} \operatorname{tr} \left[\mathbf{W}^T(\mathbf{r})\mathbf{R}\mathbf{W}(\mathbf{r}) \right], \\ &\text{subject to } \mathbf{W}^T(\mathbf{r})\mathbf{L}(\mathbf{r}) = \|\mathbf{L}(\mathbf{r})\|. \end{aligned} \quad (3.51)$$

The vector version of the array-gain minimum-variance beamformer is expressed as

$$\mathbf{W}(\mathbf{r}) = \mathbf{R}^{-1}\tilde{\mathbf{L}}(\mathbf{r})[\tilde{\mathbf{L}}^T(\mathbf{r})\mathbf{R}^{-1}\tilde{\mathbf{L}}(\mathbf{r})]^{-1}, \quad (3.52)$$

where $\tilde{\mathbf{L}}(\mathbf{r})$ is the normalized lead field matrix defined as $\tilde{\mathbf{L}}(\mathbf{r}) = \mathbf{L}(\mathbf{r})/\|\mathbf{L}(\mathbf{r})\|$. The source-vector covariance matrix is estimated as

$$\langle \hat{\mathbf{s}}(\mathbf{r}, t)\hat{\mathbf{s}}^T(\mathbf{r}, t) \rangle = [\tilde{\mathbf{L}}^T(\mathbf{r})\mathbf{R}^{-1}\tilde{\mathbf{L}}(\mathbf{r})]^{-1}. \quad (3.53)$$

The source power estimate is obtained by

$$\langle \|\hat{s}(\mathbf{r}, t)\|^2 \rangle = \langle \hat{s}(\mathbf{r}, t)^2 \rangle = \text{tr} \left[[\tilde{\mathbf{L}}^T(\mathbf{r}) \mathbf{R}^{-1} \tilde{\mathbf{L}}(\mathbf{r})]^{-1} \right]. \quad (3.54)$$

The vector-type beamformer can be formulated with the unit-noise-gain constraint, and the detail of the formulation is found in [6, 12].

3.6.2 Semi-Bayesian Formulation

The weight matrix of a vector-type adaptive beamformer can be derived using the same Bayesian formulation in Sect. 3.3. The prior distribution in Eq. (3.20) and the noise assumption in Eq. (3.19) lead to the Bayesian estimate of \mathbf{x} in Eq. (3.21), which is written as

$$\bar{\mathbf{x}}(t) = \boldsymbol{\Phi}^{-1} \mathbf{F}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t), \quad (3.55)$$

where, $\bar{\mathbf{x}}(t)$ is expressed as

$$\bar{\mathbf{x}}(t) = \begin{bmatrix} \hat{s}(\mathbf{r}_1, t) \\ \hat{s}(\mathbf{r}_2, t) \\ \vdots \\ \hat{s}(\mathbf{r}_N, t) \end{bmatrix}. \quad (3.56)$$

To derive the vector beamformer, we use the matrix $\boldsymbol{\Phi}^{-1}$ which is a $3N \times 3N$ Block diagonal matrix such that

$$\boldsymbol{\Phi}^{-1} = \begin{bmatrix} \boldsymbol{\Upsilon}_1 & 0 & \cdots & 0 \\ 0 & \boldsymbol{\Upsilon}_2 & \cdot & \vdots \\ \vdots & \cdot & \ddots & 0 \\ 0 & \cdots & 0 & \boldsymbol{\Upsilon}_N \end{bmatrix}, \quad (3.57)$$

where the 3×3 matrix $\boldsymbol{\Upsilon}_j$ is the prior covariance matrix of the source vector at the j th voxel. Thus, the estimated source vector for the j th voxel is obtained as

$$\hat{s}(\mathbf{r}_j, t) = \boldsymbol{\Upsilon}_j \mathbf{L}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t). \quad (3.58)$$

Computing the estimated source-vector covariance matrix, $\langle \hat{s}(\mathbf{r}_j, t) \hat{s}^T(\mathbf{r}_j, t) \rangle$, leads to

$$\begin{aligned}\langle \hat{s}(\mathbf{r}_j, t) \hat{s}^T(\mathbf{r}_j, t) \rangle &= \boldsymbol{\Upsilon}_j \mathbf{L}_j^T \boldsymbol{\Sigma}_y^{-1} \langle \mathbf{y}(t) \mathbf{y}^T(t) \rangle \boldsymbol{\Sigma}_y^{-1} \mathbf{L}_j \boldsymbol{\Upsilon}_j \\ &= \boldsymbol{\Upsilon}_j \mathbf{L}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{L}_j \boldsymbol{\Upsilon}_j,\end{aligned}\quad (3.59)$$

where we assume $\langle \mathbf{y}(t) \mathbf{y}^T(t) \rangle = \boldsymbol{\Sigma}_y$.

In this vector case, the unit-gain constraint is expressed as

$$\langle \hat{s}(\mathbf{r}_j, t) \hat{s}^T(\mathbf{r}_j, t) \rangle = \boldsymbol{\Upsilon}_j. \quad (3.60)$$

Imposing this relationship, we get,

$$\boldsymbol{\Upsilon}_j = \boldsymbol{\Upsilon}_j \mathbf{L}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{L}_j \boldsymbol{\Upsilon}_j,$$

and

$$\boldsymbol{\Upsilon}_j = \left[\mathbf{L}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{L}_j \right]^{-1}. \quad (3.61)$$

Substituting the above equation into Eq. (3.58), we obtain

$$\hat{s}(\mathbf{r}_j, t) = \left[\mathbf{L}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{L}_j \right]^{-1} \mathbf{L}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t). \quad (3.62)$$

Assuming that the model data covariance can be replaced by the sample data covariance \mathbf{R} , Eq. (3.62) is rewritten as

$$\hat{s}(\mathbf{r}, t) = \mathbf{W}^T(\mathbf{r}_j) \mathbf{y}(t), \quad (3.63)$$

where the weight matrix is given by

$$\mathbf{W}(\mathbf{r}_j) = \mathbf{R}^{-1} \mathbf{L}(\mathbf{r}_j) \left[\mathbf{L}^T(\mathbf{r}_j) \mathbf{R}^{-1} \mathbf{L}(\mathbf{r}_j) \right]^{-1}. \quad (3.64)$$

The weight matrix of the vector array-gain minimum-variance beamformer can be derived using the array-gain constraint,

$$\langle \hat{s}(\mathbf{r}_j, t) \hat{s}^T(\mathbf{r}_j, t) \rangle = \boldsymbol{\Upsilon}_j \|\mathbf{L}(\mathbf{r}_j)\|. \quad (3.65)$$

Substituting this relationship into Eq. (3.59) and using Eq. (3.58), the resultant weight matrix is expressed as

$$\mathbf{W}(\mathbf{r}_j) = \mathbf{R}^{-1} \tilde{\mathbf{L}}(\mathbf{r}_j) \left[\tilde{\mathbf{L}}^T(\mathbf{r}_j) \mathbf{R}^{-1} \tilde{\mathbf{L}}(\mathbf{r}_j) \right]^{-1}. \quad (3.66)$$

3.7 Narrow-Band Beamformer

3.7.1 Background

Stimulus-induced power modulation of spontaneous brain activity has been the subject of intense investigations. Such power modulation is sometimes referred to as the event-related spectral power change. When the power change is negative, it is customarily termed as event-related desynchronization (ERD), and when it is positive, it is termed as event-related synchronization (ERS) [13].

The narrow-band dual-state beamformer [7] is a powerful tool for localization of specific brain activities related to these power changes. This is because the power change is frequency-specific, and the narrow-band beamformer uses a weight tuned to a specific target frequency. In this section, we describe the time-domain and the frequency-domain implementations of the narrow-band beamformer.

3.7.2 Time-Domain Implementation

Since the induced brain activity is not time-locked to the stimulus, the sample covariance matrix should be computed from nonaveraged raw trials. We assume that total N_E trials denoted as $\mathbf{b}_1(t), \dots, \mathbf{b}_{N_E}$ are obtained. To compute a frequency-specific weight, these raw trials are band-pass filtered with a specific frequency band of interest. The band-pass filtered sensor data for the n th trial is denoted as $\tilde{\mathbf{b}}_n(t, f)$, where f represents the frequency band of interest.⁴

The sample covariance matrix is computed using the band-pass filtered data, such that

$$\widehat{\mathbf{R}}(f) = \frac{1}{N_E} \sum_{n=1}^{N_E} \sum_{k=1}^K \tilde{\mathbf{b}}_n(t_k, f) \tilde{\mathbf{b}}_n^T(t_k, f). \quad (3.67)$$

Since we use the nonaveraged trial data, the data naturally contains a significant amount of influence from brain activities that are not related to the activity of interest.

To remove the influence of such unwanted brain activities, we use dual-state datasets: one from the target time period and the other from the baseline period. We try to reconstruct source activities that cause the power change in a specific frequency band. That is, using Eq. (3.67), we compute the frequency-specific covariance matrix for the target period, $\widehat{\mathbf{R}}_T(f)$, and for the baseline period, $\widehat{\mathbf{R}}_C(f)$.

We then reconstruct frequency-specific power changes between the target and baseline periods by using the F -ratio method. That is, using $\widehat{\mathbf{R}}_T(f)$ and $\widehat{\mathbf{R}}_C(f)$, the F -ratio method computes the frequency-specific filter weight such that

⁴ The notation f may indicate the center frequency of the frequency band of interest.

$$\mathbf{w}(\mathbf{r}, f) = \frac{\widehat{\mathbf{R}}_{\text{total}}^{-1}(f)\mathbf{l}(\mathbf{r})}{\mathbf{l}^T(\mathbf{r})\widehat{\mathbf{R}}_{\text{total}}^{-1}(f)\mathbf{l}(\mathbf{r})}, \quad (3.68)$$

where $\widehat{\mathbf{R}}_{\text{total}}(f) = \widehat{\mathbf{R}}_T(f) + \widehat{\mathbf{R}}_C(f)$. Then, the F -ratio image, $F(\mathbf{r}, f)$, is obtained such that

$$F(\mathbf{r}, f) = \frac{\mathbf{w}^T(\mathbf{r}, f)\widehat{\mathbf{R}}_T(f)\mathbf{w}(\mathbf{r}, f) - \mathbf{w}^T(\mathbf{r}, f)\widehat{\mathbf{R}}_C(f)\mathbf{w}(\mathbf{r}, f)}{\mathbf{w}^T(\mathbf{r}, f)\widehat{\mathbf{R}}_C(f)\mathbf{w}(\mathbf{r}, f)}. \quad (3.69)$$

On the right-hand side of Eq. (3.69), the first term in the numerator represents the reconstruction of the source power in the target period and the second represents that in the control period. Thus, $F(\mathbf{r}, f)$ represents the ratio of the reconstructed source power change to the power of baseline activities.

3.7.3 Frequency-Domain Implementation

The narrow-band beamformer can also be implemented in the frequency domain. We first define the Fourier transform of the raw-trial vector $\mathbf{b}_n(t)$ as $\mathbf{g}_n(f)$. The sample cross-spectrum matrix $\widetilde{\mathbf{R}}(f)$ is computed using $\mathbf{g}_n(f)$, such that

$$\widetilde{\mathbf{R}}(f) = \frac{1}{N_E} \sum_{n=1}^{N_E} \mathbf{g}_n(f) \mathbf{g}_n(f)^H, \quad (3.70)$$

where the superscript H indicates the Hermitian transpose (complex conjugation plus matrix transpose). Using Eq. (3.70), we compute the frequency-specific cross-spectral matrix for the target period, $\widetilde{\mathbf{R}}_T(f)$, and for the control period, $\widetilde{\mathbf{R}}_C(f)$. The frequency-selective weight $\mathbf{w}(\mathbf{r}, f)$ is obtained such that

$$\mathbf{w}(\mathbf{r}, f) = \frac{\widetilde{\mathbf{R}}_{\text{total}}^{-1}(f)\mathbf{l}(\mathbf{r})}{\mathbf{l}^T(\mathbf{r})\widetilde{\mathbf{R}}_{\text{total}}^{-1}(f)\mathbf{l}(\mathbf{r})}, \quad (3.71)$$

where $\widetilde{\mathbf{R}}_{\text{total}}(f) = \widetilde{\mathbf{R}}_T(f) + \widetilde{\mathbf{R}}_C(f)$. The pseudo F -ratio image is computed using

$$F(\mathbf{r}, f) = \frac{\mathbf{w}^H(\mathbf{r}, f)\widetilde{\mathbf{R}}_T(f)\mathbf{w}(\mathbf{r}, f) - \mathbf{w}^H(\mathbf{r}, f)\widetilde{\mathbf{R}}_C(f)\mathbf{w}(\mathbf{r}, f)}{\mathbf{w}^H(\mathbf{r}, f)\widetilde{\mathbf{R}}_C(f)\mathbf{w}(\mathbf{r}, f)}. \quad (3.72)$$

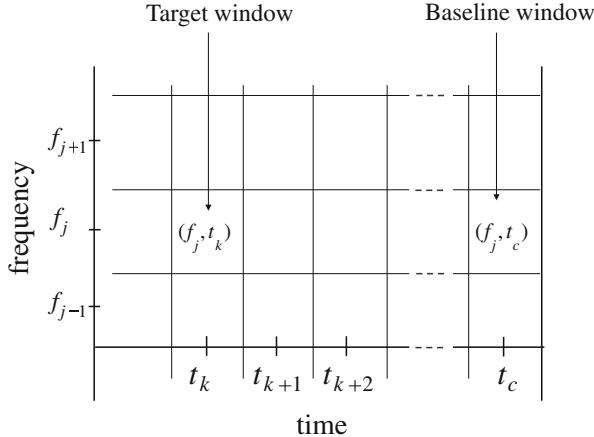


Fig. 3.1 Depiction of time–frequency domain discretization. The target window is set at the one represented by (f_j, t_k) , and the baseline window is set at the one represented by (f_j, t_c)

3.7.4 Five-Dimensional Brain Imaging

Using the narrow-band dual-state beamformer, we can implement the five-dimensional (time–frequency–space) imaging of brain activities [7]. In this implementation, we use the sliding window method in which the target window is moved along the time and frequency directions. The discretization of the time–frequency domain is depicted in Fig. 3.1.

We assign the window denoted (f_j, t_k) to the target time–frequency window, and the window denoted (f_j, t_c) to the baseline window. By implementing the narrow-band dual-state beamformer using these two time–frequency windows, we obtain the pseudo F-ratio image, $F(\mathbf{r}, f_j, t_k)$, which represents the source power difference at the frequency f_j between the time windows at t_k and t_c . If we move the target window along the time and frequency directions, we can obtain the five-dimensional source power difference map of the induced brain activity, $F(\mathbf{r}, f_j, t_k)$, where $j = 1, \dots, N_f$ and $k = 1, \dots, K$, where N_f is the number of frequency bins and K is the number of time windows.

3.8 Nonadaptive Spatial Filters

3.8.1 Minimum-Norm Filter

There is a different class of beamformers that uses the nonadaptive weight, the weight computed only using the sensor lead field. These are customarily called as

the nonadaptive spatial filters.⁵ A representative and basic nonadaptive spatial filter is the minimum-norm filter, which is the spatial filter version of the minimum-norm method described in Chap. 2.

The minimum-norm source reconstruction method is formulated as a nonadaptive spatial filter in the following manner. In Eq. (2.26), $\mathbf{F}\mathbf{F}^T$ is expressed such that [14]

$$\mathbf{F}\mathbf{F}^T = \sum_{n=1}^N \mathbf{L}(\mathbf{r}_n) \mathbf{L}^T(\mathbf{r}_n) \approx \int_{\Omega} \mathbf{L}(\mathbf{r}) \mathbf{L}^T(\mathbf{r}) d\mathbf{r} = \mathbf{G}, \quad (3.73)$$

where Ω indicates the source space. That is, $\mathbf{F}\mathbf{F}^T$ is equal to the gram matrix \mathbf{G} , if we ignore the voxel discretization error. Thus, let us rewrite Eq. (2.26) as

$$\begin{bmatrix} \hat{\mathbf{s}}(\mathbf{r}_1, t) \\ \hat{\mathbf{s}}(\mathbf{r}_2, t) \\ \vdots \\ \hat{\mathbf{s}}(\mathbf{r}_N, t) \end{bmatrix} = \mathbf{F}^T \mathbf{G}^{-1} \mathbf{y}(t) = \begin{bmatrix} \mathbf{L}^T(\mathbf{r}_1) \\ \mathbf{L}^T(\mathbf{r}_2) \\ \vdots \\ \mathbf{L}^T(\mathbf{r}_N) \end{bmatrix} \mathbf{G}^{-1} \mathbf{y}(t). \quad (3.74)$$

We can see that, at each voxel location \mathbf{r}_n , the relationship,

$$\hat{\mathbf{s}}(\mathbf{r}_n, t) = \mathbf{L}^T(\mathbf{r}_n) \mathbf{G}^{-1} \mathbf{y}(t), \quad (3.75)$$

holds.

The equation above has the same form as the vector-type beamformer in Eq. (3.46), i.e., Eq. (3.75) is rewritten as

$$\hat{\mathbf{s}}(\mathbf{r}, t) = \mathbf{W}^T(\mathbf{r}) \mathbf{y}(t), \quad (3.76)$$

where the weight matrix $\mathbf{W}(\mathbf{r})$ is given by:

$$\mathbf{W}(\mathbf{r}) = \mathbf{G}^{-1} \mathbf{L}(\mathbf{r}). \quad (3.77)$$

The two equations above indicate that the minimum-norm method is formulated as a nonadaptive spatial filter, in which the weight matrix is given by Eq. (3.77). Also, the L_2 -regularized version of the weight matrix is given by:

$$\mathbf{W}(\mathbf{r}) = (\mathbf{G} + \xi \mathbf{I})^{-1} \mathbf{L}(\mathbf{r}), \quad (3.78)$$

where a scalar ξ is the regularization constant. The spatial filter whose weight is expressed in either Eq. (3.77) or (3.78) is called the minimum-norm filter.

⁵ They might be also called as the nonadaptive beamformers, but this usage is uncommon.

3.8.2 Weight-Normalized Minimum-Norm Filter

The weight-normalized minimum-norm filter has been proposed by Dale et al. and the method is often called as dynamic statistical parametric mapping (dSPM) [15]. The idea is to normalize the minimum-norm weight with its weight's norm to ensure that the spatial distribution of the noise is uniform. The scalar-type weight is thus given by

$$\mathbf{w}(\mathbf{r}) = \frac{\mathbf{G}^{-1}\mathbf{l}(\mathbf{r})}{\|\mathbf{G}^{-1}\mathbf{l}(\mathbf{r})\|} = \frac{\mathbf{G}^{-1}\mathbf{l}(\mathbf{r})}{\sqrt{\mathbf{l}^T(\mathbf{r})\mathbf{G}^{-2}\mathbf{l}(\mathbf{r})}}. \quad (3.79)$$

The idea of weight normalization can be extended to derive the weight matrix for the vector-type spatial filter, such that

$$\mathbf{W}(\mathbf{r}) = \frac{\mathbf{G}^{-1}\mathbf{L}(\mathbf{r})}{\sqrt{\text{tr}[\mathbf{L}^T(\mathbf{r})\mathbf{G}^{-2}\mathbf{L}(\mathbf{r})]}}. \quad (3.80)$$

Using this weight matrix, the source vector is estimated as

$$\hat{\mathbf{s}}(\mathbf{r}, t) = \mathbf{W}^T(\mathbf{r})\mathbf{y}(t) = \frac{\mathbf{L}^T(\mathbf{r})\mathbf{G}^{-1}\mathbf{y}(t)}{\sqrt{\text{tr}[\mathbf{L}^T(\mathbf{r})\mathbf{G}^{-2}\mathbf{L}(\mathbf{r})]}}. \quad (3.81)$$

The weight matrix for the L_2 -regularized version is given by

$$\mathbf{W}(\mathbf{r}) = \frac{(\mathbf{G} + \xi\mathbf{I})^{-1}\mathbf{L}(\mathbf{r})}{\sqrt{\text{tr}[\mathbf{L}^T(\mathbf{r})(\mathbf{G} + \xi\mathbf{I})^{-2}\mathbf{L}(\mathbf{r})]}}. \quad (3.82)$$

3.8.3 sLORETA Filter

Standardized low resolution electromagnetic tomography (sLORETA) was originally proposed by Pascual-Marqui [16]. The method can be reformulated as a nonadaptive spatial filter. In this method, the minimum-norm filter outputs are normalized by the quantity $\sqrt{\mathbf{l}^T(\mathbf{r})\mathbf{G}^{-1}\mathbf{l}(\mathbf{r})}$. This normalization is called as the standardization. The scalar-type weight is given by

$$\mathbf{w}(\mathbf{r}) = \frac{\mathbf{l}\mathbf{G}^{-1}}{\sqrt{\mathbf{l}^T(\mathbf{r})\mathbf{G}^{-1}\mathbf{l}(\mathbf{r})}}. \quad (3.83)$$

The extension to the vector-type sLORETA filter results in the weight matrix expressed as

$$\mathbf{W}(\mathbf{r}) = \mathbf{G}^{-1}\mathbf{L}(\mathbf{r})[\mathbf{L}^T(\mathbf{r})\mathbf{G}^{-1}\mathbf{L}(\mathbf{r})]^{-1/2}, \quad (3.84)$$

and the source vector $\hat{s}(\mathbf{r}, t)$ is estimated as

$$\hat{s}(\mathbf{r}, t) = \mathbf{W}^T(\mathbf{r})\mathbf{y}(t) = [\mathbf{L}^T(\mathbf{r})\mathbf{G}^{-1}\mathbf{L}(\mathbf{r})]^{-1/2}\mathbf{L}^T(\mathbf{r})\mathbf{G}^{-1}\mathbf{y}(t). \quad (3.85)$$

The weight matrix for the L_2 -regularized version is given by,

$$\mathbf{W}(\mathbf{r}) = (\mathbf{G} + \xi\mathbf{I})^{-1}\mathbf{L}(\mathbf{r})[\mathbf{L}^T(\mathbf{r})(\mathbf{G} + \xi\mathbf{I})^{-1}\mathbf{L}(\mathbf{r})]^{-1/2}. \quad (3.86)$$

Although the standardization makes the sLORETA filter to have no localization bias [6], the theoretical basis of this standardization is not entirely clear.

3.9 Recursive Null-Steering (RENS) Beamformer

3.9.1 Beamformer Obtained Based on Beam-Response Optimization

In the beamformer formulation, the product, $\mathbf{w}^T(\mathbf{r})\mathbf{l}(\mathbf{r}')$, is called the beam response, which expresses the sensitivity of the beamformer pointing at \mathbf{r} to a source located at \mathbf{r}' . We wish to derive a beamformer that only passes the signal from a source at the pointing location and suppresses the leakage from sources at other locations. Such a beamformer may be derived by imposing a delta-function-like property on the beam response. That is, the weight vector is obtained using

$$\mathbf{w}(\mathbf{r}) = \underset{\mathbf{w}(\mathbf{r})}{\operatorname{argmin}} \int_{\Omega} \left[\mathbf{w}^T(\mathbf{r})\mathbf{l}(\mathbf{r}') - \delta(\mathbf{r} - \mathbf{r}') \right]^2 d\mathbf{r}'. \quad (3.87)$$

The weight obtained from the optimization above is expressed as

$$\mathbf{w}(\mathbf{r}) = \mathbf{G}^{-1}\mathbf{l}(\mathbf{r}). \quad (3.88)$$

The beamformer in Eq. (3.88) is exactly equal to the minimum-norm filter.

Variants of the minimum-norm filter can be obtained by adding various constraints to the optimization in Eq. (3.87) [17]. For example, the optimization

$$\begin{aligned} \mathbf{w}(\mathbf{r}) = \underset{\mathbf{w}(\mathbf{r})}{\operatorname{argmin}} \int_{\Omega} & \left[\mathbf{w}^T(\mathbf{r})\mathbf{l}(\mathbf{r}') - \delta(\mathbf{r} - \mathbf{r}') \right]^2 d\mathbf{r}' \\ & \text{subject to } \mathbf{w}^T(\mathbf{r})\mathbf{l}(\mathbf{r}) = 1, \end{aligned} \quad (3.89)$$

leads to the weight expression

$$\mathbf{w}(\mathbf{r}) = [\mathbf{l}^T(\mathbf{r})\mathbf{G}^{-1}\mathbf{l}(\mathbf{r})]^{-1}\mathbf{G}^{-1}\mathbf{l}(\mathbf{r}). \quad (3.90)$$

The beamformer using the above weight is called the unit-gain (constraint) minimum-norm filter. Using the optimization

$$\begin{aligned} \mathbf{w}(\mathbf{r}) = \operatorname{argmin}_{\mathbf{w}(\mathbf{r})} \int_{\Omega} \left[\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}') - \delta(\mathbf{r} - \mathbf{r}') \right]^2 d\mathbf{r}' \\ \text{subject to } \mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = \|\mathbf{l}(\mathbf{r})\|, \end{aligned} \quad (3.91)$$

we obtain the array-gain (constraint) minimum-norm filter, such that

$$\mathbf{w}(\mathbf{r}) = [\tilde{\mathbf{l}}^T(\mathbf{r}) \mathbf{G}^{-1} \tilde{\mathbf{l}}(\mathbf{r})]^{-1} \mathbf{G}^{-1} \tilde{\mathbf{l}}(\mathbf{r}), \quad (3.92)$$

where $\tilde{\mathbf{l}}(\mathbf{r}) = \mathbf{l}(\mathbf{r}) / \|\mathbf{l}(\mathbf{r})\|$.

3.9.2 Derivation of RENS Beamformer

The recursive null-steering (RENS) beamformer is derived from the following optimization [18]:

$$\begin{aligned} \mathbf{w}(\mathbf{r}) = \operatorname{argmin}_{\mathbf{w}(\mathbf{r})} \int_{\Omega} \left[\mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}') - \delta(\mathbf{r} - \mathbf{r}') \right]^2 s(\mathbf{r}, t)^2 d\mathbf{r}' \\ \text{subject to } \mathbf{w}^T(\mathbf{r}) \mathbf{l}(\mathbf{r}) = \|\mathbf{l}(\mathbf{r})\|, \end{aligned} \quad (3.93)$$

where $s(\mathbf{r}, t)^2$ is the instantaneous source power. The idea behind the above optimization is that we wish to impose a delta-function-like property on the beam response only over a region where sources exist, instead of imposing that property over the entire source space.

The filter weight is obtained as

$$\mathbf{w}(\mathbf{r}) = [\tilde{\mathbf{l}}^T(\mathbf{r}) \bar{\mathbf{G}}^{-1} \tilde{\mathbf{l}}(\mathbf{r})]^{-1} \bar{\mathbf{G}}^{-1} \tilde{\mathbf{l}}(\mathbf{r}), \quad (3.94)$$

where

$$\bar{\mathbf{G}} = \int_{\Omega} s(\mathbf{r}, t)^2 \mathbf{l}(\mathbf{r}) \mathbf{l}^T(\mathbf{r}) d\mathbf{r}. \quad (3.95)$$

However, to compute the weight in Eq. (3.94), we need to know the source magnitude distribution $s(\mathbf{r}, t)^2$, which is unknown. Therefore, we use an estimated source magnitude $\hat{s}(\mathbf{r}, t)^2$ when computing $\bar{\mathbf{G}}$, i.e.,

$$\bar{\mathbf{G}} = \int_{\Omega} \hat{s}(\mathbf{r}, t)^2 \mathbf{l}(\mathbf{r}) \mathbf{l}^T(\mathbf{r}) d\mathbf{r}. \quad (3.96)$$

The proposed weight is therefore derived in a recursive manner. That is, by setting an initial $\hat{s}(\mathbf{r}, t)^2$ to have a uniform value, the weight $\mathbf{w}(\mathbf{r})$ is derived using Eqs. (3.94) and (3.96). The estimated source intensity $\hat{s}(\mathbf{r}, t)$ is obtained using Eq. (3.1). This $\hat{s}(\mathbf{r}, t)$ is then used in Eqs. (3.94) and (3.96) to update $\mathbf{w}(\mathbf{r})$. These procedures are repeated until some stopping criterion is satisfied. Here we describe a derivation of the scalar-type RENS beamformer. An extension for deriving the vector-type RENS beamformer is straightforward [18].

The RENS beamformer described here provides spatial resolution higher than that of the nonadaptive spatial filters described in Sect. 3.8. Moreover, it is free from the limitations of adaptive beamformers. That is, the RENS beamformer is robust to the source correlation problem. It does not require large time samples and works even with single time point data.

References

1. J. Capon, High-resolution frequency wavenumber spectrum analysis. Proc. IEEE **57**, 1408–1419 (1969)
2. S.E. Robinson, D.F. Rose, Current source image estimation by spatially filtered MEG, in *Biomagnetism Clinical Aspects*, ed. by M. Hoke, et al. (Elsevier Science Publishers, New York, 1992), pp. 761–765
3. M.E. Spencer, R.M. Leahy, J.C. Mosher, P.S. Lewis, Adaptive filters for monitoring localized brain activity from surface potential time series, in *Conference Record for 26th Annual Asilomar Conference on Signals, Systems, and Computers*, pp. 156–161, November 1992
4. B.D. Van Veen, W. Van Drongelen, M. Yuchtman, A. Suzuki, Localization of brain electrical activity via linearly constrained minimum variance spatial filtering. IEEE Trans. Biomed. Eng. **44**, 867–880 (1997)
5. K. Sekihara, B. Scholz, Generalized Wiener estimation of three-dimensional current distribution from biomagnetic measurements, in *Biomag 96: Proceedings of the Tenth International Conference on Biomagnetism*, ed. by C.J. Aine et al. (Springer, New York, 1996), pp. 338–341
6. K. Sekihara, S.S. Nagarajan, *Adaptive Spatial Filters for Electromagnetic Brain Imaging* (Springer, Berlin, 2008)
7. S.S. Dalal, A.G. Guggisberg, E. Edwards, K. Sekihara, A.M. Findlay, R.T. Canolty, M.S. Berger, R.T. Knight, N.M. Barbaro, H.E. Kirsch, S.S. Nagarajan, Five-dimensional neuroimaging: localization of the time-frequency dynamics of cortical activity. NeuroImage **40**, 1686–1700 (2008)
8. G. Borgiotti, L.J. Kaplan, Superresolution of uncorrelated interference sources by using adaptive array technique. IEEE Trans. Antennas Propag. **27**, 842–845 (1979)
9. K. Sekihara, B. Scholz, Generalized Wiener estimation of three-dimensional current distribution from biomagnetic measurements. IEEE Trans. Biomed. Eng. **43**, 281–291 (1996)
10. H. Cox, R.M. Zeskind, M.M. Owen, Robust adaptive beamforming. IEEE Trans. Signal Process. **35**, 1365–1376 (1987)
11. M. Woolrich, L. Hunt, A. Groves, G. Barnes, MEG beamforming using Bayesian PCA for adaptive data covariance matrix regularization. NeuroImage **57**(4), 1466–1479 (2011)
12. K. Sekihara, S.S. Nagarajan, D. Poeppel, A. Marantz, Y. Miyashita, Reconstructing spatio-temporal activities of neural sources using an MEG vector beamformer technique. IEEE Trans. Biomed. Eng. **48**, 760–771 (2001)
13. G. Pfurtscheller, F.H. Lopes da Silva, Event-related EEG/MEG synchronization and desynchronization: basic principles. Clin. Neurophysiol. **110**, 1842–1857 (1999)

14. K. Sekihara, M. Sahani, S.S. Nagarajan, Location bias and spatial resolution of adaptive and non-adaptive spatial filters for MEG source reconstruction. *NeuroImage* **25**, 1056–1067 (2005)
15. A.M. Dale, A.K. Liu, B.R. Fischl, R.L. Buckner, J.W. Belliveau, J.D. Lewine, E. Halgren, Dynamic statistical parametric mapping: combining fMRI and MEG for high-resolution imaging of cortical activity. *Neuron* **26**, 55–67 (2000)
16. R.D. Pascual-Marqui, Standardized low resolution brain electromagnetic tomography (sLORETA): technical details. *Methods. Find. Exp. Clin. Pharmacol.* **24**, 5–12 (2002)
17. R.E. Greenblatt, A. Ossadtchi, M.E. Pfleger, Local linear estimators for the bioelectromagnetic inverse problem. *IEEE Trans. Signal Process.* **53**, 3403–3412 (2005)
18. I. Kumihashi, K. Sekihara, Array-gain constraint minimum-norm spatial filter with recursively updated gram matrix for biomagnetic source imaging. *IEEE Trans. Biomed. Eng.* **57**(6), 1358–1365 (2010)

Chapter 4

Sparse Bayesian (Champagne) Algorithm

4.1 Introduction

In this chapter, we provide a detailed description of an algorithm for electromagnetic brain imaging, called the Champagne algorithm [1, 2]. The Champagne algorithm is formulated based on an empirical Bayesian schema, and can provide a sparse solution, since the sparsity constraint is embedded in the algorithm. The algorithm is free from the problems that cannot be avoided in other sparse-solution methods, such as the L_1 -regularized minimum-norm method. Such problems include the difficulty in reconstructing voxel time courses or the difficulty in incorporating the source-orientation estimation.

In Sect. 2.10.2, we show that the L_2 -regularized minimum-norm method is derived using the Gaussian prior for the j th voxel value,¹

$$x_j \sim \mathcal{N}(x_j|0, \alpha^{-1}), \quad (4.1)$$

where the precision α is common to all x_j . In this chapter, we use the Gaussian prior whose precision (variance) is specific to each x_j , i.e.,

$$x_j \sim \mathcal{N}(x_j|0, \alpha_j^{-1}). \quad (4.2)$$

We show that this “slightly different” prior distribution gives a solution totally different from the L_2 -norm solution. Actually, the prior distribution in Eq. (4.2) leads to a sparse solution. The estimation method based on the prior in Eq. (4.2) is called the sparse Bayesian learning in the field of machine learning [3, 4], and the source reconstruction algorithm derived using Eq. (4.2) is called the Champagne algorithm [1].

In this chapter, we formulate the source reconstruction problem as the spatiotemporal reconstruction, i.e., the voxel time series $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$ is reconstructed using the sensor time series $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K$ where $\mathbf{y}(t_k)$ and $\mathbf{x}(t_k)$ are denoted \mathbf{y}_k and \mathbf{x}_k . We use collective expressions \mathbf{x} and \mathbf{y} , indicating the whole voxel time series

¹ We use the notational convenience $\mathcal{N}(\text{variable}|\text{mean, covariance matrix})$ throughout this book.

$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$ and the whole sensor time series $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K$. We first formulate the Champagne algorithm omitting the source orientation estimation. That is, we assume that the source orientation is predetermined at each voxel and use the measurement model in Eq. (2.15). The relationship between \mathbf{x}_k and \mathbf{y}_k is expressed as

$$\mathbf{y}_k = \mathbf{H}\mathbf{x}_k + \boldsymbol{\varepsilon}, \quad (4.3)$$

where the lead field matrix \mathbf{H} is defined in Eq. (2.11).

4.2 Probabilistic Model and Method Formulation

Defining a column vector $\boldsymbol{\alpha}$ such that $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]^T$, the prior distribution is expressed as

$$p(\mathbf{x}|\boldsymbol{\alpha}) = \prod_{k=1}^K p(\mathbf{x}_k|\boldsymbol{\alpha}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}_k|\mathbf{0}, \boldsymbol{\Phi}^{-1}), \quad (4.4)$$

where $\boldsymbol{\Phi}$ is equal to $\boldsymbol{\Phi} = \text{diag}(\boldsymbol{\alpha})$, which indicates a diagonal matrix whose diagonal elements are those of a vector in the parenthesis. Thus, we have

$$p(\mathbf{x}_k|\boldsymbol{\alpha}) = \mathcal{N}(\mathbf{x}_k|\mathbf{0}, \boldsymbol{\Phi}^{-1}) = \prod_{j=1}^N \mathcal{N}(x_j|0, \alpha_j^{-1}). \quad (4.5)$$

Since we assume that the noise $\boldsymbol{\varepsilon}$ is independent across time, the conditional probability $p(\mathbf{y}|\mathbf{x})$ is expressed as

$$p(\mathbf{y}|\mathbf{x}) = \prod_{k=1}^K p(\mathbf{y}_k|\mathbf{x}_k) = \prod_{k=1}^K \mathcal{N}(\mathbf{y}_k|\mathbf{H}\mathbf{x}_k, \beta^{-1}\mathbf{I}). \quad (4.6)$$

In this chapter, the noise precision matrix is assumed to be $\beta\mathbf{I}$ where β is known. Therefore, the parameters we must estimate is the voxel source distribution \mathbf{x} and the hyperparameter $\boldsymbol{\alpha}$.

In truly Bayesian formulation, we have to derive the joint posterior distribution of all the unknown parameters, $p(\mathbf{x}, \boldsymbol{\alpha}|\mathbf{y})$, using

$$p(\mathbf{x}, \boldsymbol{\alpha}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x}, \boldsymbol{\alpha})p(\mathbf{x}, \boldsymbol{\alpha})}{p(\mathbf{y})}. \quad (4.7)$$

However, we cannot compute²

$$p(\mathbf{y}) = \iint p(\mathbf{y}|\mathbf{x}, \boldsymbol{\alpha})p(\mathbf{x}, \boldsymbol{\alpha})d\mathbf{x}d\boldsymbol{\alpha}, \quad (4.8)$$

² In Eq. (4.8), the notation $d\mathbf{x}$ indicates $d\mathbf{x}_1 d\mathbf{x}_2 \dots d\mathbf{x}_K$.

because the integral on the right-hand side does not have a closed-form solution. Accordingly, the explicit form of $p(\mathbf{y})$, and thus the joint distribution $p(\mathbf{x}, \boldsymbol{\alpha}|\mathbf{y})$ cannot be obtained.

Therefore, instead, using

$$p(\mathbf{x}, \boldsymbol{\alpha}|\mathbf{y}) = p(\mathbf{x}|\mathbf{y}, \boldsymbol{\alpha})p(\boldsymbol{\alpha}|\mathbf{y}), \quad (4.9)$$

we can derive the posterior distribution $p(\mathbf{x}|\mathbf{y})$, such that

$$p(\mathbf{x}|\mathbf{y}) = \int p(\mathbf{x}, \boldsymbol{\alpha}|\mathbf{y})d\boldsymbol{\alpha} = \int p(\mathbf{x}|\mathbf{y}, \boldsymbol{\alpha})p(\boldsymbol{\alpha}|\mathbf{y})d\boldsymbol{\alpha}, \quad (4.10)$$

where we eliminate the unknown hyperparameter $\boldsymbol{\alpha}$ by marginalization. Denoting $\widehat{\boldsymbol{\alpha}}$ such that $\widehat{\boldsymbol{\alpha}} = \operatorname{argmax}_{\boldsymbol{\alpha}} p(\boldsymbol{\alpha}|\mathbf{y})$, and assuming that $p(\boldsymbol{\alpha}|\mathbf{y})$ has a (hopefully sharp) peak at $\widehat{\boldsymbol{\alpha}}$, we can use the approximation

$$p(\mathbf{x}|\mathbf{y}) = \int p(\mathbf{x}|\mathbf{y}, \boldsymbol{\alpha})p(\boldsymbol{\alpha}|\mathbf{y})d\boldsymbol{\alpha} \approx \int p(\mathbf{x}|\mathbf{y}, \boldsymbol{\alpha})\delta(\boldsymbol{\alpha} - \widehat{\boldsymbol{\alpha}})d\boldsymbol{\alpha} = p(\mathbf{x}|\mathbf{y}, \widehat{\boldsymbol{\alpha}}), \quad (4.11)$$

where $\delta(\boldsymbol{\alpha})$ indicates the delta function. The above equation shows that $p(\mathbf{x}|\mathbf{y}, \widehat{\boldsymbol{\alpha}})$, (the posterior distribution obtained with setting $\boldsymbol{\alpha}$ to $\widehat{\boldsymbol{\alpha}}$), approximates the true posterior distribution $p(\mathbf{x}|\mathbf{y})$.

Using the arguments in Sect. B.3 in the Appendix, when $p(\mathbf{x}|\boldsymbol{\alpha})$ and $p(\mathbf{y}|\mathbf{x})$ are expressed in Eqs. (4.4) and (4.6), the posterior distribution $p(\mathbf{x}|\mathbf{y}, \boldsymbol{\alpha})$ is given by

$$p(\mathbf{x}|\mathbf{y}, \boldsymbol{\alpha}) = \prod_{k=1}^K p(\mathbf{x}_k|\mathbf{y}_k, \boldsymbol{\alpha}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}_k|\bar{\mathbf{x}}_k, \boldsymbol{\Gamma}^{-1}), \quad (4.12)$$

where the precision and the mean are obtained as

$$\boldsymbol{\Gamma} = \boldsymbol{\Phi} + \beta \mathbf{H}^T \mathbf{H}, \quad (4.13)$$

$$\bar{\mathbf{x}}_k = \beta \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \mathbf{y}_k. \quad (4.14)$$

Therefore, once $\widehat{\boldsymbol{\alpha}}$ is obtained, we can compute $\boldsymbol{\Gamma}$ and $\bar{\mathbf{x}}_k$ by substituting $\boldsymbol{\Phi} = \operatorname{diag}(\widehat{\boldsymbol{\alpha}})$ into Eqs. (4.13) and (4.14), and we obtain $p(\mathbf{x}|\mathbf{y}, \widehat{\boldsymbol{\alpha}})$ using Eq. (4.12).

The problem here is how to estimate $\widehat{\boldsymbol{\alpha}}$. The posterior distribution of the hyperparameter $\boldsymbol{\alpha}$, $p(\boldsymbol{\alpha}|\mathbf{y})$, is expressed using the Bayes' rule,

$$p(\boldsymbol{\alpha}|\mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\alpha})p(\boldsymbol{\alpha}). \quad (4.15)$$

When we assume the flat (noninformative) prior for $p(\boldsymbol{\alpha})$, we have

$$p(\boldsymbol{\alpha}|\mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\alpha}). \quad (4.16)$$

Thus, since the relationship,

$$\operatorname{argmax}_{\boldsymbol{\alpha}} p(\boldsymbol{\alpha}|\mathbf{y}) = \operatorname{argmax}_{\boldsymbol{\alpha}} p(\mathbf{y}|\boldsymbol{\alpha})$$

holds, $\hat{\alpha}$ is obtained as the one that maximizes $p(\mathbf{y}|\alpha)$. This $p(\mathbf{y}|\alpha)$ is referred to as the data evidence or the marginal likelihood.

Let us summarize the procedure to estimate the source distribution \mathbf{x} . First, we estimate the hyperparameter α by maximizing the marginal likelihood function,

$$\hat{\alpha} = \underset{\alpha}{\operatorname{argmax}} p(\mathbf{y}|\alpha).$$

Next, this $\hat{\alpha}$ is substituted into the posterior distribution $p(\mathbf{x}|\mathbf{y}, \alpha)$ to obtain $p(\mathbf{x}|\mathbf{y}, \hat{\alpha})$. When this posterior is the Gaussian distribution in Eq.(4.12), the precision and mean are obtained by substituting $\Phi = \text{diag}(\hat{\alpha})$ into Eqs.(4.13) and (4.14). The voxel time courses are reconstructed by computing $\bar{\mathbf{x}}_k$ in Eq.(4.14) for $k = 1, \dots, K$.

4.3 Cost Function for Marginal Likelihood Maximization

As described in the preceding section, the hyperparameter α is estimated by maximizing the marginal likelihood $p(\mathbf{y}|\alpha)$. In this section, we describe the maximization of the marginal likelihood, and to do so, let us derive an explicit form of the log marginal likelihood, $\log p(\mathbf{y}|\alpha)$. Substituting³

$$p(\mathbf{x}|\alpha) = \prod_{k=1}^K p(\mathbf{x}_k|\alpha) = \left[\frac{|\Phi|^{1/2}}{(2\pi)^{N/2}} \right]^K \exp \left[-\frac{1}{2} \sum_{k=1}^K \mathbf{x}_k^T \Phi \mathbf{x}_k \right], \quad (4.17)$$

and

$$p(\mathbf{y}|\mathbf{x}) = \prod_{k=1}^K p(\mathbf{y}_k|\mathbf{x}_k) = \left[\left(\frac{\beta}{2\pi} \right)^{\frac{M}{2}} \right]^K \exp \left[-\frac{\beta}{2} \sum_{k=1}^K \|\mathbf{y}_k - \mathbf{H}\mathbf{x}_k\|^2 \right], \quad (4.18)$$

into

$$p(\mathbf{y}|\alpha) = \int p(\mathbf{y}, \mathbf{x}|\alpha) d\mathbf{x} = \int p(\mathbf{y}|\mathbf{x}) p(\mathbf{x}|\alpha) d\mathbf{x},$$

we obtain

$$p(\mathbf{y}|\alpha) = \left[\frac{|\Phi|^{1/2}}{(2\pi)^{N/2}} \left(\frac{\beta}{2\pi} \right)^{M/2} \right]^K \int \exp[-D] d\mathbf{x}, \quad (4.19)$$

where

$$D = \frac{\beta}{2} \sum_{k=1}^K \|\mathbf{y}_k - \mathbf{H}\mathbf{x}_k\|^2 + \frac{1}{2} \sum_{k=1}^K \mathbf{x}_k^T \Phi \mathbf{x}_k, \quad (4.20)$$

³ Note that M and N are respectively the sizes of \mathbf{y}_k and \mathbf{x}_k .

and this D can be rewritten as

$$D = \frac{1}{2} \sum_{k=1}^K (\mathbf{x}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Gamma} (\mathbf{x}_k - \bar{\mathbf{x}}_k) + \Delta. \quad (4.21)$$

The derivation of the above equation is presented in Sect. 4.10.1, and Δ on the right-hand side of Eq. (4.21) is given in Eq. (4.97).

Using the results in Eq. (4.21), let us compute the integral on the right-hand side of Eq. (4.19). First, we have

$$\int \exp[-D] d\mathbf{x} = \exp[-\Delta] \int \exp \left[-\frac{1}{2} \sum_{k=1}^K (\mathbf{x}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Gamma} (\mathbf{x}_k - \bar{\mathbf{x}}_k) \right] d\mathbf{x}. \quad (4.22)$$

Considering the Gaussian with its mean $\bar{\mathbf{x}}_k$ and precision $\boldsymbol{\Gamma}$, the integral on the right-hand side is computed such that,

$$\int \exp \left[-\frac{1}{2} (\mathbf{x}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Gamma} (\mathbf{x}_k - \bar{\mathbf{x}}_k) \right] d\mathbf{x}_k = \frac{(2\pi)^{N/2}}{|\boldsymbol{\Gamma}|^{1/2}}.$$

Substituting the above results into Eqs. (4.22) and (4.19), we have

$$p(\mathbf{y}|\boldsymbol{\alpha}) = \left[\frac{|\boldsymbol{\Phi}|^{1/2}}{(2\pi)^{N/2}} \left(\frac{\beta}{2\pi} \right)^{M/2} \right]^K \exp[-\Delta] \left[\frac{(2\pi)^{N/2}}{|\boldsymbol{\Gamma}|^{1/2}} \right]^K. \quad (4.23)$$

Taking the logarithm of both sides, and omitting constant terms, which are unrelated to the arguments, we get

$$\log p(\mathbf{y}|\boldsymbol{\alpha}) = K \left(-\frac{1}{2} \log |\boldsymbol{\Gamma}| + \frac{1}{2} \log |\boldsymbol{\Phi}| + \frac{M}{2} \log \beta \right) - \Delta. \quad (4.24)$$

Let us define $\boldsymbol{\Sigma}_y$ such that

$$\boldsymbol{\Sigma}_y = \beta^{-1} \mathbf{I} + \mathbf{H} \boldsymbol{\Phi}^{-1} \mathbf{H}^T. \quad (4.25)$$

Using the formula in Eq. (C.95) in the Appendix, the relationship

$$|\boldsymbol{\Phi}| |\beta^{-1} \mathbf{I} + \mathbf{H} \boldsymbol{\Phi}^{-1} \mathbf{H}^T| = |\beta^{-1} \mathbf{I}| |\boldsymbol{\Phi} + \beta \mathbf{H}^T \mathbf{H}| \quad (4.26)$$

holds. Thus, substituting Eqs. (4.25) and (4.13) into the above equation, we get

$$|\boldsymbol{\Phi}| |\boldsymbol{\Sigma}_y| = |\beta^{-1} \mathbf{I}| |\boldsymbol{\Gamma}|, \quad (4.27)$$

and

$$\log |\boldsymbol{\Sigma}_y| = \log |\boldsymbol{\Gamma}| - M \log \beta - \log |\boldsymbol{\Phi}|. \quad (4.28)$$

On the other hand, according to Sect. 4.10.2, Δ is expressed as

$$\Delta = \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k, \quad (4.29)$$

where $\boldsymbol{\Sigma}_y$ is given in Eq. (4.25). Thus, substituting Eqs. (4.28) and (4.29), into (4.24), we get

$$\log p(\mathbf{y}|\boldsymbol{\alpha}) = -\frac{1}{2} K \log |\boldsymbol{\Sigma}_y| - \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k. \quad (4.30)$$

The above equation indicates that $p(\mathbf{y}|\boldsymbol{\alpha})$ is Gaussian with the mean equal to zero and covariance matrix equal to $\boldsymbol{\Sigma}_y$. This $\boldsymbol{\Sigma}_y$ is called the model data covariance. Therefore, the estimate of $\boldsymbol{\alpha}$, $\hat{\boldsymbol{\alpha}}$, is obtained by maximizing $\log p(\mathbf{y}|\boldsymbol{\alpha})$ expressed above. Alternatively, defining the cost function such that

$$\mathcal{F}(\boldsymbol{\alpha}) = \log |\boldsymbol{\Sigma}_y| + \frac{1}{K} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k, \quad (4.31)$$

the estimate $\hat{\boldsymbol{\alpha}}$ is obtained by minimizing this cost function.

4.4 Update Equations for $\boldsymbol{\alpha}$

In this section, we derive the update equation for $\boldsymbol{\alpha}$. As will be shown, the updated equation contains the parameters of the posterior distribution. Since the value of $\boldsymbol{\alpha}$ is needed to compute the posterior distribution, the algorithm for computing $\boldsymbol{\alpha}$ is a recursive algorithm, as is the case of the EM algorithm presented in Sect. B.5 in the Appendix. That is, first setting an initial value for $\boldsymbol{\alpha}$, the posterior distribution is computed. Then, using the parameters of the posterior distribution, $\boldsymbol{\alpha}$ is updated. These procedures are repeated until a certain stopping condition is met.

Let us derive the update equation for $\boldsymbol{\alpha}$ by minimizing the cost function $\mathcal{F}(\boldsymbol{\alpha})$, i.e.,

$$\hat{\boldsymbol{\alpha}} = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \mathcal{F}(\boldsymbol{\alpha}).$$

The derivative of $\mathcal{F}(\boldsymbol{\alpha})$ with respect to $\boldsymbol{\alpha}$ is computed,

$$\frac{\partial \mathcal{F}(\boldsymbol{\alpha})}{\partial \alpha_j} = \frac{\partial}{\partial \alpha_j} \log |\boldsymbol{\Sigma}_y| + \frac{\partial}{\partial \alpha_j} \frac{1}{K} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k. \quad (4.32)$$

The first term in the right-hand side is expressed using Eq. (4.28) as

$$\begin{aligned} \frac{\partial}{\partial \alpha_j} \log |\boldsymbol{\Sigma}_y| &= \frac{\partial}{\partial \alpha_j} [-M \log \beta - \log |\boldsymbol{\Phi}| + \log |\boldsymbol{\Gamma}|] \\ &= -\frac{\partial}{\partial \alpha_j} \log |\boldsymbol{\Phi}| + \frac{\partial}{\partial \alpha_j} \log |\boldsymbol{\Gamma}|. \end{aligned} \quad (4.33)$$

First, we have

$$\frac{\partial}{\partial \alpha_j} \log |\boldsymbol{\Phi}| = \frac{\partial}{\partial \alpha_j} \sum_{j=1}^N \log \alpha_j = \alpha_j^{-1}. \quad (4.34)$$

Defining $\boldsymbol{\Pi}_{j,j}$ as an $(N \times N)$ matrix in which the (j, j) th component is equal to 1 and all other components are zero, we can compute the second term of Eq.(4.33), such that

$$\begin{aligned} \frac{\partial}{\partial \alpha_j} \log |\boldsymbol{\Gamma}| &= \text{tr} \left[\boldsymbol{\Gamma}^{-1} \frac{\partial}{\partial \alpha_j} \boldsymbol{\Gamma} \right] = \text{tr} \left[\boldsymbol{\Gamma}^{-1} \frac{\partial}{\partial \alpha_j} [\boldsymbol{\Phi} + \beta \mathbf{H}^T \mathbf{H}] \right] \\ &= \text{tr} \left[\boldsymbol{\Gamma}^{-1} \frac{\partial}{\partial \alpha_j} \boldsymbol{\Phi} \right] = \text{tr} \left[\boldsymbol{\Gamma}^{-1} \boldsymbol{\Pi}_{j,j} \right] = [\boldsymbol{\Gamma}^{-1}]_{j,j}, \end{aligned} \quad (4.35)$$

where $[\cdot]_{j,j}$ indicates the (j, j) th element of a matrix in the squared brackets. Note that $[\boldsymbol{\Gamma}^{-1}]_{j,j}$ above is the (j, j) th element of the posterior covariance matrix.

Next, we compute the second term of the right-hand side of Eq.(4.32). Using Eqs.(4.98) and (4.99), we can derive,

$$\begin{aligned} \frac{\partial}{\partial \alpha_j} \frac{1}{K} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k &= \frac{\partial}{\partial \alpha_j} \frac{1}{K} \sum_{k=1}^K \left[\beta \|\mathbf{y}_k - \mathbf{H} \bar{\mathbf{x}}_k\|^2 + \bar{\mathbf{x}}_k^T \boldsymbol{\Phi} \bar{\mathbf{x}}_k \right] \\ &= \frac{1}{K} \sum_{k=1}^K \bar{\mathbf{x}}_k^T \left[\frac{\partial}{\partial \alpha_j} \boldsymbol{\Phi} \right] \bar{\mathbf{x}}_k = \frac{1}{K} \sum_{k=1}^K \bar{\mathbf{x}}_k^T \boldsymbol{\Pi}_{j,j} \bar{\mathbf{x}}_k \\ &= \frac{1}{K} \sum_{k=1}^K \bar{x}_j^2(t_k). \end{aligned} \quad (4.36)$$

Therefore, denoting $[\boldsymbol{\Gamma}^{-1}]_{j,j}$ as $\Sigma_{j,j}$, the relationship

$$\frac{\partial \mathcal{F}(\boldsymbol{\alpha})}{\partial \alpha_j} = \Sigma_{j,j} - \alpha_j^{-1} + \frac{1}{K} \sum_{k=1}^K \bar{x}_j^2(t_k) = 0 \quad (4.37)$$

holds, and we get

$$\alpha_j^{-1} = \Sigma_{j,j} + \frac{1}{K} \sum_{k=1}^K \bar{x}_j^2(t_k). \quad (4.38)$$

This is equal to the update equation in the EM algorithm derived in Eq.(B.42).

On the other hand, we can derive a different update equation. To do so, we rewrite the expression for the posterior precision matrix in Eq.(4.13) to $\mathbf{I} - \boldsymbol{\Gamma}^{-1} \boldsymbol{\Phi} = \beta \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \mathbf{H}$, where the (j, j) th element of this equation is

$$1 - \alpha_j [\boldsymbol{\Gamma}^{-1}]_{j,j} = \left[\beta \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \mathbf{H} \right]_{j,j}. \quad (4.39)$$

We then rewrite Eq. (4.37) as

$$1 - \alpha_j [\boldsymbol{\Gamma}^{-1}]_{j,j} = \alpha_j \frac{1}{K} \sum_{k=1}^K \bar{x}_j^2(t_k). \quad (4.40)$$

Thus, we can derive the following equation:

$$\alpha_j = \frac{[\beta \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \mathbf{H}]_{j,j}}{\frac{1}{K} \sum_{k=1}^K \bar{x}_j^2(t_k)}. \quad (4.41)$$

Equation (4.41) is called the MacKay update equation [5]. The MacKay update is known to be faster than the EM update in Eq. (4.38) particularly when the estimation problem is highly ill-posed, i.e., $M < N$, although there is no theoretical proof that guarantees the convergence of the MacKay update.

Let us summarize the algorithm to estimate the source distribution \mathbf{x}_k . First, α is set to an appropriate initial value and the parameters of the posterior distribution, $\bar{\mathbf{x}}_k$ and $\boldsymbol{\Gamma}$, are computed using Eqs. (4.13) and (4.14). Then, the hyperparameter α is updated using Eq. (4.41) with the values of $\bar{\mathbf{x}}_k$ and $\boldsymbol{\Gamma}$ obtained in the preceding step. The algorithm is similar to the EM algorithm and the only difference is the update equation for α .

4.5 Modified Algorithm Integrating Interference Suppression

A simple modification of the algorithm described in the preceding sections leads to an algorithm robust to the interference overlapped onto the sensor data \mathbf{y} . The idea is similar to the one for the PFA algorithm described in Sect. 5.4, and the prerequisite is that a control measurement, which contains the interference but not the signal of interest, be available. Using the factor analysis model, the data is expressed such that

$$\mathbf{y}_k = \mathbf{B}\mathbf{u}_k + \varepsilon \quad \text{for control data}, \quad (4.42)$$

$$\mathbf{y}_k = \mathbf{H}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \varepsilon \quad \text{for target data}. \quad (4.43)$$

In the above equations, $L \times 1$ column vector \mathbf{u}_k is the factor activity and \mathbf{B} is an $M \times L$ mixing matrix. As in Sect. 5.4, $\mathbf{B}\mathbf{u}_k$ represents the interference.

Similar to the PFA algorithm, this modified Champagne algorithm has a two-step procedure. The first step applies the VBFA (or BFA) algorithm to the control data, and estimates the interference mixing matrix \mathbf{B} and the sensor-noise precision β . The target data is modeled as

$$\mathbf{y}_k = \mathbf{H}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \varepsilon = [\mathbf{H}, \mathbf{B}] \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix} + \varepsilon = \mathbf{H}_c \mathbf{z}_k + \varepsilon, \quad (4.44)$$

where

$$\mathbf{H}_c = [\mathbf{H}, \mathbf{B}] \quad \text{and} \quad \mathbf{z}_k = \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}.$$

Here, \mathbf{H}_c and \mathbf{z}_k can respectively be called the extended lead field matrix and the extended source vector. The second step applies the Champagne algorithm to the target data using the extended lead field matrix \mathbf{H}_c where \mathbf{B} is given from the first step.

The prior probability for the extended source vector \mathbf{z}_k is:

$$\begin{aligned} p(\mathbf{z}_k) &= p(\mathbf{x}_k)p(\mathbf{u}_k) = \mathcal{N}(\mathbf{x}_k|\mathbf{0}, \boldsymbol{\Phi}^{-1})\mathcal{N}(\mathbf{u}_k|\mathbf{0}, \mathbf{I}) \\ &= \left| \frac{\boldsymbol{\Phi}}{2\pi} \right|^{1/2} \exp \left[-\frac{1}{2} \mathbf{x}_k^T \boldsymbol{\Phi} \mathbf{x}_k \right] \left| \frac{\mathbf{I}}{2\pi} \right|^{1/2} \exp \left[-\frac{1}{2} \mathbf{u}_k^T \mathbf{u}_k \right] \\ &= \left| \frac{\tilde{\boldsymbol{\Phi}}}{2\pi} \right|^{1/2} \exp \left[-\frac{1}{2} \mathbf{z}_k^T \tilde{\boldsymbol{\Phi}} \mathbf{z}_k \right] = \mathcal{N}(\mathbf{z}_k|\mathbf{0}, \tilde{\boldsymbol{\Phi}}), \end{aligned} \quad (4.45)$$

where the $(N + L) \times (N + L)$ matrix $\tilde{\boldsymbol{\Phi}}$ is the prior precision matrix of the extended source vector, expressed as

$$\tilde{\boldsymbol{\Phi}} = \begin{bmatrix} \boldsymbol{\Phi} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \alpha_1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \alpha_N & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 1 & \cdots & 0 \\ \vdots & \cdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 1 \end{bmatrix}. \quad (4.46)$$

In this modified version of the Champagne algorithm, we use the same update equation for α with $\boldsymbol{\Phi}$ replaced with $\tilde{\boldsymbol{\Phi}}$. Thus, when updating the hyperparameter α , only components up to the N th diagonal element of $\tilde{\boldsymbol{\Phi}}$ are updated, and the rest of the elements are fixed to 1.

4.6 Convexity-Based Algorithm

In this section, we describe an alternative algorithm that minimizes the cost function in Eq. (4.31). The algorithm is called the convexity-based algorithm [1, 2]. It is faster than the EM algorithm. Unlike the MacKay update, this algorithm is guaranteed to converge. The algorithm also provides a theoretical basis for the sparsity analysis described in Sect. 4.7.

4.6.1 Deriving an Alternative Cost Function

The convexity-based algorithm makes use of the fact that $\log |\Sigma_y|$ is a concave function of $1/\alpha_1, \dots, 1/\alpha_N$, which is the voxel variance of the prior probability

distribution. Because of this, we use the voxel variance instead of the voxel precision in this section. We define prior voxel variance of the j th voxel as ν_j , which is equal to $1/\alpha_j$, and define the column vector $\boldsymbol{\nu}$ such that $\boldsymbol{\nu} = [\nu_1, \dots, \nu_N]^T$. We rewrite the cost function in Eq.(4.31) using $\boldsymbol{\nu}$,

$$\mathcal{F}(\boldsymbol{\nu}) = \log |\boldsymbol{\Sigma}_y| + \frac{1}{K} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k, \quad (4.47)$$

where the model data covariance $\boldsymbol{\Sigma}_y$ is expressed as

$$\boldsymbol{\Sigma}_y = \beta^{-1} \mathbf{I} + \mathbf{H} \boldsymbol{\Upsilon} \mathbf{H}^T, \quad (4.48)$$

and $\boldsymbol{\Upsilon}$ is the covariance matrix of the voxel prior distribution, defined as $\boldsymbol{\Upsilon} = \text{diag}([\nu_1, \dots, \nu_N])$.

The first term in Eq.(4.47), $\log |\boldsymbol{\Sigma}_y|$, is a concave function of ν_1, \dots, ν_N . Thus, for an arbitrary $\boldsymbol{\nu}$, we can find \mathbf{z} that satisfies the relationship

$$\mathbf{z}^T \boldsymbol{\nu} - z_o \geq \log |\boldsymbol{\Sigma}_y|, \quad (4.49)$$

where \mathbf{z} is the column vector $\mathbf{z} = [z_1, \dots, z_N]$, which contains auxiliary variables z_j ($j = 1, \dots, N$), and z_o is a scalar term that depends on \mathbf{z} .

The second term in Eq.(4.47) can be rewritten using

$$\mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k = \min_{\mathbf{x}_k} \left[\beta \|\mathbf{y}_k - \mathbf{H} \mathbf{x}_k\|^2 + \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k \right]. \quad (4.50)$$

The proof for the equation above is presented in Sect. 4.10.3. Therefore, we have the relationship

$$\frac{1}{K} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k = \min_{\mathbf{x}} \frac{1}{K} \sum_{k=1}^K \left[\beta \|\mathbf{y}_k - \mathbf{H} \mathbf{x}_k\|^2 + \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k \right]. \quad (4.51)$$

Using auxiliary variables \mathbf{z} and \mathbf{x} (where \mathbf{x} collectively expresses $\mathbf{x}_1, \dots, \mathbf{x}_K$), we define a new cost function $\tilde{\mathcal{F}}(\boldsymbol{\nu}, \mathbf{x}, \mathbf{z})$, such that

$$\tilde{\mathcal{F}}(\boldsymbol{\nu}, \mathbf{x}, \mathbf{z}) = \frac{1}{K} \sum_{k=1}^K \left[\beta \|\mathbf{y}_k - \mathbf{H} \mathbf{x}_k\|^2 + \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k \right] + \mathbf{z}^T \boldsymbol{\nu} - z_o. \quad (4.52)$$

Equations (4.49) and (4.51) guarantee that the relationship

$$\tilde{\mathcal{F}}(\boldsymbol{\nu}, \mathbf{x}, \mathbf{z}) \geq \mathcal{F}(\boldsymbol{\nu}),$$

always holds. That is, the alternative cost function $\tilde{\mathcal{F}}(\boldsymbol{\nu}, \mathbf{x}, \mathbf{z})$ forms an upper bound of the true cost function $\mathcal{F}(\boldsymbol{\nu})$. When we minimize $\tilde{\mathcal{F}}(\boldsymbol{\nu}, \mathbf{x}, \mathbf{z})$ with respect to $\boldsymbol{\nu}$, \mathbf{x} and \mathbf{z} , such $\boldsymbol{\nu}$ also minimizes the true cost function $\mathcal{F}(\boldsymbol{\nu})$.

4.6.2 Update Equation for z

Let us derive the update equation for the auxiliary variable z . The update value of z , \hat{z} , minimizes the cost function $\tilde{\mathcal{F}}(\boldsymbol{\nu}, \mathbf{x}, z)$ and at the same time satisfies the constraint in Eq. (4.49). Since $\mathbf{z}^T \boldsymbol{\nu} - z_o$ are the only terms that depend on z in $\tilde{\mathcal{F}}(\boldsymbol{\nu}, \mathbf{x}, z)$, this minimization problem is expressed as

$$\hat{z} = \operatorname{argmin}_z \left(\mathbf{z}^T \boldsymbol{\nu} - z_o \right) \quad \text{subject to } \mathbf{z}^T \boldsymbol{\nu} - z_o \geq \log |\boldsymbol{\Sigma}_y|. \quad (4.53)$$

That is, the minimization problem is equivalent to finding the hyperplane $\mathbf{z}^T \boldsymbol{\nu} - z_o$ that forms a closest upper bound of $\log |\boldsymbol{\Sigma}_y|$. Such a hyperplane is found as the plane that is tangential to $\log |\boldsymbol{\Sigma}_y|$ [6]. Therefore, the update value \hat{z} is given by

$$\hat{z} = \frac{\partial}{\partial \boldsymbol{\nu}} \log |\boldsymbol{\Sigma}_y|. \quad (4.54)$$

This \hat{z} forms a bound tightest to the concave function $\log |\boldsymbol{\Sigma}_y|$, and minimizes the cost function $\tilde{\mathcal{F}}(\boldsymbol{\nu}, \mathbf{x}, z)$ with respect to z . To compute \hat{z} , we use

$$\hat{z}_j = \frac{\partial}{\partial \nu_j} \log |\boldsymbol{\Sigma}_y| = \operatorname{tr} \left[\boldsymbol{\Sigma}_y^{-1} \frac{\partial}{\partial \nu_j} \boldsymbol{\Sigma}_y \right], \quad (4.55)$$

and

$$\begin{aligned} \frac{\partial}{\partial \nu_j} \boldsymbol{\Sigma}_y &= \frac{\partial}{\partial \nu_j} \left[\beta^{-1} \mathbf{I} + \mathbf{H} \boldsymbol{\Upsilon} \mathbf{H}^T \right] \\ &= \mathbf{H} \frac{\partial}{\partial \nu_j} \boldsymbol{\Upsilon} \mathbf{H}^T = \mathbf{H} \boldsymbol{\Pi}_{j,j} \mathbf{H}^T = \mathbf{l}_j \mathbf{l}_j^T, \end{aligned} \quad (4.56)$$

where \mathbf{l}_j is the lead field vector at the j th voxel, which is the j th column of \mathbf{H} . Substituting Eqs. (4.56) into (4.55), we get

$$\frac{\partial}{\partial \nu_j} \log |\boldsymbol{\Sigma}_y| = \operatorname{tr} \left[\boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j \mathbf{l}_j^T \right] = \operatorname{tr} \left[\mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j \right] = \mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j. \quad (4.57)$$

Using the equation above, we arrive at

$$\hat{z}_j = \operatorname{argmin}_{z_j} \tilde{\mathcal{F}}(\boldsymbol{\nu}, \mathbf{x}, z) = \mathbf{l}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_j \quad (4.58)$$

for the update equation of \hat{z}_j .

4.6.3 Update Equation for \mathbf{x}_k

Now, let us derive the update equation for the auxiliary variable \mathbf{x} . Since the relationship

$$\hat{\mathbf{x}}_k = \operatorname{argmin}_{\mathbf{x}_k} \left[\beta \|\mathbf{y}_k - \mathbf{H} \mathbf{x}_k\|^2 + \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k \right] \quad (4.59)$$

holds, according to Eq. (4.101) in Sect. 4.10.3, $\hat{\mathbf{x}}_k$ is obtained as

$$\hat{\mathbf{x}}_k = \beta \left(\boldsymbol{\Upsilon}^{-1} + \beta \mathbf{H}^T \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{y}_k. \quad (4.60)$$

This is the update equation for \mathbf{x}_k . Using the matrix inversion formula in Eq. (C.92), this equation can be rewritten as

$$\hat{\mathbf{x}}_k = \boldsymbol{\Upsilon} \mathbf{H}^T \left(\beta^{-1} \mathbf{I} + \mathbf{H} \boldsymbol{\Upsilon} \mathbf{H}^T \right)^{-1} \mathbf{y}_k = \boldsymbol{\Upsilon} \mathbf{H}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k. \quad (4.61)$$

Since this expression uses $\boldsymbol{\Upsilon}$ and $\boldsymbol{\Sigma}_y$, it is conveniently used in the convexity-based algorithm. Note that the auxiliary variable \mathbf{x}_k is equal to the posterior mean of the voxel source distribution because the update equation above is exactly equal to Eq. (B.26).

4.6.4 Update Equation for ν

Let us derive ν that minimizes the cost function $\tilde{\mathcal{F}}(\nu, z, \mathbf{x})$. Since only the second and the third terms in $\tilde{\mathcal{F}}(\nu, z, \mathbf{x})$ contains ν (as shown in Eq. (4.52)), we have:

$$\begin{aligned} \hat{\nu} &= \underset{\nu}{\operatorname{argmin}} \tilde{\mathcal{F}}(\nu, \mathbf{x}, z) = \underset{\nu}{\operatorname{argmin}} \left[\frac{1}{K} \sum_{k=1}^K \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k + z^T \nu \right] \\ &= \underset{\nu}{\operatorname{argmin}} \sum_{j=1}^N \left[z_j \nu_j + \frac{\frac{1}{K} \sum_{k=1}^K x_j^2(t_k)}{\nu_j} \right]. \end{aligned} \quad (4.62)$$

Thus, considering the relationship

$$\frac{\partial}{\partial \nu_j} \sum_{k=1}^N \left[z_j \nu_j + \frac{\frac{1}{K} \sum_{k=1}^K x_j^2(t_k)}{\nu_j} \right] = z_j - \frac{\frac{1}{K} \sum_{k=1}^K x_j^2(t_k)}{\nu_j^2} = 0,$$

we can derive

$$\hat{\nu}_j = \underset{\nu_j}{\operatorname{argmin}} \tilde{\mathcal{F}}(\nu, z, \mathbf{x}) = \sqrt{\frac{\frac{1}{K} \sum_{k=1}^K x_j^2(t_k)}{z_j}}. \quad (4.63)$$

4.6.5 Summary of the Convexity-Based Algorithm

The update for z in Eq. (4.58) and the update for \mathbf{x} in Eq. (4.61) require ν to be known. The update equation for ν in Eq. (4.63) requires \mathbf{x} and z to be known. Therefore, the convexity-based algorithm updates z , \mathbf{x} and ν using Eqs. (4.58), (4.61) and (4.63), in

a recursive manner. Since the auxiliary variable \mathbf{x}_k is the posterior mean of the voxel source distribution, the value of the auxiliary variable $\widehat{\mathbf{x}}_k$ is equal to the Bayesian estimate of the source distribution, $\bar{\mathbf{x}}_k$, after the iterative procedure is terminated.

4.7 The Origin of the Sparsity

Why does the Champagne algorithm produce sparse solutions? This section tries to answer this question, and consider the origin of the sparsity by analyzing the cost function of the Champagne algorithm [7]. For simplicity, we set $K = 1$, and by omitting the time index, \mathbf{x}_1 and \mathbf{y}_1 are denoted \mathbf{x} and \mathbf{y} . Using Eq.(4.47), we have

$$\widehat{\boldsymbol{\nu}} = \operatorname{argmin}_{\boldsymbol{\nu}} \left(\log |\boldsymbol{\Sigma}_y| + \mathbf{y}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y} \right).$$

As shown in Eq.(4.50), we have

$$\mathbf{y}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y} = \min_{\mathbf{x}} \left[\beta \|\mathbf{y} - \mathbf{Hx}\|^2 + \mathbf{x}^T \boldsymbol{\gamma}^{-1} \mathbf{x} \right].$$

Combining the two equations above, the cost function for estimating \mathbf{x} is given by

$$\widehat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \mathcal{F}: \quad \mathcal{F} = \beta \|\mathbf{y} - \mathbf{Hx}\|^2 + \phi(\mathbf{x}), \quad (4.64)$$

where the constraint $\phi(\mathbf{x})$ is expressed as

$$\phi(\mathbf{x}) = \min_{\boldsymbol{\nu}} \left(\mathbf{x}^T \boldsymbol{\gamma}^{-1} \mathbf{x} + \log |\boldsymbol{\Sigma}_y| \right) = \min_{\boldsymbol{\nu}} \left(\sum_{j=1}^N \frac{x_j^2}{\nu_j} + \log |\boldsymbol{\Sigma}_y| \right). \quad (4.65)$$

However, computing $\phi(\mathbf{x})$ in Eq.(4.65) is not so easy because $\log |\boldsymbol{\Sigma}_y|$ contains $\boldsymbol{\nu}$ as in Eq.(4.48). Therefore, we further introduce a simplification by assuming that the columns in the matrix \mathbf{H} are orthogonal, i.e., the relationship $\mathbf{l}_i^T \mathbf{l}_j = I_{i,j}$ is assumed to hold. In this case, using Eq.(4.48), we get

$$\log |\boldsymbol{\Sigma}_y| = \sum_{j=1}^N \log(\beta^{-1} + \nu_j),$$

and, accordingly,

$$\phi(\mathbf{x}) = \min_{\boldsymbol{\nu}} \sum_{j=1}^N \left(\frac{x_j^2}{\nu_j} + \log(\beta^{-1} + \nu_j) \right). \quad (4.66)$$

By computing the minimum on the right-hand side, we finally have

$$\phi(\mathbf{x}) = \sum_{j=1}^N \varphi(x_j), \quad (4.67)$$

where

$$\varphi(x) = \frac{2|x|}{\sqrt{x^2 + 4\beta^{-1}} + |x|} + \log \left[\beta^{-1} + \frac{x^2}{2} + \frac{1}{2}|x|\sqrt{x^2 + 4\beta^{-1}} \right]. \quad (4.68)$$

The constraint $\varphi(x)$ in Eq. (4.68) is plotted in Fig. 4.1. In Fig. 4.1a, the plot of $\varphi(x)$ is shown by the solid line. For comparison, the constraint for the L_1 -norm solution, $|x|$, is also shown by the broken line. The plots in Fig. 4.1a show that the constraint of the Champagne algorithm $\varphi(x)$ is very similar to (but sharper than) the L_1 -norm constraint $|x|$, suggesting that the Champagne algorithm produces sparse solutions.

In Fig. 4.1b, the plots of $\varphi(x)$ when $\beta^{-1} = 0.1$, $\beta^{-1} = 1$ and $\beta^{-1} = 10$ are shown by the dot-and-dash, broken, and solid lines, respectively. The vertical broken line at $x = 0$ shows the L_0 -norm constraint, $\|x\|_0$, for comparison. It is shown here that the shape of the constraint $\varphi(x)$ depends on β^{-1} —namely the noise variance. When the noise variance is small (i.e., a high SNR), $\varphi(x)$ is much sharper than $|x|$, and becomes closer to the L_0 -norm constraint. That is, the Champagne algorithm uses an adaptive constraint. When the SNR of the sensor data is high, it gives solutions with enhanced sparsity. In contrast, when the sensor data is noisy, the shape of $\varphi(x)$ becomes similar to the shape of the L_1 -norm constraint and the algorithm gives solutions with mild sparsity.

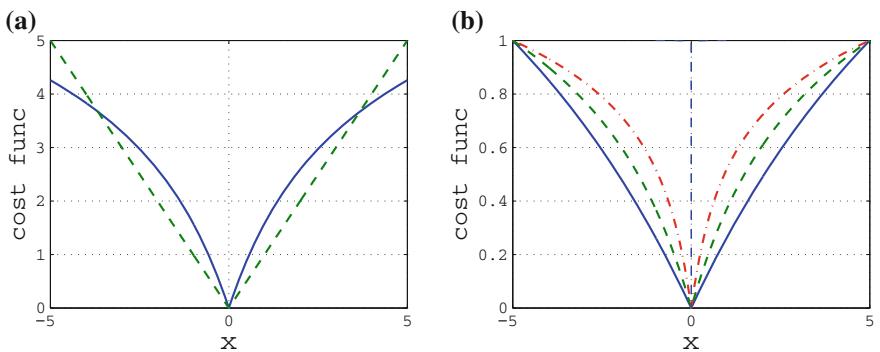


Fig. 4.1 Plots of the cost function $\varphi(x)$ shown in Eq.(4.68). **a** The *solid line* shows the plot of $\varphi(x)$ with $\beta^{-1} = 1$, and the *broken line* the plot of $|x|$, which is the constraint for the L_1 -norm minimum-norm solution. **b** Plots of $\varphi(x)$ when $\beta^{-1} = 0.1$, $\beta^{-1} = 1$ and $\beta^{-1} = 10$ shown by the *dot-and-dash*, *broken*, and *solid lines*, respectively. In this figure, the *vertical broken line* at $x = 0$ shows the constraint for the L_0 -norm solution

4.8 Extension to Include Source Vector Estimation

In Sect. 4.6, we derive the convexity-based algorithm when the source orientation at each voxel is known, i.e., when the relationship between the sensor data and the voxel source distribution is expressed as in Eq. (4.3).

In this section, we derive the convexity-based algorithm when the source orientation at each voxel is unknown and the relationship,

$$\mathbf{y}_k = \mathbf{F}\mathbf{x}_k + \boldsymbol{\varepsilon}, \quad (4.69)$$

holds, where \mathbf{x}_k is given by

$$\mathbf{x}_k = \left[\mathbf{s}_1^T(t_k), \dots, \mathbf{s}_N^T(t_k) \right]^T, \quad (4.70)$$

and $\mathbf{s}_j(t_k)$ is the 3×1 source vector at the j th voxel. In other words, we describe an extension that enables estimating the source vector at each voxel. The algorithms in the preceding sections use a diagonal prior covariance (or precision) matrix, and the use of the diagonal matrix is possible because the source orientation is known. Therefore, a naïve application of those algorithms to cases where the source orientation is unknown generally results in the shrinkage of source vector components. Namely, the naïve application possibly leads to a solution where only a single component of a source vector has nonzero value and other components have values close to zero, leading to incorrect estimation of source orientations.

Since the unknown parameter at each voxel is not a scalar quantity but the 3×1 vector quantity in this section, the algorithm being developed here uses the nondiagonal covariance matrix of the prior distribution. That is, the prior distribution is assumed to be [1]:

$$p(\mathbf{x}_k | \boldsymbol{\Upsilon}) = \prod_{j=1}^N \mathcal{N}(\mathbf{s}_j(t_k) | \mathbf{0}, \boldsymbol{\Upsilon}_j), \quad (4.71)$$

where $\boldsymbol{\Upsilon}_j$ is a 3×3 covariance matrix of the source vector $\mathbf{s}_j(t_k)$. Thus, the covariance matrix of the voxel source distribution, $\boldsymbol{\Upsilon}$, is a $3N \times 3N$ block diagonal matrix expressed as

$$\boldsymbol{\Upsilon} = \begin{bmatrix} \boldsymbol{\Upsilon}_1 & 0 & \cdots & 0 \\ 0 & \boldsymbol{\Upsilon}_2 & \cdots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & 0 & \cdots & \boldsymbol{\Upsilon}_N \end{bmatrix}. \quad (4.72)$$

The cost function in this case has the same form:

$$\mathcal{F}(\boldsymbol{\Upsilon}) = \log |\boldsymbol{\Sigma}_y| + \frac{1}{K} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k. \quad (4.73)$$

However, Σ_y is given by

$$\Sigma_y = \beta^{-1} I + \sum_{j=1}^N \mathbf{L}_j \boldsymbol{\Upsilon}_j \mathbf{L}_j^T, \quad (4.74)$$

where the $M \times 3$ lead field matrix at the j th voxel, $\mathbf{L}(\mathbf{r}_j)$, is denoted \mathbf{L}_j for simplicity.

In agreement with Eq. (4.49), since $\log |\Sigma_y|$ is a concave function, we can find 3×3 auxiliary parameter matrices, \mathbf{Z}_j , ($j = 1, \dots, N$) that satisfy

$$\sum_{j=1}^N \text{tr} (\mathbf{Z}_j^T \boldsymbol{\Upsilon}_j) - Z_o \geq \log |\Sigma_y|, \quad (4.75)$$

where Z_o is a scalar term that depends on \mathbf{Z}_j . Regarding the second term in the right-hand side of Eq. (4.73), we have the relationship

$$\begin{aligned} \mathbf{y}_k^T \Sigma_y^{-1} \mathbf{y}_k &= \min_{\mathbf{x}_k} \left[\beta \|\mathbf{y}_k - \mathbf{F}\mathbf{x}_k\|^2 + \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k \right] \\ &= \min_{\mathbf{x}_k} \left[\beta \|\mathbf{y}_k - \mathbf{F}\mathbf{x}_k\|^2 + \sum_{j=1}^N \mathbf{s}_j^T(t_k) \boldsymbol{\Upsilon}_j^{-1} \mathbf{s}_j(t_k) \right]. \end{aligned} \quad (4.76)$$

Let us use \mathbf{Z} to collectively express $\mathbf{Z}_1, \dots, \mathbf{Z}_N$, and use \mathbf{s} to collectively express $\mathbf{s}_j(t_k)$, where $j = 1, \dots, N$ and $k = 1, \dots, K$. The alternative cost function, $\tilde{\mathcal{F}}(\boldsymbol{\Upsilon}, \mathbf{s}, \mathbf{Z})$, is obtained as

$$\begin{aligned} \tilde{\mathcal{F}}(\boldsymbol{\Upsilon}, \mathbf{s}, \mathbf{Z}) &= \frac{1}{K} \sum_{k=1}^K \left[\beta \|\mathbf{y}_k - \mathbf{F}\mathbf{x}_k\|^2 + \sum_{j=1}^N \mathbf{s}_j(t_k) \boldsymbol{\Upsilon}_j^{-1} \mathbf{s}_j^T(t_k) \right] \\ &\quad + \sum_{j=1}^N \text{tr} (\mathbf{Z}_j^T \boldsymbol{\Upsilon}_j) - Z_o. \end{aligned} \quad (4.77)$$

This $\tilde{\mathcal{F}}(\boldsymbol{\Upsilon}, \mathbf{s}, \mathbf{Z})$ forms an upper bound of the true cost function in Eq. (4.73). Accordingly, when we minimize $\tilde{\mathcal{F}}(\boldsymbol{\Upsilon}, \mathbf{s}, \mathbf{Z})$ with respect to $\boldsymbol{\Upsilon}$, \mathbf{s} , and \mathbf{Z} , we can minimize the true cost function $\mathcal{F}(\boldsymbol{\Upsilon})$.

4.8.1 Update Equation for \mathbf{Z}_j

Using the same arguments in Sect. 4.6.2, the update equation for \mathbf{Z}_j is derived as

$$\widehat{\mathbf{Z}}_j = \frac{\partial}{\partial \boldsymbol{\Upsilon}_j} \log |\Sigma_y|. \quad (4.78)$$

Here, the hyperplane $\sum_{j=1}^N \text{tr}(\widehat{\mathbf{Z}}_j^T \boldsymbol{\Upsilon}_j) - Z_o$ forms a tightest upper bound of the concave function $\log |\boldsymbol{\Sigma}_y|$. Let us express the lead field matrix at the j th voxel using its column vectors,

$$\mathbf{L}_j = [\mathbf{l}_x(\mathbf{r}_j), \mathbf{l}_y(\mathbf{r}_j), \mathbf{l}_z(\mathbf{r}_j)] = [\mathbf{l}_x, \mathbf{l}_y, \mathbf{l}_z],$$

where explicit notation of the voxel location (\mathbf{r}_j) is omitted. Using exactly the same derivation from Eqs. (4.55) to (4.58), we can derive the update equation

$$\begin{aligned} \widehat{\mathbf{Z}}_j &= \frac{\partial}{\partial \boldsymbol{\Upsilon}_j} \log |\boldsymbol{\Sigma}_y| \\ &= \begin{bmatrix} \partial \log |\boldsymbol{\Sigma}_y| / \partial [\boldsymbol{\Upsilon}_j]_{1,1} & \partial \log |\boldsymbol{\Sigma}_y| / \partial [\boldsymbol{\Upsilon}_j]_{1,2} & \partial \log |\boldsymbol{\Sigma}_y| / \partial [\boldsymbol{\Upsilon}_j]_{1,3} \\ \partial \log |\boldsymbol{\Sigma}_y| / \partial [\boldsymbol{\Upsilon}_j]_{2,1} & \partial \log |\boldsymbol{\Sigma}_y| / \partial [\boldsymbol{\Upsilon}_j]_{2,2} & \partial \log |\boldsymbol{\Sigma}_y| / \partial [\boldsymbol{\Upsilon}_j]_{2,3} \\ \partial \log |\boldsymbol{\Sigma}_y| / \partial [\boldsymbol{\Upsilon}_j]_{3,1} & \partial \log |\boldsymbol{\Sigma}_y| / \partial [\boldsymbol{\Upsilon}_j]_{3,2} & \partial \log |\boldsymbol{\Sigma}_y| / \partial [\boldsymbol{\Upsilon}_j]_{3,3} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{l}_x^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_x & \mathbf{l}_x^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_y & \mathbf{l}_x^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_z \\ \mathbf{l}_y^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_x & \mathbf{l}_y^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_y & \mathbf{l}_y^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_z \\ \mathbf{l}_z^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_x & \mathbf{l}_z^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_y & \mathbf{l}_z^T \boldsymbol{\Sigma}_y^{-1} \mathbf{l}_z \end{bmatrix} = \mathbf{L}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{L}_j, \end{aligned} \quad (4.79)$$

where $[\boldsymbol{\Upsilon}_j]_{\ell,m}$ indicates the (ℓ, m) th element of the matrix $\boldsymbol{\Upsilon}_j$.

4.8.2 Update Equation for $s_j(t_k)$

The update equation for \mathbf{x}_k can be obtained using

$$\widehat{\mathbf{x}}_k = \underset{\mathbf{x}_k}{\operatorname{argmin}} \left[\beta \|\mathbf{y}_k - \mathbf{F}\mathbf{x}_k\|^2 + \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k \right]. \quad (4.80)$$

The solution is given by Eq. (4.61) with replacing \mathbf{H} with \mathbf{F} , which is rewritten as

$$\begin{bmatrix} \widehat{s}_1(t_k) \\ \widehat{s}_2(t_k) \\ \vdots \\ \widehat{s}_N(t_k) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Upsilon}_1 & 0 & \cdots & 0 \\ 0 & \boldsymbol{\Upsilon}_2 & \cdots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & 0 & \cdots & \boldsymbol{\Upsilon}_N \end{bmatrix} \begin{bmatrix} \mathbf{L}_1^T \\ \mathbf{L}_2^T \\ \vdots \\ \mathbf{L}_N^T \end{bmatrix} \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k. \quad (4.81)$$

Therefore, the source vector at the j th voxel $\widehat{s}_j(t_k)$ is given by

$$\widehat{s}_j(t_k) = \boldsymbol{\Upsilon}_j \mathbf{L}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k, \quad (4.82)$$

and Eq. (4.82) is the update equation for $s_j(t_k)$.

4.8.3 Update Equation for $\boldsymbol{\Upsilon}_j$

The update equation for $\boldsymbol{\Upsilon}_j$ is obtained using

$$\widehat{\boldsymbol{\Upsilon}}_j = \underset{\boldsymbol{\Upsilon}_j}{\operatorname{argmin}} \left[\frac{1}{K} \sum_{k=1}^K \mathbf{s}_j^T(t_k) \boldsymbol{\Upsilon}_j^{-1} \mathbf{s}_j(t_k) + \operatorname{tr}(\mathbf{Z}_j^T \boldsymbol{\Upsilon}_j) \right]. \quad (4.83)$$

Taking

$$\frac{\partial}{\partial \boldsymbol{\Upsilon}_j} \mathbf{s}_j^T(t_k) \boldsymbol{\Upsilon}_j^{-1} \mathbf{s}_j(t_k) = -\boldsymbol{\Upsilon}_j^{-1} \mathbf{s}_j(t_k) \mathbf{s}_j^T(t_k) \boldsymbol{\Upsilon}_j^{-1} \quad (4.84)$$

into consideration, we have

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\Upsilon}_j} & \left[\frac{1}{K} \sum_{k=1}^K \mathbf{s}_j^T(t_k) \boldsymbol{\Upsilon}_j^{-1} \mathbf{s}_j(t_k) + \operatorname{tr}(\mathbf{Z}_j^T \boldsymbol{\Upsilon}_j) \right] \\ & = -\boldsymbol{\Upsilon}_j^{-1} \left[\frac{1}{K} \sum_{k=1}^K \mathbf{s}_j(t_k) \mathbf{s}_j^T(t_k) \right] \boldsymbol{\Upsilon}_j^{-1} + \mathbf{Z}_j. \end{aligned} \quad (4.85)$$

Setting the right-hand side to zero, we get the equation,

$$\boldsymbol{\Upsilon}_j \mathbf{Z}_j \boldsymbol{\Upsilon}_j = \left[\frac{1}{K} \sum_{k=1}^K \mathbf{s}_j(t_k) \mathbf{s}_j^T(t_k) \right]. \quad (4.86)$$

However, there are multiple solutions for $\boldsymbol{\Upsilon}_j$ that satisfy Eq. (4.86). We should find a positive semidefinite matrix that satisfies Eq. (4.86). Defining $\mathcal{E} = (1/K) \sum_{k=1}^K \mathbf{s}_j(t_k) \mathbf{s}_j^T(t_k)$, and using

$$\begin{aligned} \mathcal{E} & = \mathbf{Z}_j^{-1/2} (\mathbf{Z}_j^{1/2} \mathcal{E} \mathbf{Z}_j^{1/2}) \mathbf{Z}_j^{-1/2} \\ & = \mathbf{Z}_j^{-1/2} (\mathbf{Z}_j^{1/2} \mathcal{E} \mathbf{Z}_j^{1/2})^{1/2} (\mathbf{Z}_j^{1/2} \mathcal{E} \mathbf{Z}_j^{1/2})^{1/2} \mathbf{Z}_j^{-1/2} \\ & = \mathbf{Z}_j^{-1/2} (\mathbf{Z}_j^{1/2} \mathcal{E} \mathbf{Z}_j^{1/2})^{1/2} \mathbf{Z}_j^{-1/2} \mathbf{Z}_j \mathbf{Z}_j^{-1/2} (\mathbf{Z}_j^{1/2} \mathcal{E} \mathbf{Z}_j^{1/2})^{1/2} \mathbf{Z}_j^{-1/2}, \end{aligned} \quad (4.87)$$

the solution for $\boldsymbol{\Upsilon}_j$ that is a positive semidefinite matrix is derived such that

$$\widehat{\boldsymbol{\Upsilon}}_j = \mathbf{Z}_j^{-1/2} \left[\mathbf{Z}_j^{1/2} \left[\frac{1}{K} \sum_{k=1}^K \mathbf{s}_j(t_k) \mathbf{s}_j^T(t_k) \right] \mathbf{Z}_j^{1/2} \right]^{1/2} \mathbf{Z}_j^{-1/2}. \quad (4.88)$$

Equation (4.88) is the update equation for $\boldsymbol{\Upsilon}_j$.

In summary, \mathbf{Z}_j , $\mathbf{s}_j(t_k)$, and $\boldsymbol{\Upsilon}_j$ are updated using Eqs. (4.79), (4.82) and (4.88), respectively. Since the auxiliary variable \mathbf{x}_k is equal to the posterior mean of the

voxel source distribution, the Bayes estimate of the voxel source distribution, $\bar{\mathbf{x}}_k$, can be obtained as the final updated results of $\hat{\mathbf{s}}_j(t_k)$.

4.9 Source Vector Estimation Using Hyperparameter Tying

In Sect. 4.8, we describe an algorithm that reconstructs the source distribution when the source orientation is unknown. The algorithm uses the block diagonal covariance matrix $\boldsymbol{\Upsilon}$ expressed in Eq. (4.72). However, because the matrix $\boldsymbol{\Upsilon}$ is a nondiagonal matrix, the resultant algorithm is rather complex and computationally expensive, compared to the algorithm that uses a diagonal covariance matrix. In this section, we describe an extension of the Champagne algorithm that still uses a diagonal covariance matrix for the prior distribution but enables us to estimate the voxel source vectors.

The method uses a technique called hyperparameter tying, and simple modification of the algorithm described in Sects. 4.2–4.6 leads to a successful estimation of voxel source vectors. We use the forward relationship given in Eq. (4.69), and a $3N \times 3N$ diagonal covariance matrix for the Gaussian prior:

$$\boldsymbol{\Upsilon} = \begin{bmatrix} \ddots & & & & \vdots \\ & \ddots & & & \vdots \\ \cdots & \nu_{3n+1} & 0 & 0 & \cdots \\ \cdots & 0 & \nu_{3n+2} & 0 & \cdots \\ \cdots & 0 & 0 & \nu_{3n+3} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (4.89)$$

where we show a diagonal block matrix corresponding to the n th voxel. The equation in Eq. (4.63) is applied to update ν_j . In the method proposed here, the three hyperparameters ν_{3n+1} , ν_{3n+2} , and ν_{3n+3} from the same voxel are tied together to have the same new update value. That is, for $\hat{\nu}_{3n+1}$, $\hat{\nu}_{3n+2}$, and $\hat{\nu}_{3n+3}$, we compute a new single update value by using

$$\hat{\nu}_{3n+m}^{\text{new}} = \frac{1}{3} \sum_{i=1}^3 \hat{\nu}_{3n+i}, \quad \text{for } m = 1, 2, 3 \quad (4.90)$$

where $\hat{\nu}_{3n+i}$ (for $i = 1, 2, 3$) are the update values obtained from Eq. (4.63), and $\hat{\nu}_{3n+m}^{\text{new}}$ ($m = 1, 2, 3$) is a new update value for these hyperparameters.

The hyperparameter tying method can be implemented in exactly the same manner when using the MacKay update in Eq. (4.41). A new single update value for the voxel precision corresponding to the n th voxel is computed by using

$$\hat{\alpha}_{3n+m}^{\text{new}} = \frac{1}{3} \sum_{i=1}^3 \hat{\alpha}_{3n+i}, \quad \text{for } m = 1, 2, 3 \quad (4.91)$$

where $\widehat{\alpha}_{3n+i}$ (for $i = 1, 2, 3$) is the update value from Eq.(4.41) and $\widehat{\alpha}_{3n+m}^{\text{new}}$ (for $m = 1, 2, 3$) is a new update value for these hyperparameters.

The rationale for this hyperparameter tying can be explained using the cost function analysis described in Sect. 4.7. Let us compute the constraint function in Eq.(4.66) for a two-dimensional case in which the unknown parameters are denoted x_1 and x_2 . The constraint is rewritten in this case as,

$$\phi(x_1, x_2) = \min_{\nu_1, \nu_2} \sum_{j=1}^2 \left(\frac{x_j^2}{\nu_j} + \log(\beta^{-1} + \nu_j) \right). \quad (4.92)$$

When the hyperparameters ν_1 and ν_2 are tied together, i.e., when we set these hyperparameters at the same value ν , the constraint function is changed to

$$\phi(x_1, x_2) = \min_{\nu} \left(\frac{x_1^2 + x_2^2}{\nu} + 2 \log(\beta^{-1} + \nu) \right). \quad (4.93)$$

By implementing this minimization, the value of ν that minimizes the right-hand side of the above equation, $\widehat{\nu}$, is derived as

$$\widehat{\nu} = \frac{a + \sqrt{a^2 + 8\beta^{-1}a}}{4},$$

where $a = x_1^2 + x_2^2$. Substituting this $\widehat{\nu}$ into Eq.(4.93), we derive the constraint function,

$$\phi(x_1, x_2) = \frac{4a}{a + \sqrt{a^2 + 8\beta^{-1}a}} + 2 \log \left[\beta^{-1} + \frac{a + \sqrt{a^2 + 8\beta^{-1}a}}{4} \right]. \quad (4.94)$$

The plot of the constraint function in Eq.(4.94) is shown in Fig. 4.2a. For comparison, the Champagne constraint function when untying ν_1 and ν_2 (Eq.4.92) is shown in Fig. 4.2b. The constraint functions for the L_2 and L_1 -norm regularizations are also shown in Fig. 4.2c, d for comparison. These plots show that the Champagne constraint function when tying ν_1 and ν_2 has a shape very similar to the constraint function for the L_2 regularization. According to the arguments in Sect. 2.9.2, this type of constraint does not generate a sparse solution. Thus, when tying the hyper-parameter update values, the sparsity is lost among the solutions of x_{2n+1} , x_{2n+2} , and x_{2n+3} and there is no shrinkage over the source vector components. However, since the sparsity is maintained across voxels, a sparse source distribution can still be reconstructed.

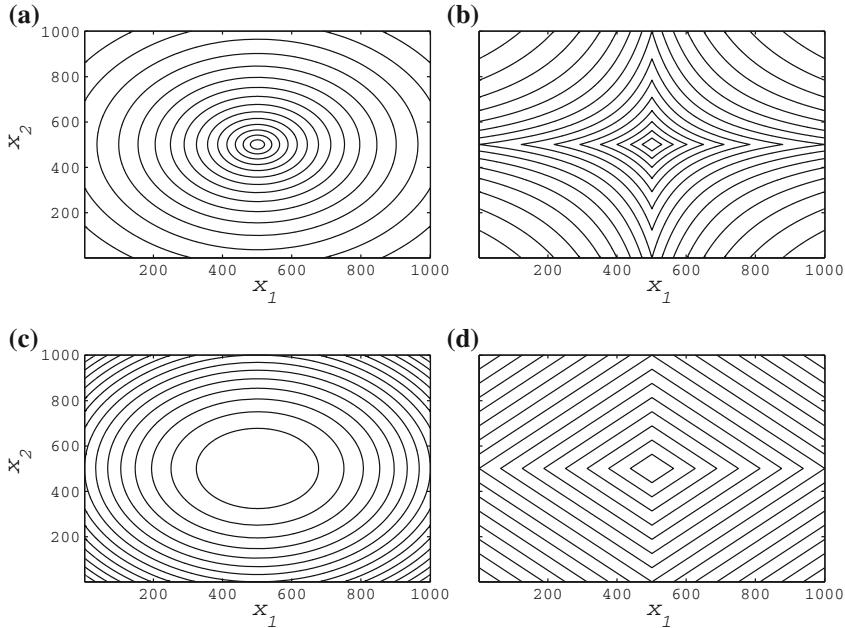


Fig. 4.2 Plots of the constraint functions for the two-dimensional case. **a** The plot of $\phi(x_1, x_2)$ in Eq. (4.94). **b** The plot of $\phi(x_1, x_2)$ in Eq. (4.92). **c** The plot of the constraint function for the L_2 -norm constraint: $\phi(x_1, x_2) = x_1^2 + x_2^2$. **d** The plot of the constraint function for the L_1 -norm constraint: $\phi(x_1, x_2) = |x_1| + |x_2|$. The parameter β was set to 1 when computing the plots in **a** and **b**

4.10 Appendix to This Chapter

4.10.1 Derivation of Eq. (4.21)

D in Eq. (4.20) is rewritten as

$$D = \sum_{k=1}^K \left[\frac{1}{2} \mathbf{x}_k^T [\boldsymbol{\Phi} + \beta \mathbf{H}^T \mathbf{H}] \mathbf{x}_k - \beta \mathbf{x}_k^T \mathbf{H}^T \mathbf{y}_k \right] + \mathcal{C}. \quad (4.95)$$

Here, terms not containing \mathbf{x}_k are expressed as \mathcal{C} . We apply the completion of the square with respect to \mathbf{x}_k , i.e., we change Eq. (4.95) to have a form,

$$D = \frac{1}{2} \sum_{k=1}^K (\mathbf{x}_k - \mathbf{m}_k)^T \mathbf{A} (\mathbf{x}_k - \mathbf{m}_k) + \Delta,$$

where \mathbf{A} is a real symmetric matrix and \mathbf{m}_k is a column vector, and Δ represents remaining terms. The first term on the right-hand side is rewritten as

$$\frac{1}{2} \sum_{k=1}^K (\mathbf{x}_k - \mathbf{m}_k)^T \mathbf{A} (\mathbf{x}_k - \mathbf{m}_k) = \sum_{k=1}^K \left[\frac{1}{2} \mathbf{x}_k^T \mathbf{A} \mathbf{x}_k - \mathbf{x}_k^T \mathbf{A} \mathbf{m}_k + \cdots \right].$$

Comparing the right-hand side of the equation above with that of Eq. (4.95), we get

$$\begin{aligned} \mathbf{A} &= \boldsymbol{\Phi} + \beta \mathbf{H}^T \mathbf{H}, \\ \mathbf{A} \mathbf{m}_k &= \beta \mathbf{H}^T \mathbf{y}_k \quad \text{namely} \quad \mathbf{m}_k = \beta [\boldsymbol{\Phi} + \beta \mathbf{H}^T \mathbf{H}]^{-1} \mathbf{H}^T \mathbf{y}_k. \end{aligned}$$

Comparing the equations above with Eqs. (4.13) and (4.14), we have the relationships $\mathbf{A} = \boldsymbol{\Gamma}$ and $\mathbf{m}_k = \bar{\mathbf{x}}_k$, giving

$$D = \frac{1}{2} \sum_{k=1}^K (\mathbf{x}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Gamma} (\mathbf{x}_k - \bar{\mathbf{x}}_k) + \Delta. \quad (4.96)$$

This equation shows that D reaches the minimum at $\mathbf{x}_k = \bar{\mathbf{x}}_k$, and the minimum value is equal to Δ . The value of Δ is obtained by substituting $\mathbf{x}_k = \bar{\mathbf{x}}_k$ into Eq. (4.20), such that

$$\Delta = \frac{\beta}{2} \sum_{k=1}^K \|\mathbf{y}_k - \mathbf{H} \bar{\mathbf{x}}_k\|^2 + \frac{1}{2} \sum_{k=1}^K \bar{\mathbf{x}}_k^T \boldsymbol{\Phi} \bar{\mathbf{x}}_k. \quad (4.97)$$

4.10.2 Derivation of Eq. (4.29)

The derivation starts from Δ in Eq. (4.97). On the right-hand side of this equation, the k th terms containing \mathbf{x}_k are

$$\begin{aligned} \frac{\beta}{2} \|\mathbf{y}_k - \mathbf{H} \bar{\mathbf{x}}_k\|^2 + \frac{1}{2} \bar{\mathbf{x}}_k^T \boldsymbol{\Phi} \bar{\mathbf{x}}_k &= \frac{\beta}{2} \left(\mathbf{y}_k^T \mathbf{y}_k - 2 \bar{\mathbf{x}}_k^T \mathbf{H}^T \mathbf{y}_k + \bar{\mathbf{x}}_k^T \mathbf{H}^T \mathbf{H} \bar{\mathbf{x}}_k \right) + \frac{1}{2} \bar{\mathbf{x}}_k^T \boldsymbol{\Phi} \bar{\mathbf{x}}_k \\ &= \frac{1}{2} \left[\beta \mathbf{y}_k^T \mathbf{y}_k - 2 \bar{\mathbf{x}}_k^T \beta \mathbf{H}^T \mathbf{y}_k + \bar{\mathbf{x}}_k^T (\boldsymbol{\Phi} + \beta \mathbf{H}^T \mathbf{H}) \bar{\mathbf{x}}_k \right] \\ &= \frac{1}{2} \left[\beta \mathbf{y}_k^T \mathbf{y}_k - 2 \bar{\mathbf{x}}_k^T \beta \mathbf{H}^T \mathbf{y}_k + \bar{\mathbf{x}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{x}}_k \right]. \end{aligned} \quad (4.98)$$

Using the relationship $\beta \mathbf{H}^T \mathbf{y}_k = \boldsymbol{\Gamma} \bar{\mathbf{x}}_k$, we have

$$\begin{aligned} \frac{1}{2} \left[\beta \mathbf{y}_k^T \mathbf{y}_k - 2 \bar{\mathbf{x}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{x}}_k + \bar{\mathbf{x}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{x}}_k \right] &= \frac{1}{2} \left[\beta \mathbf{y}_k^T \mathbf{y}_k - \bar{\mathbf{x}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{x}}_k \right] \\ &= \frac{1}{2} \left[\beta \mathbf{y}_k^T \mathbf{y}_k - (\beta \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \mathbf{y}_k)^T \boldsymbol{\Gamma} (\beta \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \mathbf{y}_k) \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \mathbf{y}_k^T \left[\beta \mathbf{I} - \beta \mathbf{H} \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \beta \right] \mathbf{y}_k \\
&= \frac{1}{2} \mathbf{y}_k^T \left[\beta \mathbf{I} - \beta \mathbf{H} (\boldsymbol{\Phi} + \beta \mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \beta \right] \mathbf{y}_k \\
&= \frac{1}{2} \mathbf{y}_k^T \left[\beta^{-1} \mathbf{I} + \mathbf{H} \boldsymbol{\Phi}^{-1} \mathbf{H}^T \right]^{-1} \mathbf{y}_k. \tag{4.99}
\end{aligned}$$

In the equation above, we use the matrix inversion formula in Eq. (C.91). Using the model covariance matrix $\boldsymbol{\Sigma}_y$ defined in Eq. (4.25), we get

$$\Delta = \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \left[\beta^{-1} \mathbf{I} + \mathbf{H} \boldsymbol{\Phi}^{-1} \mathbf{H}^T \right]^{-1} \mathbf{y}_k = \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k.$$

The above equation is equal to Eq. (4.29).

4.10.3 Proof of Eq. (4.50)

The proof of Eq. (4.50) begins with

$$\hat{\mathbf{x}}_k = \underset{\mathbf{x}_k}{\operatorname{argmin}} \left[\beta \|\mathbf{y}_k - \mathbf{H}\mathbf{x}_k\|^2 + \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k \right]. \tag{4.100}$$

The solution of this minimization, $\hat{\mathbf{x}}_k$, is known as the weighted minimum-norm solution. To derive it, we define the cost function \mathcal{F} ,

$$\mathcal{F} = \beta \|\mathbf{y}_k - \mathbf{H}\mathbf{x}_k\|^2 + \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k.$$

Let us differentiate \mathcal{F} with respect to \mathbf{x}_k , and set it to zero,

$$\frac{\partial}{\partial \mathbf{x}_k} \mathcal{F} = -2\beta \mathbf{H}^T (\mathbf{y}_k - \mathbf{H}\mathbf{x}_k) + 2\boldsymbol{\Upsilon}^{-1} \mathbf{x}_k = 0.$$

Thus, the weighted minimum-norm solution is given by

$$\hat{\mathbf{x}}_k = \beta \left(\boldsymbol{\Upsilon}^{-1} + \beta \mathbf{H}^T \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{y}_k. \tag{4.101}$$

The above $\hat{\mathbf{x}}_k$ is exactly the same as the posterior mean $\bar{\mathbf{x}}_k$ in Eq. (4.14). Therefore, according to the arguments in the preceding subsection, we have the relationship,

$$\begin{aligned}
\min_{\mathbf{x}_k} \left[\beta \|\mathbf{y}_k - \mathbf{H}\mathbf{x}_k\|^2 + \mathbf{x}_k^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}_k \right] &= \left[\beta \|\mathbf{y}_k - \mathbf{H}\bar{\mathbf{x}}_k\|^2 + \bar{\mathbf{x}}_k^T \boldsymbol{\Upsilon}^{-1} \bar{\mathbf{x}}_k \right] \\
&= \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k. \tag{4.102}
\end{aligned}$$

References

1. D.P. Wipf, J.P. Owen, H.T. Attias, K. Sekihara, S.S. Nagarajan, Robust Bayesian estimation of the location, orientation, and time course of multiple correlated neural sources using MEG. *NeuroImage* **49**, 641–655 (2010)
2. D. Wipf, S. Nagarajan, A unified Bayesian framework for MEG/EEG source imaging. *NeuroImage* **44**(3), 947–966 (2009)
3. M.E. Tipping, Sparse Bayesian learning and relevance vector machine. *J. Mach. Learn. Res.* **1**, 211–244 (2001)
4. M.E. Tipping, Bayesian inference: an introduction to principles and practice in machine learning, in *Advanced Lectures on Machine Learning* (Springer, New York, 2004), pp. 41–62
5. D.J.C. MacKay, Bayesian interpolation. *Neural Comput.* **4**, 415–447 (1992)
6. M.I. Jordan, Z. Ghahramani, T.S. Jaakkola, L.K. Saul, An introduction to variational methods for graphical models. *Mach. Learn.* **37**, 183–233 (1999)
7. D.P. Wipf, S.S. Nagarajan, A new view of automatic relevance determination, in *Advances in Neural Information Processing Systems*, pp. 1625–1632 (2008)

Chapter 5

Bayesian Factor Analysis: A Versatile Framework for Denoising, Interference Suppression, and Source Localization

5.1 Introduction

This chapter describes Bayesian factor analysis (BFA), which is a technique that can decompose multiple sensor time courses into time courses of independent factor activities, where the number of factors is much smaller than the number of sensors. The factors are artificially-introduced “abstract” causes that explain the temporal behavior of the sensor data, and do not necessarily correspond to physical sources. Since the factor activities cannot directly be observed, they are considered latent variables.

The Bayesian factor analysis is a versatile framework. It can be used for selectively extracting signal components from noise-contaminated sensor data. It can provide an estimate of the data covariance matrix better than a sample covariance matrix, and such covariance estimate can be used in source reconstruction algorithms such as adaptive beamformers. It is extended to suppress interference components in interference-overlapped sensor data [1]. It is also extended to a virtual-sensor type source localization method, which is called the Saketini algorithm [2].

We start this chapter by explaining the basic form of the Bayesian factor analysis and then extend it to the variational Bayesian factor analysis (VBFA) [3] in which the model order (the number of factors) is determined by the algorithm itself. Following VBFA arguments, we describe the extensions for interference suppression and source localization.

5.2 Bayesian Factor Analysis

5.2.1 Factor Analysis Model

As in the previous chapters, let us define the output of the m th sensor at time t as $y_m(t)$, and the data vector as $\mathbf{y}(t)$, such that

$$\mathbf{y}(t) = \begin{bmatrix} y_1(t) \\ y_2(t) \\ \vdots \\ y_M(t) \end{bmatrix}, \quad (5.1)$$

where M is the total number of sensors. This column vector $\mathbf{y}(t)$ expresses the outputs of the whole sensor array. We assume that $\mathbf{y}(t)$ is measured at discrete time points $t = t_1, \dots, t_K$ where K is the total number of time points, and $\mathbf{y}(t_k)$ is denoted \mathbf{y}_k for simplicity. In the factor analysis model, the data vector \mathbf{y}_k is expressed as

$$\mathbf{y}_k = \mathbf{A}\mathbf{u}_k + \boldsymbol{\varepsilon}, \quad (5.2)$$

where $\boldsymbol{\varepsilon}$ is the sensor noise. On the right-hand side of the equation above, the first term, $\mathbf{A}\mathbf{u}_k$, expresses a signal component which is a product of an unknown matrix \mathbf{A} and unknown variables \mathbf{u}_k . The matrix \mathbf{A} is an $M \times L$ matrix, referred to as the mixing matrix, and \mathbf{u}_k is an L -dimensional column vector,

$$\mathbf{u}_k = \begin{bmatrix} u_1(t_k) \\ \vdots \\ u_L(t_k) \end{bmatrix}, \quad (5.3)$$

where the number of factors L is much smaller than the number of sensors M ; i.e., $L \ll M$. In the factor analysis model, the temporal variation of the sensor data \mathbf{y}_k is explained by the temporal variation of a small number of L independent factors, $u_1(t_k), \dots, u_L(t_k)$. Here, since the number of factors L specifies the model complexity, this number is called the model order. The factor is a latent variable, which cannot be directly observed, and each factor, in general, does not correspond to any physical source. The Bayesian factor analysis estimates the mixing matrix \mathbf{A} and the factor activity \mathbf{u}_k from the data \mathbf{y}_k where $k = 1, \dots, K$.

5.2.2 Probability Model

The prior probability distribution of the factor \mathbf{u}_k is assumed to be the zero-mean Gaussian with its precision matrix equal to the identity matrix, i.e.,

$$p(\mathbf{u}_k) = \mathcal{N}(\mathbf{u}_k | \mathbf{0}, \mathbf{I}). \quad (5.4)$$

The factor activity is assumed to be independent across time. Thus, the joint prior probability distribution for all time points is given by

$$p(\mathbf{u}) = p(\mathbf{u}_1, \dots, \mathbf{u}_K) = \prod_{k=1}^K p(\mathbf{u}_k) = \prod_{k=1}^K \mathcal{N}(\mathbf{u}_k | \mathbf{0}, \mathbf{I}), \quad (5.5)$$

where $\mathbf{u}_1, \dots, \mathbf{u}_K$ are collectively denoted \mathbf{u} . The noise $\boldsymbol{\varepsilon}$ is assumed to be Gaussian with the mean of zero, i.e.,

$$p(\boldsymbol{\varepsilon}) = \mathcal{N}(\boldsymbol{\varepsilon} | \mathbf{0}, \boldsymbol{\Lambda}^{-1}), \quad (5.6)$$

where $\boldsymbol{\Lambda}$ is a diagonal precision matrix. With this assumption, the conditional probability $p(\mathbf{y}_k | \mathbf{u}_k)$ is expressed as

$$p(\mathbf{y}_k | \mathbf{u}_k) = \mathcal{N}(\mathbf{y}_k | A\mathbf{u}_k, \boldsymbol{\Lambda}^{-1}). \quad (5.7)$$

The noise $\boldsymbol{\varepsilon}$ is also assumed to be independent across time. Thus, we have

$$p(\mathbf{y} | \mathbf{u}) = p(\mathbf{y}_1, \dots, \mathbf{y}_K | \mathbf{u}_1, \dots, \mathbf{u}_K) = \prod_{k=1}^K p(\mathbf{y}_k | \mathbf{u}_k) = \prod_{k=1}^K \mathcal{N}(\mathbf{y}_k | A\mathbf{u}_k, \boldsymbol{\Lambda}^{-1}), \quad (5.8)$$

where $\mathbf{y}_1, \dots, \mathbf{y}_K$ are collectively denoted \mathbf{y} . Using the probability distributions defined above, the Bayesian factor analysis can factorize the sensor data into L independent factor activity and additive sensor noise. This factorization is achieved using the EM algorithm [4, 5]. Explanation of the basics of the EM algorithm is provided in Sect. B.5 in the Appendix.

5.2.3 EM Algorithm

The E-step of the EM algorithm derives the posterior distribution $p(\mathbf{u}_k | \mathbf{y}_k)$. Derivation of the posterior distribution in the Gaussian model is described in Sect. B.3 in the Appendix. Since the posterior distribution is also Gaussian, we define the posterior distribution $p(\mathbf{u}_k | \mathbf{y}_k)$ such that

$$p(\mathbf{u}_k | \mathbf{y}_k) = \mathcal{N}(\mathbf{u}_k | \bar{\mathbf{u}}_k, \boldsymbol{\Gamma}^{-1}), \quad (5.9)$$

where $\bar{\mathbf{u}}_k$ is the mean and $\boldsymbol{\Gamma}$ is the precision matrix. Using Eqs. (B.24) and (B.25) with setting $\boldsymbol{\Phi}$ to \mathbf{I} and \mathbf{H} to A in these equations, we get

$$\boldsymbol{\Gamma} = (A^T \boldsymbol{\Lambda} A + \mathbf{I}), \quad (5.10)$$

$$\bar{\mathbf{u}}_k = (A^T \boldsymbol{\Lambda} A + \mathbf{I})^{-1} A^T \boldsymbol{\Lambda} \mathbf{y}_k. \quad (5.11)$$

The equations above are the E-step update equations in the Bayesian factor analysis.

Let us derive the M-step update equations for the mixing matrix A and the noise precision matrix $\boldsymbol{\Lambda}$. To do so, we derive the average log likelihood, $\Theta(A, \boldsymbol{\Lambda})$, and according to Eqs. (B.34), (5.5), and (5.8), it is expressed as

$$\begin{aligned}
\Theta(\mathbf{A}, \boldsymbol{\Lambda}) &= E_{\mathbf{u}} [\log p(\mathbf{y}_1, \dots, \mathbf{y}_K, \mathbf{u}_1, \dots, \mathbf{u}_K)] \\
&= E_{\mathbf{u}} [\log p(\mathbf{y}_1, \dots, \mathbf{y}_K | \mathbf{u}_1, \dots, \mathbf{u}_K)] + E_{\mathbf{u}} [\log p(\mathbf{u}_1, \dots, \mathbf{u}_K)] \\
&= E_{\mathbf{u}} \left[\sum_{k=1}^K \log p(\mathbf{y}_k | \mathbf{u}_k) \right] + E_{\mathbf{u}} \left[\sum_{k=1}^K \log p(\mathbf{u}_k) \right], \tag{5.12}
\end{aligned}$$

where $E_{\mathbf{u}}[\cdot]$ indicates the expectation with respect to the posterior probability $p(\mathbf{u}|\mathbf{y})$. Taking a look at Eqs. (5.4) and (5.7), only the first term on the right-hand side of Eq. (5.12) contains \mathbf{A} and $\boldsymbol{\Lambda}$. Thus, we have

$$\begin{aligned}
\Theta(\mathbf{A}, \boldsymbol{\Lambda}) &= E_{\mathbf{u}} \left[\sum_{k=1}^K [\log p(\mathbf{y}_k | \mathbf{u}_k)] \right] + \mathcal{C} \\
&= \frac{K}{2} \log |\boldsymbol{\Lambda}| - \frac{1}{2} E_{\mathbf{u}} \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k) \right] + \mathcal{C}, \tag{5.13}
\end{aligned}$$

where \mathcal{C} expresses terms not containing \mathbf{A} and $\boldsymbol{\Lambda}$.

The derivative of $\Theta(\mathbf{A}, \boldsymbol{\Lambda})$ with respect to \mathbf{A} is given by

$$\frac{\partial \Theta(\mathbf{A}, \boldsymbol{\Lambda})}{\partial \mathbf{A}} = \boldsymbol{\Lambda} E_{\mathbf{u}} \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k) \mathbf{u}_k^T \right] = \boldsymbol{\Lambda} (\mathbf{R}_{yu} - \mathbf{A}\mathbf{R}_{uu}), \tag{5.14}$$

where

$$\mathbf{R}_{uu} = E_{\mathbf{u}} \left[\sum_{k=1}^K \mathbf{u}_k \mathbf{u}_k^T \right] = \sum_{k=1}^K \bar{\mathbf{u}}_k \bar{\mathbf{u}}_k^T + K \boldsymbol{\Gamma}^{-1}, \tag{5.15}$$

$$\mathbf{R}_{yu} = E_{\mathbf{u}} \left[\sum_{k=1}^K \mathbf{y}_k \mathbf{u}_k^T \right] = \sum_{k=1}^K \mathbf{y}_k \bar{\mathbf{u}}_k^T, \tag{5.16}$$

$$\mathbf{R}_{uy} = \mathbf{R}_{yu}^T. \tag{5.17}$$

Setting the right-hand side of Eq. (5.14) to zero gives $\mathbf{R}_{yu} = \mathbf{A}\mathbf{R}_{uu}$. That is, the M-step update equation for \mathbf{A} is derived as

$$\mathbf{A} = \mathbf{R}_{yu} \mathbf{R}_{uu}^{-1}. \tag{5.18}$$

Next, the update equation for $\boldsymbol{\Lambda}$ is derived. The partial derivative of $\Theta(\mathbf{A}, \boldsymbol{\Lambda})$ with respect to $\boldsymbol{\Lambda}$ is expressed as

$$\frac{\partial \Theta(\mathbf{A}, \boldsymbol{\Lambda})}{\partial \boldsymbol{\Lambda}} = \frac{K}{2} \boldsymbol{\Lambda}^{-1} - \frac{1}{2} E_{\mathbf{u}} \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k)(\mathbf{y}_k - \mathbf{A}\mathbf{u}_k)^T \right]. \tag{5.19}$$

Setting the right-hand side to zero gives

$$\boldsymbol{\Lambda}^{-1} = \frac{1}{K} E_{\boldsymbol{u}} \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\boldsymbol{u}_k)(\mathbf{y}_k - \mathbf{A}\boldsymbol{u}_k)^T \right]. \quad (5.20)$$

The right-hand side is further changed to

$$\begin{aligned} & E_{\boldsymbol{u}} \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\boldsymbol{u}_k)(\mathbf{y}_k - \mathbf{A}\boldsymbol{u}_k)^T \right] \\ &= \sum_{k=1}^K E_{\boldsymbol{u}}[\mathbf{y}_k \mathbf{y}_k^T] - \sum_{k=1}^K E_{\boldsymbol{u}}[\mathbf{y}_k \boldsymbol{u}_k^T] \mathbf{A}^T - \sum_{k=1}^K \mathbf{A} E_{\boldsymbol{u}}[\boldsymbol{u}_k \mathbf{y}_k^T] + \sum_{k=1}^K \mathbf{A} E_{\boldsymbol{u}}[\boldsymbol{u}_k \boldsymbol{u}_k^T] \mathbf{A}^T \\ &= \sum_{k=1}^K \mathbf{y}_k \mathbf{y}_k^T - \sum_{k=1}^K \mathbf{y}_k \bar{\boldsymbol{u}}_k^T \mathbf{A}^T - \sum_{k=1}^K \mathbf{A} \bar{\boldsymbol{u}}_k \mathbf{y}_k^T + \sum_{k=1}^K \mathbf{A} E_{\boldsymbol{u}}[\boldsymbol{u}_k \boldsymbol{u}_k^T] \mathbf{A}^T \\ &= \mathbf{R}_{yy} - \mathbf{R}_{yu} \mathbf{A}^T - \mathbf{A} \mathbf{R}_{uy} + \mathbf{A} \mathbf{R}_{uu} \mathbf{A}^T, \end{aligned} \quad (5.21)$$

where $\mathbf{R}_{yy} = \sum_{k=1}^K \mathbf{y}_k \mathbf{y}_k^T$. Using $\mathbf{R}_{yu} = \mathbf{A} \mathbf{R}_{uu}$ derived from Eq.(5.18), we get

$$\frac{1}{K} E_{\boldsymbol{u}} \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\boldsymbol{u}_k)(\mathbf{y}_k - \mathbf{A}\boldsymbol{u}_k)^T \right] = \frac{1}{K} [\mathbf{R}_{yy} - \mathbf{A} \mathbf{R}_{uy}]. \quad (5.22)$$

Considering that $\boldsymbol{\Lambda}^{-1}$ is a diagonal matrix, we finally derive

$$\boldsymbol{\Lambda}^{-1} = \frac{1}{K} \text{diag}(\mathbf{R}_{yy} - \mathbf{A} \mathbf{R}_{uy}), \quad (5.23)$$

where $\text{diag}(\cdot)$ indicates a diagonal matrix obtained using the diagonal entries of a matrix between the brackets.

5.2.4 Computation of Marginal Likelihood

Let us derive an expression to compute the marginal likelihood $p(\mathbf{y}|\boldsymbol{\theta})$, which can be used for monitoring the progress of the EM algorithm. Here, $\boldsymbol{\theta}$ is used to collectively express the hyperparameters \mathbf{A} and $\boldsymbol{\Lambda}$. Using the marginalization, we obtain $p(\mathbf{y}|\boldsymbol{\theta})$, such that

$$p(\mathbf{y}|\boldsymbol{\theta}) = \int_{-\infty}^{\infty} p(\mathbf{y}|\boldsymbol{u}, \boldsymbol{\theta}) p(\boldsymbol{u}) d\boldsymbol{u}. \quad (5.24)$$

Based on the arguments in Sect. B.4 of the Appendix, $p(\mathbf{y}|\boldsymbol{\theta})$ is derived such that

$$p(\mathbf{y}|\boldsymbol{\theta}) = \prod_{k=1}^K p(\mathbf{y}_k|\boldsymbol{\theta}) \quad \text{where} \quad p(\mathbf{y}_k|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}_k|\mathbf{0}, \boldsymbol{\Sigma}_y), \quad (5.25)$$

where $\boldsymbol{\Sigma}_y$ is the model data covariance. Setting \mathbf{A} to \mathbf{H} and \mathbf{I} to $\boldsymbol{\Phi}^{-1}$ in Eq. (B.30), we derive $\boldsymbol{\Sigma}_y$, such that

$$\boldsymbol{\Sigma}_y = \boldsymbol{\Lambda}^{-1} + \mathbf{A}\mathbf{A}^T. \quad (5.26)$$

According to Eq. (B.29), the log likelihood function is obtained as

$$\log p(\mathbf{y}|\boldsymbol{\theta}) = \sum_{k=1}^K p(\mathbf{y}_k|\boldsymbol{\theta}) = \frac{K}{2} \log |\boldsymbol{\Sigma}_y^{-1}| - \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k. \quad (5.27)$$

Using the matrix inversion formula in Eq. (C.91), we have

$$\begin{aligned} \boldsymbol{\Sigma}_y^{-1} &= \left(\boldsymbol{\Lambda}^{-1} + \mathbf{A}\mathbf{A}^T \right)^{-1} \\ &= \boldsymbol{\Lambda} - \boldsymbol{\Lambda}\mathbf{A} \left(\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} + \mathbf{I} \right)^{-1} \mathbf{A}^T \boldsymbol{\Lambda} = \boldsymbol{\Lambda} - \boldsymbol{\Lambda}\mathbf{A}\boldsymbol{\Gamma}^{-1}\mathbf{A}^T\boldsymbol{\Lambda}, \end{aligned} \quad (5.28)$$

Using the equation above and Eq. (5.11), we get

$$\mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k = \mathbf{y}_k^T \left[\boldsymbol{\Lambda} - \boldsymbol{\Lambda}\mathbf{A}\boldsymbol{\Gamma}^{-1}\mathbf{A}^T\boldsymbol{\Lambda} \right] \mathbf{y}_k = \mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{y}_k - \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{u}}_k.$$

Also, since the relationship

$$|\boldsymbol{\Sigma}_y^{-1}| = |\boldsymbol{\Lambda} - \boldsymbol{\Lambda}\mathbf{A}\boldsymbol{\Gamma}^{-1}\mathbf{A}^T\boldsymbol{\Lambda}| = |\mathbf{I} - \boldsymbol{\Lambda}\mathbf{A}\boldsymbol{\Gamma}^{-1}\mathbf{A}^T| |\boldsymbol{\Lambda}| = |\boldsymbol{\Gamma}^{-1}| |\boldsymbol{\Lambda}| \quad (5.29)$$

holds,¹ by substituting the above two equations into Eq. (5.27), we finally obtain the expression for the likelihood:

$$\log p(\mathbf{y}|\boldsymbol{\theta}) = \sum_{k=1}^K p(\mathbf{y}_k|\boldsymbol{\theta}) = \frac{K}{2} \log \frac{|\boldsymbol{\Lambda}|}{|\boldsymbol{\Gamma}|} - \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{y}_k + \frac{1}{2} \sum_{k=1}^K \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{u}}_k. \quad (5.30)$$

A more informative derivation of Eq. (5.30) uses the free energy expression in Eq. (B.63) from the Appendix, which claims the following relationship holds:

¹ Defining $\mathbf{C} = \mathbf{I} - \boldsymbol{\Lambda}\mathbf{A}\boldsymbol{\Gamma}^{-1}\mathbf{A}^T$, we have

$$\mathbf{A}^T \mathbf{C} = \mathbf{A}^T - \mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} \boldsymbol{\Gamma}^{-1} \mathbf{A}^T = \mathbf{A}^T - (\boldsymbol{\Gamma} - \mathbf{I}) \boldsymbol{\Gamma}^{-1} \mathbf{A}^T = \boldsymbol{\Gamma}^{-1} \mathbf{A}^T.$$

Taking the determinant of the equation above, we get $|\mathbf{C}| = |\boldsymbol{\Gamma}^{-1}|$.

$$\begin{aligned}\log p(\mathbf{y}|\boldsymbol{\theta}) &= \int d\mathbf{u} p(\mathbf{u}|\mathbf{y})[\log p(\mathbf{u}, \mathbf{y}|\boldsymbol{\theta}) - \log p(\mathbf{u}|\mathbf{y})] \\ &= E_{\mathbf{u}} [\log p(\mathbf{u}, \mathbf{y}|\boldsymbol{\theta})] + \mathcal{H}[p(\mathbf{u}|\mathbf{y})],\end{aligned}\quad (5.31)$$

where the second term on the right-hand side is the entropy regarding $p(\mathbf{u}|\mathbf{y})$. Using Eqs. (5.9) and (C.9) in the Appendix, we get

$$\mathcal{H}[p(\mathbf{u}|\mathbf{y})] = \sum_{k=1}^K \mathcal{H}[p(\mathbf{u}_k|\mathbf{y}_k)] = -\frac{K}{2} \log |\boldsymbol{\Gamma}|. \quad (5.32)$$

Using Eqs. (5.5) and (5.8), the first term in Eq. (5.31) can be rewritten as

$$\begin{aligned}E_{\mathbf{u}} [\log p(\mathbf{u}, \mathbf{y}|\boldsymbol{\theta})] &= E_{\mathbf{u}} [\log p(\mathbf{y}|\mathbf{u}, \boldsymbol{\theta}) + \log p(\mathbf{u})] \\ &= -\frac{1}{2} \sum_{k=1}^K E_{\mathbf{u}} \left[(\mathbf{y}_k - \mathbf{A}\mathbf{u}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k) + \mathbf{u}_k^T \mathbf{u}_k \right] \\ &\quad + \frac{K}{2} \log |\boldsymbol{\Lambda}|,\end{aligned}\quad (5.33)$$

where constant terms are omitted. On the right-hand side of the equation above,

$$\begin{aligned}E_{\mathbf{u}} [\mathbf{u}_k^T \mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} \mathbf{u}_k + \mathbf{u}_k^T \mathbf{u}_k] &= E_{\mathbf{u}} [\mathbf{u}_k^T (\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} + \mathbf{I}) \mathbf{u}_k] = E_{\mathbf{u}} [\mathbf{u}_k^T \boldsymbol{\Gamma} \mathbf{u}_k] \\ &= E_{\mathbf{u}} [\text{tr}(\mathbf{u}_k \mathbf{u}_k^T \boldsymbol{\Gamma})] = \text{tr}[E_{\mathbf{u}} (\mathbf{u}_k \mathbf{u}_k^T) \boldsymbol{\Gamma}] \\ &= \text{tr}[(\bar{\mathbf{u}}_k \bar{\mathbf{u}}_k^T + \boldsymbol{\Gamma}^{-1}) \boldsymbol{\Gamma}] = \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{u}}_k,\end{aligned}\quad (5.34)$$

where a constant term is again omitted. Also, we can show that

$$\begin{aligned}E_{\mathbf{u}} [\mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{A} \mathbf{u}_k + \mathbf{u}_k^T \mathbf{A}^T \boldsymbol{\Lambda} \mathbf{y}_k] &= \mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{A} \bar{\mathbf{u}}_k + \bar{\mathbf{u}}_k^T \mathbf{A}^T \boldsymbol{\Lambda} \mathbf{y}_k \\ &= \mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{A} \boldsymbol{\Gamma}^{-1} \boldsymbol{\Gamma} \bar{\mathbf{u}}_k + \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \boldsymbol{\Gamma}^{-1} \mathbf{A}^T \boldsymbol{\Lambda} \mathbf{y}_k = 2\bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{u}}_k,\end{aligned}\quad (5.35)$$

where Eq. (5.11) is used. Substituting the above four equations into Eq. (5.31), we can get Eq. (5.30).

5.2.5 Summary of the BFA Algorithm

Let us summarize the BFA algorithm. The algorithm is based on the factor analysis model presented in Eq. (5.2). The algorithm estimates the factor activity \mathbf{u}_k , the mixing matrix \mathbf{A} , and the sensor-noise precision matrix $\boldsymbol{\Lambda}$. In the estimation, first

the model order L is set according to some prior knowledge, and appropriate initial values are given to \mathbf{A} and $\boldsymbol{\Lambda}$. In the E-step, $\boldsymbol{\Gamma}$ and $\bar{\mathbf{u}}_k$ ($k = 1, \dots, L$) are updated according to Eqs. (5.10) and (5.11). In the M-step, parameters \mathbf{A} and $\boldsymbol{\Lambda}$ are updated using Eqs. (5.18) and (5.23). The marginal likelihood $p(\mathbf{y}|\mathbf{A}, \boldsymbol{\Lambda})$ in Eq. (5.30) can be used for monitoring the progress of the EM iteration. If the likelihood increase becomes very small with respect to the iteration count, the EM iteration may be stopped.

The BFA can be applied to denoising the sensor data. Using the BFA algorithm, the estimate of the signal component, $\hat{\mathbf{y}}_k^S$, is given by

$$\hat{\mathbf{y}}_k^S = E_{\mathbf{u}}[\mathbf{A}\mathbf{u}_k] = \mathbf{A}E_{\mathbf{u}}[\mathbf{u}_k] = \mathbf{A}\bar{\mathbf{u}}_k. \quad (5.36)$$

This $\hat{\mathbf{y}}_k^S$ may be used for further analysis such as source localization. We can compute the sample data covariance using only the signal component $\mathbf{A}\mathbf{u}_k$, which is

$$\begin{aligned} \bar{\mathbf{R}} &= \frac{1}{K} \sum_{k=1}^K E_{\mathbf{u}}[(\mathbf{A}\mathbf{u}_k)(\mathbf{A}\mathbf{u}_k)^T] = \frac{1}{K} \sum_{k=1}^K \mathbf{A}E_{\mathbf{u}}[\mathbf{u}_k\mathbf{u}_k^T]\mathbf{A}^T \\ &= \frac{1}{K} \sum_{k=1}^K \mathbf{A} \left[\bar{\mathbf{u}}_k \bar{\mathbf{u}}_k^T + \boldsymbol{\Gamma}^{-1} \right] \mathbf{A}^T = \mathbf{A} \left[\frac{1}{K} \sum_{k=1}^K \bar{\mathbf{u}}_k \bar{\mathbf{u}}_k^T \right] \mathbf{A}^T + \mathbf{A}\boldsymbol{\Gamma}^{-1}\mathbf{A}^T. \end{aligned} \quad (5.37)$$

This $\bar{\mathbf{R}}$ can be used in source imaging algorithms such as adaptive beamformers. Note that, on the right-hand side of Eq. (5.37), the second term $\mathbf{A}\boldsymbol{\Gamma}^{-1}\mathbf{A}^T$ works as a regularization term and $\bar{\mathbf{R}}$ has a form in which the regularization is already incorporated.

5.3 Variational Bayes Factor Analysis (VBFA)

5.3.1 Prior Distribution for Mixing Matrix

In the BFA algorithm described in the preceding section, a user must determine the model order L , according to prior knowledge of the measurement. However, determination of the model order is not easy in most practical applications. We describe here an extension of the Bayesian factor analysis based on the variational Bayesian method [6]. The method is called the variational Bayesian factor analysis (VBFA), in which the model order determination is embedded in the algorithm. On the basis of the factor analysis model in Eq. (5.2), the VBFA algorithm estimates the posterior probability distributions not only for the factor activity \mathbf{u}_k but also for the mixing matrix \mathbf{A} .

For convenience in the following arguments, we define the j th row of the mixing matrix \mathbf{A} as the column vector \mathbf{a}_j , i.e., the $M \times L$ matrix \mathbf{A} is expressed as

$$\mathbf{A} = \begin{bmatrix} A_{1,1} & \dots & A_{1,L} \\ A_{2,1} & \dots & A_{2,L} \\ \vdots & \ddots & \vdots \\ A_{M,1} & \dots & A_{M,L} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_M^T \end{bmatrix}, \quad (5.38)$$

where

$$\mathbf{a}_j = [A_{j,1}, \dots, A_{j,L}]^T. \quad (5.39)$$

The prior distribution of \mathbf{a}_j is assumed to be²

$$p(\mathbf{a}_j) = \mathcal{N}(\mathbf{a}_j | \mathbf{0}, (\lambda_j \boldsymbol{\alpha})^{-1}), \quad (5.40)$$

where the j th diagonal element of the noise precision matrix $\boldsymbol{\Lambda}$ is denoted λ_j , and $\boldsymbol{\alpha}$ is a diagonal matrix given by

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 & 0 & \dots & 0 \\ 0 & \alpha_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \alpha_L \end{bmatrix}. \quad (5.41)$$

Here, note that we assume that elements of \mathbf{a}_j are statistically independent because the precision matrix in Eq. (5.40) is diagonal. We further assume that \mathbf{a}_i and \mathbf{a}_j are statistically independent when $i \neq j$. Thus, the prior probability distribution for the whole mixing matrix \mathbf{A} is expressed as

$$p(\mathbf{A}) = \prod_{j=1}^M \mathcal{N}(\mathbf{a}_j | \mathbf{0}, (\lambda_j \boldsymbol{\alpha})^{-1}). \quad (5.42)$$

This equation is equivalent to

$$p(\mathbf{A}) = \prod_{j=1}^M \prod_{\ell=1}^L \mathcal{N}(A_{j,\ell} | 0, (\lambda_j \alpha_\ell)^{-1}). \quad (5.43)$$

² In the prior distribution in Eq. (5.40), the precision matrix has a form of a diagonal matrix $\boldsymbol{\alpha}$ multiplied by a scalar λ_j . The inclusion of this scalar, λ_j , is just for convenience in the mathematical expressions for the update equations of Ψ and $\bar{\mathbf{a}}_j$. The inclusion of λ_j actually makes these update equations significantly simpler.

Since $p(\mathbf{A})$ is a so-called conjugate prior, the posterior probability distribution has the form of the Gaussian distribution:

$$p(\mathbf{a}_j|\mathbf{y}) = \mathcal{N}(\mathbf{a}_j|\bar{\mathbf{a}}_j, (\lambda_j \boldsymbol{\Psi})^{-1}),$$

and

$$p(\mathbf{A}|\mathbf{y}) = p(\mathbf{a}_1, \dots, \mathbf{a}_M|\mathbf{y}) = \prod_{j=1}^M p(\mathbf{a}_j|\mathbf{y}) = \prod_{j=1}^M \mathcal{N}(\mathbf{a}_j|\bar{\mathbf{a}}_j, (\lambda_j \boldsymbol{\Psi})^{-1}), \quad (5.44)$$

where $\bar{\mathbf{a}}_j$ and $\lambda_j \boldsymbol{\Psi}$ are the mean and precision matrix of the posterior distribution. Namely, the posterior distribution $p(\mathbf{A}|\mathbf{y})$ has a form identical to the prior distribution $p(\mathbf{A})$ with the diagonal $\boldsymbol{\alpha}$ replaced by the non-diagonal $\boldsymbol{\Psi}$. We define for later use the matrix $\bar{\mathbf{A}}$, such that

$$\bar{\mathbf{A}} = \begin{bmatrix} \bar{\mathbf{a}}_1^T \\ \bar{\mathbf{a}}_2^T \\ \vdots \\ \bar{\mathbf{a}}_M^T \end{bmatrix} = \begin{bmatrix} \bar{A}_{1,1} & \dots & \bar{A}_{1,L} \\ \bar{A}_{2,1} & \dots & \bar{A}_{2,L} \\ \vdots & \ddots & \vdots \\ \bar{A}_{M,1} & \dots & \bar{A}_{M,L} \end{bmatrix}. \quad (5.45)$$

In the VBFA algorithm, an overspecified model order L is used, i.e., the value of L is set greater than the true model order L_0 . The posterior mean of \mathbf{a}_j , $\bar{\mathbf{a}}_j$:

$$\bar{\mathbf{a}}_j = [\bar{A}_{j,1}, \dots, \bar{A}_{j,L_0}, \bar{A}_{j,L_0+1}, \dots, \bar{A}_{j,L}]^T$$

is the Bayes estimate of the j th row of the mixing matrix. In this estimated mixing matrix, the matrix elements $\bar{A}_{j,L_0+1}, \dots, \bar{A}_{j,L}$ are those corresponding to non-existing factors. In the VBFA algorithm, those elements are estimated to be significantly small, and the influence of the overspecified factors are automatically eliminated in the final estimation results.

5.3.2 Variational Bayes EM Algorithm (VBEM)

5.3.2.1 E-Step

We follow the arguments in Sect. B.6 in the Appendix, and derive the variational Bayes EM algorithm. The posterior distribution $p(\mathbf{u}, \mathbf{A}|\mathbf{y})$ is approximated by

$$p(\mathbf{u}, \mathbf{A}|\mathbf{y}) = p(\mathbf{u}|\mathbf{y})p(\mathbf{A}|\mathbf{y}), \quad (5.46)$$

The E-step of the VBEM algorithm is a step that estimates $p(\mathbf{u}|\mathbf{y})$, and according to the arguments in Sect. B.6, the estimate of the posterior distribution, $\widehat{p}(\mathbf{u}|\mathbf{y})$, is obtained as

$$\log \widehat{p}(\mathbf{u}|\mathbf{y}) = E_A [\log p(\mathbf{u}, \mathbf{y}, A)], \quad (5.47)$$

where $E_A [\cdot]$ indicates the expectation with respect to $p(A|\mathbf{y})$. To obtain $\widehat{p}(\mathbf{u}|\mathbf{y})$, we substitute

$$p(\mathbf{u}, \mathbf{y}, A) = p(\mathbf{y}|\mathbf{u}, A)p(\mathbf{u}, A) = p(\mathbf{y}|\mathbf{u}, A)p(\mathbf{u})p(A), \quad (5.48)$$

into Eq. (5.47). Neglecting constant terms, we can derive,

$$\begin{aligned} \log \widehat{p}(\mathbf{u}|\mathbf{y}) &= E_A [\log p(\mathbf{y}|\mathbf{u}, A) + \log p(\mathbf{u}) + \log p(A)] \\ &= E_A [\log p(\mathbf{y}|\mathbf{u}, A) + \log p(\mathbf{u})] \\ &= \sum_{k=1}^K E_A [\log p(\mathbf{y}_k|\mathbf{u}_k, A) + \log p(\mathbf{u}_k)]. \end{aligned} \quad (5.49)$$

In the equation above, the term $\log p(A)$ is omitted because it does not contain \mathbf{u} .

Since we have assumed that the prior and the noise probability distributions are independent across time, we have the independence of the posterior with respect to time, i.e.,

$$\widehat{p}(\mathbf{u}|\mathbf{y}) = \prod_{k=1}^K \widehat{p}(\mathbf{u}_k|\mathbf{y}_k). \quad (5.50)$$

Using Eqs. (5.49) and (5.50), we obtain

$$\log \widehat{p}(\mathbf{u}_k|\mathbf{y}_k) = E_A [\log p(\mathbf{y}_k|\mathbf{u}_k, A) + \log p(\mathbf{u}_k)]. \quad (5.51)$$

Substituting Eqs. (5.4) and (5.7) into (5.51), we get

$$\log \widehat{p}(\mathbf{u}_k|\mathbf{y}_k) = E_A \left[-\frac{1}{2}(\mathbf{y}_k - A\mathbf{u}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - A\mathbf{u}_k) - \frac{1}{2}\mathbf{u}_k^T \mathbf{u}_k \right]. \quad (5.52)$$

Since the posterior $\widehat{p}(\mathbf{u}_k|\mathbf{y}_k)$ is Gaussian, we assume

$$\widehat{p}(\mathbf{u}_k|\mathbf{y}_k) = \mathcal{N}(\mathbf{u}_k|\bar{\mathbf{u}}_k, \boldsymbol{\Gamma}^{-1}), \quad (5.53)$$

where $\bar{\mathbf{u}}_k$ and $\boldsymbol{\Gamma}$ are the mean and the precision of this posterior distribution.

The E-step of the VBEM algorithm estimates $\hat{p}(\mathbf{u}_k | \mathbf{y}_k)$, i.e., it estimates $\bar{\mathbf{u}}_k$, and $\boldsymbol{\Lambda}$. For this estimation, we compute the derivative $\frac{\partial}{\partial \mathbf{u}_k} \log \hat{p}(\mathbf{u}_k | \mathbf{y}_k)$. Using Eq. (5.52), the derivative is given by

$$\begin{aligned}\frac{\partial}{\partial \mathbf{u}_k} \log \hat{p}(\mathbf{u}_k | \mathbf{y}_k) &= E_{\mathbf{A}} \left[\mathbf{A}^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{A} \mathbf{u}_k) \right] - \mathbf{u}_k \\ &= \bar{\mathbf{A}}^T \boldsymbol{\Lambda} \mathbf{y}_k - E_{\mathbf{A}} \left[\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} \right] \mathbf{u}_k - \mathbf{u}_k,\end{aligned}\quad (5.54)$$

where $E_{\mathbf{A}}(\mathbf{A}) = \bar{\mathbf{A}}$, i.e., $\bar{\mathbf{A}}$ (defined in Eq. (5.45)) is the mean of the posterior $p(\mathbf{A} | \mathbf{y})$. Let us compute $E_{\mathbf{A}} [\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A}]$ contained in the right-hand side of Eq. (5.54). Noting that the diagonal elements of $\boldsymbol{\Lambda}$ are denoted $\lambda_1, \dots, \lambda_M$, $\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A}$ is rewritten as

$$\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_M] \begin{bmatrix} \lambda_1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & \lambda_M \end{bmatrix} \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_M^T \end{bmatrix} = \sum_{j=1}^M \lambda_j \mathbf{a}_j \mathbf{a}_j^T,$$

giving

$$E_{\mathbf{A}} \left[\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} \right] = \sum_{j=1}^M \lambda_j E_{\mathbf{A}} \left[\mathbf{a}_j \mathbf{a}_j^T \right]. \quad (5.55)$$

Since the precision matrix of the posterior distribution of \mathbf{a}_j is $\lambda_j \boldsymbol{\Psi}$, we have

$$E_{\mathbf{A}} \left[\mathbf{a}_j \mathbf{a}_j^T \right] = \bar{\mathbf{a}}_j \bar{\mathbf{a}}_j^T + \frac{1}{\lambda_j} \boldsymbol{\Psi}^{-1}.$$

Therefore, the relationship

$$\begin{aligned}E_{\mathbf{A}} \left[\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} \right] &= \sum_{j=1}^M \lambda_j \left(\bar{\mathbf{a}}_j \bar{\mathbf{a}}_j^T + \frac{1}{\lambda_j} \boldsymbol{\Psi}^{-1} \right) \\ &= \sum_{j=1}^M \lambda_j \bar{\mathbf{a}}_j \bar{\mathbf{a}}_j^T + M \boldsymbol{\Psi}^{-1} = \bar{\mathbf{A}}^T \boldsymbol{\Lambda} \bar{\mathbf{A}} + M \boldsymbol{\Psi}^{-1}\end{aligned}\quad (5.56)$$

holds.

Substituting Eq. (5.56) into (5.54), we get

$$\frac{\partial}{\partial \mathbf{u}_k} \log \hat{p}(\mathbf{u}_k | \mathbf{y}_k) = \bar{\mathbf{A}}^T \boldsymbol{\Lambda} \mathbf{y}_k - \left[\bar{\mathbf{A}}^T \boldsymbol{\Lambda} \bar{\mathbf{A}} + M \boldsymbol{\Psi}^{-1} + \mathbf{I} \right] \mathbf{u}_k. \quad (5.57)$$

The posterior precision matrix $\boldsymbol{\Gamma}$ is obtained as the coefficient of \mathbf{u}_k in the right-hand side of Eq. (5.57), i.e., $\boldsymbol{\Gamma}$ is given by

$$\boldsymbol{\Gamma} = \bar{\mathbf{A}}^T \boldsymbol{\Lambda} \bar{\mathbf{A}} + M \boldsymbol{\Psi}^{-1} + \mathbf{I}. \quad (5.58)$$

The mean $\bar{\mathbf{u}}_k$ is obtained as the value of \mathbf{u}_k that makes the right-hand side of Eq. (5.57) equal to zero. That is, we have

$$\bar{\mathbf{u}}_k = \left[\bar{\mathbf{A}}^T \boldsymbol{\Lambda} \bar{\mathbf{A}} + M \boldsymbol{\Psi}^{-1} + \mathbf{I} \right]^{-1} \bar{\mathbf{A}}^T \boldsymbol{\Lambda} \mathbf{y}_k. \quad (5.59)$$

5.3.2.2 M-Step

The M-step of the VBEM algorithm estimates $p(\mathbf{A}|\mathbf{y})$. According to the arguments in Sect. B.6 in the Appendix, the estimate of the posterior distribution, $\hat{p}(\mathbf{A}|\mathbf{y})$, is obtained as

$$\log \hat{p}(\mathbf{A}|\mathbf{y}) = E_{\mathbf{u}} [\log p(\mathbf{u}, \mathbf{y}, \mathbf{A})], \quad (5.60)$$

where $E_{\mathbf{u}} [\cdot]$ indicates the expectation with respect to the posterior distribution, $p(\mathbf{u}|\mathbf{y})$. To obtain $\hat{p}(\mathbf{A}|\mathbf{y})$, we substitute Eq. (5.48) into (5.60). Neglecting terms (such as $\log p(\mathbf{u})$) that do not contain \mathbf{A} , we derive,

$$\log \hat{p}(\mathbf{A}|\mathbf{y}) = E_{\mathbf{u}} [\log p(\mathbf{y}|\mathbf{u}, \mathbf{A}) + \log p(\mathbf{A})]. \quad (5.61)$$

Substituting Eqs. (5.8) and (5.42) into (5.61), omitting terms not containing \mathbf{A} , we get

$$\log \hat{p}(\mathbf{A}|\mathbf{y}) = \sum_{k=1}^K E_{\mathbf{u}} \left[-\frac{1}{2} (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k) \right] - \frac{1}{2} \sum_{j=1}^M \lambda_j \mathbf{a}_j^T \boldsymbol{\alpha} \mathbf{a}_j. \quad (5.62)$$

As in Eq. (5.44), the posterior $p(\mathbf{a}_j|\mathbf{y})$ is Gaussian with the mean $\bar{\mathbf{a}}_j$ and the precision $\lambda_j \boldsymbol{\Psi}$. The precision $\lambda_j \boldsymbol{\Psi}$ can be obtained by the coefficient of \mathbf{a}_j in the derivative $\frac{\partial}{\partial \mathbf{A}} \log \hat{p}(\mathbf{A}|\mathbf{y})$, and $\bar{\mathbf{a}}_j$ can be obtained as \mathbf{a}_j that makes this derivative equal to zero. To compute the derivative, we first rewrite:

$$\sum_{j=1}^M \lambda_j \mathbf{a}_j^T \boldsymbol{\alpha} \mathbf{a}_j = \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_{\ell} A_{j,\ell}^2, \quad (5.63)$$

and we can find

$$\frac{\partial}{\partial A_{j,\ell}} \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_{\ell} A_{j,\ell}^2 = \lambda_j \alpha_{\ell} A_{j,\ell} = [\boldsymbol{\Lambda} \mathbf{A} \boldsymbol{\alpha}]_{j,\ell}. \quad (5.64)$$

Therefore, we have

$$\frac{\partial}{\partial \Lambda} \frac{1}{2} \sum_{j=1}^M \lambda_j \mathbf{a}_j^T \boldsymbol{\alpha} \mathbf{a}_j = \Lambda \Lambda \boldsymbol{\alpha}. \quad (5.65)$$

We also have

$$\frac{\partial}{\partial \Lambda} \sum_{k=1}^K -\frac{1}{2} (\mathbf{y}_k - \Lambda \mathbf{u}_k)^T \Lambda (\mathbf{y}_k - \Lambda \mathbf{u}_k) = \Lambda \sum_{k=1}^K (\mathbf{y}_k - \Lambda \mathbf{u}_k) \mathbf{u}_k^T. \quad (5.66)$$

Consequently, we can derive

$$\begin{aligned} \frac{\partial}{\partial \Lambda} \log \hat{p}(\Lambda | \mathbf{y}) &= E_{\mathbf{u}} \left[\Lambda \sum_{k=1}^K (\mathbf{y}_k - \Lambda \mathbf{u}_k) \mathbf{u}_k^T \right] - \Lambda \Lambda \boldsymbol{\alpha} \\ &= \Lambda \mathbf{R}_{yu} - \Lambda \Lambda (\mathbf{R}_{uu} + \boldsymbol{\alpha}). \end{aligned} \quad (5.67)$$

Setting the right-hand side of the equation above to zero, we get

$$\bar{\Lambda} = \mathbf{R}_{yu} (\mathbf{R}_{uu} + \boldsymbol{\alpha})^{-1}. \quad (5.68)$$

where \mathbf{R}_{uu} and \mathbf{R}_{yu} are defined in Eqs. (5.15) and (5.16). The precision $\lambda_j \Psi$ is obtained as the coefficient of \mathbf{a}_j in the right-hand side of Eq. (5.67). The second term in the right-hand side of this equation can be rewritten as

$$\Lambda \Lambda (\mathbf{R}_{uu} + \boldsymbol{\alpha}) = \begin{bmatrix} \lambda_1 \mathbf{a}_1^T \\ \vdots \\ \lambda_M \mathbf{a}_M^T \end{bmatrix} (\mathbf{R}_{uu} + \boldsymbol{\alpha}) = \begin{bmatrix} \lambda_1 \mathbf{a}_1^T (\mathbf{R}_{uu} + \boldsymbol{\alpha}) \\ \vdots \\ \lambda_M \mathbf{a}_M^T (\mathbf{R}_{uu} + \boldsymbol{\alpha}) \end{bmatrix}. \quad (5.69)$$

Thus, Ψ is obtained as

$$\Psi = \mathbf{R}_{uu} + \boldsymbol{\alpha}. \quad (5.70)$$

Equations (5.68) and (5.70) are the M-step update equations in the VBFA algorithm. However, to compute these equations, we need to know the hyperparameters $\boldsymbol{\alpha}$ and Λ . The next subsection deals with the estimation of $\boldsymbol{\alpha}$.

5.3.2.3 Update Equation for Hyperparameter $\boldsymbol{\alpha}$

The hyperparameter $\boldsymbol{\alpha}$ is estimated by maximizing the free energy.³ According to the arguments in Sect. B.6, the free energy is expressed as

³ An estimate that maximizes the free energy is the MAP estimate under the assumption of the non-informative prior for the hyperparameter.

$$\begin{aligned}\mathcal{F}[\boldsymbol{\alpha}, \boldsymbol{\Lambda}] &= E_{(\boldsymbol{A}, \boldsymbol{u})} [\log p(\boldsymbol{u}, \boldsymbol{y}, \boldsymbol{A}) - \log \hat{p}(\boldsymbol{u}|\boldsymbol{y}) - \log \hat{p}(\boldsymbol{A}|\boldsymbol{y})] \\ &= E_{(\boldsymbol{A}, \boldsymbol{u})} [\log p(\boldsymbol{y}|\boldsymbol{u}, \boldsymbol{A}) + \log p(\boldsymbol{u}) \\ &\quad + \log p(\boldsymbol{A}) - \log \hat{p}(\boldsymbol{u}|\boldsymbol{y}) - \log \hat{p}(\boldsymbol{A}|\boldsymbol{y})],\end{aligned}\quad (5.71)$$

where $E_{(\boldsymbol{A}, \boldsymbol{u})}[\cdot]$ indicates the average regarding both $\hat{p}(\boldsymbol{A}|\boldsymbol{y})$ and $\hat{p}(\boldsymbol{u}|\boldsymbol{y})$. On the right-most-side of the equation above, the terms containing $\boldsymbol{\alpha}$ are only $\log p(\boldsymbol{A})$. Thus, omitting the terms not containing $\boldsymbol{\alpha}$, we can rewrite the free energy as

$$\mathcal{F}[\boldsymbol{\alpha}, \boldsymbol{\Lambda}] = E_{(\boldsymbol{A}, \boldsymbol{u})} [\log p(\boldsymbol{A})]. \quad (5.72)$$

Considering

$$\begin{aligned}\log p(\boldsymbol{A}) &= \log \prod_{j=1}^M \mathcal{N}(\boldsymbol{a}_j | \mathbf{0}, \lambda_j \boldsymbol{\alpha}) = \frac{1}{2} \sum_{j=1}^M \log |\lambda_j \boldsymbol{\alpha}| - \frac{1}{2} \sum_{j=1}^M \boldsymbol{a}_j^T \lambda_j \boldsymbol{\alpha} \boldsymbol{a}_j \\ &= \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \log(\lambda_j \alpha_\ell) - \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell A_{j,\ell}^2,\end{aligned}\quad (5.73)$$

The free energy in Eq. (5.72) is rewritten as

$$\mathcal{F}[\boldsymbol{\alpha}, \boldsymbol{\Lambda}] = E_{(\boldsymbol{A}, \boldsymbol{u})} \left[\frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \log(\lambda_j \alpha_\ell) - \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell A_{j,\ell}^2 \right]. \quad (5.74)$$

To derive the estimate of $\boldsymbol{\alpha}$, we differentiate the free energy in (5.74) with respect to $\boldsymbol{\alpha}$, resulting in

$$\frac{\partial}{\partial \boldsymbol{\alpha}} \mathcal{F}[\boldsymbol{\alpha}, \boldsymbol{\Lambda}] = E_{\boldsymbol{A}} \left[\frac{\partial}{\partial \boldsymbol{\alpha}} \left[\frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \log(\lambda_j \alpha_\ell) - \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell A_{j,\ell}^2 \right] \right], \quad (5.75)$$

where the expectation $E_{\boldsymbol{u}}[\cdot]$ is omitted because the terms on the right-hand side of the equation above do not contain \boldsymbol{u} . To compute the equation above, we consider the relationship

$$\frac{\partial}{\partial \alpha_\ell} \left[\frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \log(\lambda_j \alpha_\ell) - \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell A_{j,\ell}^2 \right] = \frac{M}{2} \alpha_\ell^{-1} - \frac{1}{2} \sum_{j=1}^M \lambda_j A_{j,\ell}^2, \quad (5.76)$$

and obtain

$$\begin{aligned}
& \frac{\partial}{\partial \boldsymbol{\alpha}} \left[\frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \log(\lambda_j \alpha_\ell) - \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell A_{j,\ell}^2 \right] \\
&= \frac{M}{2} \boldsymbol{\alpha}^{-1} - \frac{1}{2} \begin{bmatrix} \sum_{j=1}^M \lambda_j A_{j,1}^2 & 0 & \dots & 0 \\ 0 & \sum_{j=1}^M \lambda_j A_{j,2}^2 & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & \sum_{j=1}^M \lambda_j A_{j,L}^2 \end{bmatrix} \\
&= \frac{1}{2} (M \boldsymbol{\alpha}^{-1} - \text{diag}[\mathbf{A} \boldsymbol{\Lambda} \mathbf{A}^T]), \tag{5.77}
\end{aligned}$$

where $\text{diag}[\cdot]$ indicates a diagonal matrix whose diagonal entries are equal to those of a matrix in the brackets.⁴ The relationship

$$E_{\mathbf{A}} \left[\sum_{j=1}^M \lambda_j A_{j,\ell}^2 \right] = \sum_{j=1}^M \lambda_j \bar{A}_{j,\ell}^2 + \sum_{j=1}^M \lambda_j \frac{1}{\lambda_j} [\boldsymbol{\Psi}^{-1}]_{\ell,\ell} = \sum_{j=1}^M \lambda_j \bar{A}_{j,\ell}^2 + M [\boldsymbol{\Psi}^{-1}]_{\ell,\ell}$$

also holds. (Note that Eq. (5.168) shows the general case of computing $E_{\mathbf{A}}[A_{i,k} A_{j,\ell}]$.) Thus, we finally obtain

$$\begin{aligned}
\frac{\partial}{\partial \boldsymbol{\alpha}} \mathcal{F}[\boldsymbol{\alpha}, \boldsymbol{\Lambda}] &= E_{\mathbf{A}} \left[\frac{\partial}{\partial \boldsymbol{\alpha}} \left[\frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \log(\lambda_j \alpha_\ell) - \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell A_{j,\ell}^2 \right] \right] \\
&= \frac{1}{2} M \boldsymbol{\alpha}^{-1} - \frac{1}{2} \begin{bmatrix} \sum_{j=1}^M \lambda_j \bar{A}_{j,1}^2 & 0 & \dots & 0 \\ 0 & \sum_{j=1}^M \lambda_j \bar{A}_{j,2}^2 & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & \sum_{j=1}^M \lambda_j \bar{A}_{j,L}^2 \end{bmatrix} \\
&\quad - \frac{1}{2} M \begin{bmatrix} [\boldsymbol{\Psi}^{-1}]_{1,1} & 0 & \dots & 0 \\ 0 & [\boldsymbol{\Psi}^{-1}]_{2,2} & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & [\boldsymbol{\Psi}^{-1}]_{L,L} \end{bmatrix} \\
&= \frac{1}{2} (M \boldsymbol{\alpha}^{-1} - \text{diag}[\bar{\mathbf{A}}^T \boldsymbol{\Lambda} \bar{\mathbf{A}}] - M \text{diag}[\boldsymbol{\Psi}^{-1}]). \tag{5.78}
\end{aligned}$$

⁴ Here computing the derivative of a scalar X with a diagonal matrix \mathbf{A} is equal to creating a diagonal matrix whose (j, j) th diagonal element is equal to $\partial X / \partial A_{j,j}$ where $A_{j,j}$ is the (j, j) th diagonal element of \mathbf{A} .

Therefore, setting the right-hand side of the equation above to zero, we have the update equation for α , such that,

$$\alpha^{-1} = \text{diag} \left[\frac{1}{M} \bar{\mathbf{A}}^T \Lambda \bar{\mathbf{A}} + \Psi^{-1} \right]. \quad (5.79)$$

5.3.2.4 Update Equation for the Noise Precision Λ

We next derive the update equation for Λ . To do so, we maximize the free energy in Eq. (5.71) with respect to Λ . On the right-hand-side of Eq. (5.71), the terms containing Λ are $\log p(\mathbf{y}|\mathbf{u}, \mathbf{A})$ and $\log p(\mathbf{A})$. Thus, omitting the terms not containing Λ , the free energy is expressed as

$$\begin{aligned} \mathcal{F}[\alpha, \Lambda] &= E_{(\mathbf{A}, \mathbf{u})} [\log p(\mathbf{y}|\mathbf{u}, \mathbf{A}) + \log p(\mathbf{A})] \\ &= E_{(\mathbf{A}, \mathbf{u})} \left[\frac{K}{2} \log |\Lambda| - \frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k)^T \Lambda (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k) \right. \\ &\quad \left. + \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \log(\lambda_j \alpha_\ell) - \frac{1}{2} \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell A_{j,\ell}^2 \right]. \end{aligned} \quad (5.80)$$

To maximize $\mathcal{F}[\alpha, \Lambda]$ with respect to Λ , we consider the relationship

$$\begin{aligned} \frac{1}{\partial \lambda_j} \sum_{j=1}^M \sum_{\ell=1}^L \left[\frac{1}{2} \log(\lambda_j \alpha_\ell) - \frac{1}{2} \lambda_j \alpha_\ell A_{j,\ell}^2 \right] \\ = \sum_{\ell=1}^L \left[\frac{1}{2} \frac{1}{\lambda_j} - \frac{1}{2} \alpha_\ell A_{j,\ell}^2 \right] = \frac{L}{2} \frac{1}{\lambda_j} - \frac{1}{2} \sum_{\ell=1}^L \alpha_\ell A_{j,\ell}^2, \end{aligned}$$

and obtain

$$\begin{aligned} \frac{1}{\partial \Lambda} \sum_{j=1}^M \sum_{\ell=1}^L \left[\frac{1}{2} \log(\lambda_j \alpha_\ell) - \frac{1}{2} \lambda_j \alpha_\ell A_{j,\ell}^2 \right] \\ = \frac{L}{2} \begin{bmatrix} 1/\lambda_1 & 0 & \dots & 0 \\ 0 & 1/\lambda_2 & \dots & 0 \\ 0 & \vdots & \ddots & 0 \\ 0 & 0 & \dots & 1/\lambda_M \end{bmatrix} \\ - \frac{1}{2} \begin{bmatrix} \sum_{\ell=1}^L \alpha_\ell A_{1,\ell}^2 & 0 & \dots & 0 \\ 0 & \sum_{\ell=1}^L \alpha_\ell A_{2,\ell}^2 & \dots & 0 \\ 0 & \vdots & \ddots & 0 \\ 0 & 0 & \dots & \sum_{\ell=1}^L \alpha_\ell A_{M,\ell}^2 \end{bmatrix} \end{aligned}$$

$$= \frac{L}{2} \boldsymbol{\Lambda}^{-1} - \frac{1}{2} \text{diag}[\boldsymbol{A}\boldsymbol{\alpha}\boldsymbol{A}^T]. \quad (5.81)$$

Then, using Eq. (5.80), we derive the derivative of $\mathcal{F}(\boldsymbol{\Lambda}, \boldsymbol{\alpha})$ with respect to $\boldsymbol{\Lambda}$ such that

$$\begin{aligned} & \frac{\partial}{\partial \boldsymbol{\Lambda}} \mathcal{F}(\boldsymbol{\Lambda}, \boldsymbol{\alpha}) \\ &= E_{(\boldsymbol{A}, \boldsymbol{u})} \left[\frac{K}{2} \boldsymbol{\Lambda}^{-1} - \frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \boldsymbol{A}\mathbf{u}_k)(\mathbf{y}_k - \boldsymbol{A}\mathbf{u}_k)^T + \frac{L}{2} \boldsymbol{\Lambda}^{-1} - \frac{1}{2} \boldsymbol{A}\boldsymbol{\alpha}\boldsymbol{A}^T \right] \\ &= \frac{K}{2} \boldsymbol{\Lambda}^{-1} + \frac{1}{2} E_{\boldsymbol{A}}[-\boldsymbol{R}_{yy} + \boldsymbol{R}_{yu}\boldsymbol{A}^T + \boldsymbol{A}\boldsymbol{R}_{uy} - \boldsymbol{A}\boldsymbol{R}_{uu}\boldsymbol{A}^T + L\boldsymbol{\Lambda}^{-1} - \boldsymbol{A}\boldsymbol{\alpha}\boldsymbol{A}^T] \\ &= \frac{K}{2} \boldsymbol{\Lambda}^{-1} - \frac{1}{2} (\boldsymbol{R}_{yy} - \boldsymbol{R}_{yu}\bar{\boldsymbol{A}}^T - \bar{\boldsymbol{A}}\boldsymbol{R}_{uy}) + \frac{L}{2} \boldsymbol{\Lambda}^{-1} - \frac{1}{2} E_{\boldsymbol{A}}[\boldsymbol{A}\boldsymbol{R}_{uu}\boldsymbol{A} + \boldsymbol{A}\boldsymbol{\alpha}\boldsymbol{A}^T]. \end{aligned} \quad (5.82)$$

The following relationship holds:

$$E_{\boldsymbol{A}}[\boldsymbol{A}\boldsymbol{R}_{uu}\boldsymbol{A} + \boldsymbol{A}\boldsymbol{\alpha}\boldsymbol{A}^T] = E_{\boldsymbol{A}}[\boldsymbol{A}(\boldsymbol{R}_{uu} + \boldsymbol{\alpha})\boldsymbol{A}^T] = E_{\boldsymbol{A}}[\boldsymbol{A}\boldsymbol{\Psi}\boldsymbol{A}^T] \quad (5.83)$$

Also, we can prove the relationship

$$E_{\boldsymbol{A}}[\boldsymbol{A}\boldsymbol{\Psi}\boldsymbol{A}^T] = \bar{\boldsymbol{A}}\boldsymbol{\Psi}\bar{\boldsymbol{A}}^T + L\boldsymbol{\Lambda}^{-1} = \boldsymbol{R}_{yu}\bar{\boldsymbol{A}}^T + L\boldsymbol{\Lambda}^{-1}. \quad (5.84)$$

The proof of this equation is provided in Sect. 5.7.1. Thus, the relationship

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\Lambda}} \mathcal{F}(\boldsymbol{\Lambda}, \boldsymbol{\alpha}) &= \frac{K}{2} \boldsymbol{\Lambda}^{-1} - \frac{1}{2} (\boldsymbol{R}_{yy} - \boldsymbol{R}_{yu}\bar{\boldsymbol{A}}^T - \bar{\boldsymbol{A}}\boldsymbol{R}_{uy}) \\ &\quad + \frac{L}{2} \boldsymbol{\Lambda}^{-1} - \frac{1}{2} (\boldsymbol{R}_{yu}\bar{\boldsymbol{A}}^T + L\boldsymbol{\Lambda}^{-1}) = \frac{K}{2} \boldsymbol{\Lambda}^{-1} - \frac{1}{2} (\boldsymbol{R}_{yy} - \bar{\boldsymbol{A}}\boldsymbol{R}_{uy}) \end{aligned} \quad (5.85)$$

holds. Setting the right-hand side of the equation above zero, we obtain the update equation for $\boldsymbol{\Lambda}$ as

$$\boldsymbol{\Lambda}^{-1} = \frac{1}{K} \text{diag}(\boldsymbol{R}_{yy} - \bar{\boldsymbol{A}}\boldsymbol{R}_{uy}). \quad (5.86)$$

5.3.3 Computation of Free Energy

Although we cannot compute the exact marginal likelihood in the VBEM algorithm, we can compute its lower bound using the free energy. The free energy after the E-step of the VBEM algorithm is expressed as

$$\begin{aligned}\mathcal{F} &= E_{(A, \mathbf{u})} [\log p(\mathbf{y}, \mathbf{u}, A) - \log \hat{p}(\mathbf{u}|\mathbf{y}) - \log \hat{p}(A|\mathbf{y})] \\ &= E_{(A, \mathbf{u})} [\log p(\mathbf{y}|\mathbf{u}, A) + \log p(\mathbf{u}) + \log p(A)] \\ &\quad + \mathcal{H}(\hat{p}(\mathbf{u}|\mathbf{y})) + \mathcal{H}(\hat{p}(A|\mathbf{y})).\end{aligned}\tag{5.87}$$

Substituting

$$\log p(\mathbf{y}|\mathbf{u}, A) = \frac{K}{2} \log |\Lambda| - \frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - A\mathbf{u}_k)^T \Lambda (\mathbf{y}_k - A\mathbf{u}_k),\tag{5.88}$$

$$\log p(\mathbf{u}) = -\frac{1}{2} \sum_{k=1}^K \mathbf{u}_k \mathbf{u}_k^T,\tag{5.89}$$

$$\log p(A) = \frac{1}{2} \sum_{j=1}^M \log |\lambda_j \boldsymbol{\alpha}| - \frac{1}{2} \sum_{j=1}^M \mathbf{a}_j^T \lambda_j \boldsymbol{\alpha} \mathbf{a}_j,\tag{5.90}$$

and

$$\mathcal{H}(\hat{p}(\mathbf{u}|\mathbf{y})) = -\frac{K}{2} \log |\boldsymbol{\Gamma}|,\tag{5.91}$$

$$\mathcal{H}(\hat{p}(A|\mathbf{y})) = -\frac{1}{2} \sum_{j=1}^M \log |\lambda_j \boldsymbol{\Psi}|,\tag{5.92}$$

into Eq. (5.87), we obtain

$$\begin{aligned}\mathcal{F} &= \frac{K}{2} \log |\Lambda| + E_{(A, \mathbf{u})} \left[-\frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - A\mathbf{u}_k)^T \Lambda (\mathbf{y}_k - A\mathbf{u}_k) - \frac{1}{2} \sum_{k=1}^K \mathbf{u}_k \mathbf{u}_k^T \right] \\ &\quad + \frac{1}{2} \sum_{j=1}^M \log |\lambda_j \boldsymbol{\alpha}| - \frac{1}{2} E_A \left[\sum_{j=1}^M \mathbf{a}_j^T \lambda_j \boldsymbol{\alpha} \mathbf{a}_j \right] - \frac{K}{2} \log |\boldsymbol{\Gamma}| - \frac{1}{2} \sum_{j=1}^M \log |\lambda_j \boldsymbol{\Psi}|.\end{aligned}\tag{5.93}$$

Note that the following relationships:

$$\begin{aligned}E_{(A, \mathbf{u})} &\left[-\frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - A\mathbf{u}_k)^T \Lambda (\mathbf{y}_k - A\mathbf{u}_k) - \frac{1}{2} \sum_{k=1}^K \mathbf{u}_k \mathbf{u}_k^T \right] \\ &= -\frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \Lambda \mathbf{y}_k + \frac{1}{2} \sum_{k=1}^K \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{u}}_k\end{aligned}\tag{5.94}$$

and

$$\sum_{j=1}^M \log |\lambda_j \boldsymbol{\alpha}| - \sum_{j=1}^M \log |\lambda_j \boldsymbol{\Psi}| = \sum_{j=1}^M \log \frac{\lambda_j^M |\boldsymbol{\alpha}|}{\lambda_j^M |\boldsymbol{\Psi}|} = M \log \frac{|\boldsymbol{\alpha}|}{|\boldsymbol{\Psi}|} \quad (5.95)$$

hold. The proof of Eq. (5.94) is presented in Sect. 5.7.2. Also, we can show

$$\begin{aligned} E_A \left[\sum_{j=1}^M \bar{\mathbf{a}}_j^T \lambda_j \boldsymbol{\alpha} \bar{\mathbf{a}}_j \right] &= E_A \left[\sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell A_{j,\ell}^2 \right] \\ &= \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell \bar{A}_{j,\ell}^2 + \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell \frac{1}{\lambda_j} [\boldsymbol{\Psi}^{-1}]_{\ell,\ell} \\ &= \sum_{j=1}^M \sum_{\ell=1}^L \lambda_j \alpha_\ell \bar{A}_{j,\ell}^2 + M \sum_{\ell=1}^L \alpha_\ell [\boldsymbol{\Psi}^{-1}]_{\ell,\ell} \\ &= \sum_{j=1}^M \bar{\mathbf{a}}_j^T \lambda_j \boldsymbol{\alpha} \bar{\mathbf{a}}_j + M \operatorname{tr}(\boldsymbol{\alpha} \boldsymbol{\Psi}^{-1}). \end{aligned} \quad (5.96)$$

By substituting the above three equations into Eq. (5.93), we obtain the expression for the free energy:

$$\begin{aligned} \mathcal{F} &= \frac{K}{2} \log \frac{|\boldsymbol{\Lambda}|}{|\boldsymbol{\Gamma}|} - \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{y}_k + \frac{1}{2} \sum_{k=1}^K \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{u}}_k \\ &\quad + \frac{M}{2} \log \frac{|\boldsymbol{\alpha}|}{|\boldsymbol{\Psi}|} - \frac{1}{2} \sum_{j=1}^M \bar{\mathbf{a}}_j^T \lambda_j \boldsymbol{\alpha} \bar{\mathbf{a}}_j - \frac{M}{2} \operatorname{tr}(\boldsymbol{\alpha} \boldsymbol{\Psi}^{-1}). \end{aligned} \quad (5.97)$$

Note here that

$$\sum_{j=1}^M \bar{\mathbf{a}}_j^T \lambda_j \boldsymbol{\alpha} \bar{\mathbf{a}}_j = \operatorname{tr}[\bar{\mathbf{A}} \boldsymbol{\alpha} \bar{\mathbf{A}}^T \boldsymbol{\Lambda}] = \operatorname{tr}[\boldsymbol{\alpha} \bar{\mathbf{A}}^T \boldsymbol{\Lambda} \bar{\mathbf{A}}]. \quad (5.98)$$

Thus, we have

$$\begin{aligned} \sum_{j=1}^M \bar{\mathbf{a}}_j^T \lambda_j \boldsymbol{\alpha} \bar{\mathbf{a}}_j + M \operatorname{tr}(\boldsymbol{\alpha} \boldsymbol{\Psi}^{-1}) &= \operatorname{tr}[\boldsymbol{\alpha} (\bar{\mathbf{A}}^T \boldsymbol{\Lambda} \bar{\mathbf{A}} + M \boldsymbol{\Psi}^{-1})] \\ &= \operatorname{tr}[\boldsymbol{\alpha} M \boldsymbol{\alpha}^{-1}] = ML. \end{aligned} \quad (5.99)$$

In the equation above, we use the update equation for α in Eq.(5.79). Therefore, ignoring constant terms, the free energy is finally expressed as

$$\mathcal{F} = \frac{K}{2} \log \frac{|\Lambda|}{|\Gamma|} - \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \Lambda \mathbf{y}_k + \frac{1}{2} \sum_{k=1}^K \bar{\mathbf{u}}_k^T \Gamma \bar{\mathbf{u}}_k + \frac{M}{2} \log \frac{|\alpha|}{|\Psi|}. \quad (5.100)$$

Since the free energy is limited (i.e., upper bounded) by the likelihood $\log p(\mathbf{y}|\theta)$ (where θ collectively expresses the hyperparameters), increasing the free energy increases the likelihood.

5.3.4 Summary of the VBFA Algorithm

Let us summarize the VBFA algorithm. The algorithm is based on the factor analysis model in Eq.(5.2). The algorithm estimates the factor activity \mathbf{u}_k , the mixing matrix Λ , and the sensor-noise precision matrix Λ . In the estimation, an overspecified value is set for the model order L , and appropriate initial values are set for $\bar{\Lambda}$, Ψ , and Λ . In the E-step, Γ and $\bar{\mathbf{u}}_k$ ($k = 1, \dots, K$) are updated according to Eqs.(5.58) and (5.59). In the M-step, values of $\bar{\Lambda}$, and Ψ are updated using Eqs.(5.68), and (5.70).

The hyperparameters, α and Λ are updated using Eqs.(5.79) and (5.86). Since we cannot compute the marginal likelihood, the free energy in Eq.(5.100)—which is the lower limit of the marginal likelihood—is used for monitoring the progress of the VBEM iteration. That is, if the increase of \mathcal{F} in Eq.(5.100) becomes very small with respect to the iteration count, the VBEM iteration can be terminated.

Using the VBFA algorithm, the estimate of the signal component in the sensor data, $\hat{\mathbf{y}}_k^S$, is given by

$$\hat{\mathbf{y}}_k^S = E_{(A, u)}[Au] = EA[A]E_u[u] = \bar{\Lambda}\bar{\mathbf{u}}_k. \quad (5.101)$$

In the equation above, we use the values of $\bar{\mathbf{u}}_k$ and $\bar{\Lambda}$ obtained when the VBEM iteration is terminated. The sample covariance matrix can be computed using only the signal component, and such a covariance matrix is derived as

$$\begin{aligned} \bar{\mathbf{R}} &= \frac{1}{K} \sum_{k=1}^K E_{(u, A)} \left[(Au_k)(Au_k)^T \right] = \frac{1}{K} \sum_{k=1}^K E_u E_A \left[[Au_k \mathbf{u}_k^T A^T] \right] \\ &= \frac{1}{K} \sum_{k=1}^K E_A \left[AE_u \left[\mathbf{u}_k \mathbf{u}_k^T \right] A^T \right] = \frac{1}{K} E_A \left[A \mathbf{R}_{uu} A^T \right], \end{aligned} \quad (5.102)$$

where \mathbf{R}_{uu} is defined in Eq.(5.15). We can show

$$E_A \left[A \mathbf{R}_{uu} A^T \right] = \bar{\Lambda} \mathbf{R}_{uu} \bar{\Lambda}^T + \Lambda^{-1} \text{tr} \left(\mathbf{R}_{uu} \Psi^{-1} \right), \quad (5.103)$$

where the proof is presented in Sect. 5.7.3. Using Eq. (5.103), we have

$$\bar{\mathbf{R}} = \frac{1}{K} \bar{\mathbf{A}} \mathbf{R}_{uu} \bar{\mathbf{A}}^T + \frac{1}{K} \mathbf{A}^{-1} \operatorname{tr}(\mathbf{R}_{uu} \Psi^{-1}). \quad (5.104)$$

This $\bar{\mathbf{R}}$ can be used in source imaging algorithms such as the adaptive beamformers.

5.4 Partitioned Factor Analysis (PFA)

5.4.1 Factor Analysis Model

The bioelectromagnetic data is often contaminated not only by sensor noise but also by various types of interference of biological and nonbiological origins. A simple modification of the VBFA algorithm enables the removal of such interferences. The modified algorithm is called the partitioned factor analysis (PFA) [1]. The prerequisite for the PFA algorithm is that a control measurement, which contains the interferences but does not contain the signal of interest, be available. The factor analysis model for PFA is expressed as

$$\mathbf{y}_k = \mathbf{B}\mathbf{v}_k + \boldsymbol{\varepsilon} \quad \text{for control data,} \quad (5.105)$$

$$\mathbf{y}_k = \mathbf{A}\mathbf{u}_k + \mathbf{B}\mathbf{v}_k + \boldsymbol{\varepsilon} \quad \text{for target data.} \quad (5.106)$$

In the equations above, $L \times 1$ column vector \mathbf{u}_k is the factor activity that represents the signal of interest and \mathbf{A} is an $M \times L$ mixing matrix. Also, $L_v \times 1$ column vector \mathbf{v}_k is the factor activity that represents the interference and \mathbf{B} is an $M \times L_v$ interference mixing matrix.

The PFA algorithm has a two-step procedure. The first step applies the VBFA algorithm to the control data and estimates the mixing matrix \mathbf{B} and the sensor-noise precision \mathbf{A} . The second step estimates the mixing matrix \mathbf{A} and the signal factor activity \mathbf{u}_k . In the second step, the interference mixing matrix \mathbf{B} and the noise precision \mathbf{A} are fixed, and the VBFA algorithm is applied to the target data, which is expressed as

$$\mathbf{y}_k = \mathbf{A}\mathbf{u}_k + \mathbf{B}\mathbf{v}_k + \boldsymbol{\varepsilon} = [\mathbf{A}, \mathbf{B}] \begin{bmatrix} \mathbf{u}_k \\ \mathbf{v}_k \end{bmatrix} + \boldsymbol{\varepsilon} = \mathbf{A}_c \mathbf{z}_k + \boldsymbol{\varepsilon}, \quad (5.107)$$

where

$$\mathbf{A}_c = [\mathbf{A}, \mathbf{B}] \quad \text{and} \quad \mathbf{z}_k = \begin{bmatrix} \mathbf{u}_k \\ \mathbf{v}_k \end{bmatrix}.$$

The equation above indicates that the second step can also be expressed by the factor analysis model using the augmented factor vector \mathbf{z}_k and the augmented mixing matrix \mathbf{A}_c .

5.4.2 Probability Model

The prior probability distributions for the factor vectors are:

$$P(\mathbf{u}_k) = \mathcal{N}(\mathbf{u}_k | \mathbf{0}, \mathbf{I}), \quad (5.108)$$

$$P(\mathbf{v}_k) = \mathcal{N}(\mathbf{v}_k | \mathbf{0}, \mathbf{I}). \quad (5.109)$$

Thus, we have the prior probability distribution for the augmented factor vector \mathbf{z}_k , such that

$$\begin{aligned} p(\mathbf{z}_k) &= p(\mathbf{u}_k)p(\mathbf{v}_k) = \mathcal{N}(\mathbf{u}_k | \mathbf{0}, \mathbf{I})\mathcal{N}(\mathbf{v}_k | \mathbf{0}, \mathbf{I}) \\ &= \left| \frac{\mathbf{I}}{2\pi} \right|^{1/2} \exp \left[-\frac{1}{2} \mathbf{u}_k^T \mathbf{u}_k \right] \left| \frac{\mathbf{I}}{2\pi} \right|^{1/2} \exp \left[-\frac{1}{2} \mathbf{v}_k^T \mathbf{v}_k \right] \\ &= \left| \frac{\mathbf{I}}{2\pi} \right|^{1/2} \exp \left[-\frac{1}{2} \mathbf{z}_k^T \mathbf{z}_k \right] = \mathcal{N}(\mathbf{z}_k | \mathbf{0}, \mathbf{I}). \end{aligned} \quad (5.110)$$

The noise is assumed to be Gaussian,

$$P(\boldsymbol{\varepsilon}) = \mathcal{N}(\boldsymbol{\varepsilon} | \mathbf{0}, \boldsymbol{\Lambda}^{-1}).$$

Therefore, we have the conditional probability, $P(\mathbf{y}_k | \mathbf{u}_k, \mathbf{v}_k)$, such that

$$P(\mathbf{y}_k | \mathbf{u}_k, \mathbf{v}_k) = P(\mathbf{y}_k | \mathbf{z}_k) = \mathcal{N}(\mathbf{y}_k | \mathbf{A}\mathbf{u}_k + \mathbf{B}\mathbf{v}_k, \boldsymbol{\Lambda}^{-1}) = \mathcal{N}(\mathbf{y}_k | \mathbf{A}_c \mathbf{z}_k, \boldsymbol{\Lambda}^{-1}). \quad (5.111)$$

We assume the same prior for \mathbf{A} as in Eq. (5.42).

5.4.3 VBEM Algorithm for PFA

5.4.3.1 E-Step

The E-step estimates the posterior probability $\hat{p}(\mathbf{z}_k | \mathbf{y}_k)$. Using the same arguments in Sect. 5.3.2.1, we have,

$$\begin{aligned} \log \hat{p}(\mathbf{z}_k | \mathbf{y}_k) &= E_{\mathbf{A}} [\log p(\mathbf{z}_k, \mathbf{y}_k, \mathbf{A}_c)] \\ &= E_{\mathbf{A}} [\log p(\mathbf{y}_k | \mathbf{z}_k, \mathbf{A}_c) + \log p(\mathbf{z}_k) + \log p(\mathbf{A}_c)]. \end{aligned} \quad (5.112)$$

Note that, since \mathbf{B} is fixed, $p(\mathbf{A}_c)$ is the same as $p(\mathbf{A})$, which is expressed in Eq. (5.42). Omitting the terms not containing \mathbf{z}_k , and using Eqs. (5.110) and (5.111), we obtain

$$\log \hat{p}(\mathbf{z}_k | \mathbf{y}_k) = E_{\mathbf{A}} \left[-\frac{1}{2} (\mathbf{y}_k - \mathbf{A}_c \mathbf{z}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{A}_c \mathbf{z}_k) \right] - \frac{1}{2} \mathbf{z}_k^T \mathbf{z}_k. \quad (5.113)$$

Since $\hat{p}(z_k|y_k)$ is Gaussian, we assume that

$$\hat{p}(z_k|y_k) = \mathcal{N}(z_k|\bar{z}_k, \boldsymbol{\Gamma}^{-1}).$$

The mean \bar{z}_k is obtained as z_k that makes $\frac{\partial}{\partial z_k} \log \hat{p}(z_k|y_k)$ equal to zero, and the precision $\boldsymbol{\Gamma}$ is obtained from the coefficient of z_k in this derivative. Using Eq. (5.113), the derivative is given by

$$\begin{aligned} \frac{\partial}{\partial z_k} \log \hat{p}(z_k|y_k) &= E_A[A_c^T \boldsymbol{\Lambda} (y_k - A_c z_k)] - z_k \\ &= \bar{A}_c^T \boldsymbol{\Lambda} y_k - E_A[A_c^T \boldsymbol{\Lambda} A_c] z_k - z_k. \end{aligned} \quad (5.114)$$

In the expression above, using Eq. (5.56), $E_A[A_c^T \boldsymbol{\Lambda} A_c]$ is expressed as

$$\begin{aligned} E_A[A_c^T \boldsymbol{\Lambda} A_c] &= E_A \left[\begin{bmatrix} \mathbf{A}^T \\ \mathbf{B}^T \end{bmatrix} \boldsymbol{\Lambda} [\mathbf{A}^T \mathbf{B}^T] \right] = \begin{bmatrix} E_A[\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A}] & \bar{A}_c^T \boldsymbol{\Lambda} \mathbf{B} \\ \mathbf{B}^T \boldsymbol{\Lambda} \bar{A} & \mathbf{B}^T \boldsymbol{\Lambda} \mathbf{B} \end{bmatrix} \\ &= \begin{bmatrix} \bar{A}_c^T \boldsymbol{\Lambda} \bar{A} + M \boldsymbol{\Psi}^{-1} & \bar{A}_c^T \boldsymbol{\Lambda} \mathbf{B} \\ \mathbf{B}^T \boldsymbol{\Lambda} \bar{A} & \mathbf{B}^T \boldsymbol{\Lambda} \mathbf{B} \end{bmatrix} = \bar{A}_c^T \boldsymbol{\Lambda} \bar{A}_c + M \boldsymbol{\Psi}_c^{-1}, \end{aligned} \quad (5.115)$$

where

$$\boldsymbol{\Psi}_c^{-1} = \begin{bmatrix} \boldsymbol{\Psi}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad (5.116)$$

and

$$\bar{A}_c = [\bar{A}, \mathbf{B}]. \quad (5.117)$$

Therefore, we finally have

$$\frac{\partial}{\partial z_k} \log \hat{p}(z_k|y_k) = \bar{A}_c^T \boldsymbol{\Lambda} y_k - (\bar{A}_c^T \boldsymbol{\Lambda} \bar{A}_c z_k + M \boldsymbol{\Psi}_c^{-1} + \mathbf{I}) z_k. \quad (5.118)$$

The precision matrix of the posterior distribution is obtained as

$$\boldsymbol{\Gamma} = \bar{A}_c^T \boldsymbol{\Lambda} \bar{A}_c + M \boldsymbol{\Psi}_c^{-1} + \mathbf{I}, \quad (5.119)$$

and the mean of the posterior as

$$\bar{z}_k = \begin{bmatrix} \bar{u}_k \\ \bar{v}_k \end{bmatrix} = \boldsymbol{\Gamma}^{-1} \begin{bmatrix} \bar{A}^T \\ \mathbf{B}^T \end{bmatrix} \boldsymbol{\Lambda} y_k. \quad (5.120)$$

Equations (5.119) and (5.120) are the E-step update equations in the PFA algorithm.

5.4.3.2 M-Step

The M-step estimates the posterior probability $\hat{p}(A_c | \mathbf{y}_k)$. Using the same arguments as in Sect. 5.3.2.2, we have

$$\begin{aligned} & \log \hat{p}(A_c | \mathbf{y}) \\ &= E_z [\log p(\mathbf{y}, \mathbf{z}, A_c)] = E_z [\log p(\mathbf{y} | \mathbf{z}, A_c) + \log p(\mathbf{z}) + \log p(A_c)] \\ &= E_z \left[\sum_{k=1}^K \log p(\mathbf{y}_k | \mathbf{z}_k, A_c) + \log p(A_c) \right] \\ &= E_z \left[-\frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k - \mathbf{B}\mathbf{v}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k - \mathbf{B}\mathbf{v}_k) - \frac{1}{2} \sum_{j=1}^M \mathbf{a}_j^T \lambda_j \boldsymbol{\alpha} \mathbf{a}_j \right], \end{aligned} \quad (5.121)$$

where \mathbf{z} collectively expresses $\mathbf{z}_1, \dots, \mathbf{z}_K$, and terms not containing \mathbf{A} are omitted. The form of the posterior distribution in Eq. (5.44) is also assumed:

$$p(A_c | \mathbf{y}) = p(\mathbf{A} | \mathbf{y}) = \prod_{j=1}^M \mathcal{N}(\mathbf{a}_j | \bar{\mathbf{a}}_j, (\lambda_j \boldsymbol{\Psi})^{-1}).$$

The mean $\bar{\mathbf{A}}$ is obtained as the \mathbf{A} that makes $\frac{\partial}{\partial \mathbf{A}} \log \hat{p}(\mathbf{A} | \mathbf{y})$ equal to zero, and the precision $\lambda_j \boldsymbol{\Psi}$ as the coefficient of \mathbf{a}_j in that derivative. Using Eqs. (5.65) and (5.66) the derivative is computed as

$$\begin{aligned} \frac{\partial}{\partial \mathbf{A}} \log \hat{p}(\mathbf{A} | \mathbf{y}_k) &= E_z \left[\boldsymbol{\Lambda} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k - \mathbf{B}\mathbf{v}_k) \mathbf{u}_k^T \right] - \boldsymbol{\Lambda} \mathbf{A} \boldsymbol{\alpha} \\ &= \boldsymbol{\Lambda} \mathbf{R}_{yu} - \boldsymbol{\Lambda} \mathbf{B} \mathbf{R}_{vu} - \boldsymbol{\Lambda} \mathbf{A} (\mathbf{R}_{uu} + \boldsymbol{\alpha}). \end{aligned} \quad (5.122)$$

In Eq. (5.122), the coefficient of \mathbf{a}_j is $\lambda_j (\mathbf{R}_{uu} + \boldsymbol{\alpha})$, and thus the matrix $\boldsymbol{\Psi}$ is equal to

$$\boldsymbol{\Psi} = \mathbf{R}_{uu} + \boldsymbol{\alpha}, \quad (5.123)$$

and $\bar{\mathbf{A}}$ is obtained as

$$\bar{\mathbf{A}} = (\mathbf{R}_{yu} - \mathbf{B} \mathbf{R}_{vu}) \boldsymbol{\Psi}^{-1}, \quad (5.124)$$

where \mathbf{R}_{yu} is obtained as

$$\mathbf{R}_{yu} = \sum_{k=1}^K \mathbf{y}_k \bar{\mathbf{u}}_k^T. \quad (5.125)$$

Also, \mathbf{R}_{vu} and \mathbf{R}_{uu} are obtained in the following manner. We first define

$$\boldsymbol{\Gamma}^{-1} = \begin{bmatrix} \tilde{\boldsymbol{\Gamma}}_{uu} & \tilde{\boldsymbol{\Gamma}}_{uv} \\ \tilde{\boldsymbol{\Gamma}}_{vu} & \tilde{\boldsymbol{\Gamma}}_{vv} \end{bmatrix}, \quad (5.126)$$

and have

$$\begin{aligned} \begin{bmatrix} \mathbf{R}_{uu} & \mathbf{R}_{uv} \\ \mathbf{R}_{vu} & \mathbf{R}_{vv} \end{bmatrix} &= E_z \left[\sum_{k=1}^K \begin{bmatrix} \mathbf{u}_k \\ \mathbf{v}_k \end{bmatrix} [\mathbf{u}_k^T, \mathbf{v}_k^T] \right] \\ &= \begin{bmatrix} \sum_{k=1}^K \bar{\mathbf{u}}_k \bar{\mathbf{u}}_k^T & \sum_{k=1}^K \bar{\mathbf{u}}_k \bar{\mathbf{v}}_k^T \\ \sum_{k=1}^K \bar{\mathbf{v}}_k \bar{\mathbf{u}}_k^T & \sum_{k=1}^K \bar{\mathbf{v}}_k \bar{\mathbf{v}}_k^T \end{bmatrix} + K \begin{bmatrix} \tilde{\boldsymbol{\Gamma}}_{uu} & \tilde{\boldsymbol{\Gamma}}_{uv} \\ \tilde{\boldsymbol{\Gamma}}_{vu} & \tilde{\boldsymbol{\Gamma}}_{vv} \end{bmatrix}. \end{aligned} \quad (5.127)$$

Thus, we obtain,

$$\mathbf{R}_{uu} = \sum_{k=1}^K \bar{\mathbf{u}}_k \bar{\mathbf{u}}_k^T + K \tilde{\boldsymbol{\Gamma}}_{uu}, \quad (5.128)$$

$$\mathbf{R}_{vu} = \sum_{k=1}^K \bar{\mathbf{v}}_k \bar{\mathbf{u}}_k^T + K \tilde{\boldsymbol{\Gamma}}_{vu}. \quad (5.129)$$

The hyperparameter α is obtained by maximizing the free energy, which is expressed as

$$\begin{aligned} \mathcal{F} &= E_{(\mathbf{A}, \mathbf{u})} [\log p(\mathbf{y}|\mathbf{z}, \mathbf{A}_c) + \log p(\mathbf{z}) + \log p(\mathbf{A}_c) \\ &\quad - \log \hat{p}(\mathbf{z}|\mathbf{y}) - \log \hat{p}(\mathbf{A}_c|\mathbf{y})]. \end{aligned} \quad (5.130)$$

However, since α is contained only in $\log p(\mathbf{A}_c)$, (which is equal to $\log p(\mathbf{A})$), the update equation for α is exactly the same as that in Eq. (5.79).

5.4.4 Summary of the PFA Algorithm

The PFA algorithm is summarized as follows. The first step estimates the interference mixing matrix \mathbf{B} and the diagonal noise precision matrix \mathbf{A} by applying the VBFA algorithm to the control data. The second step applies the PFA-VBEM algorithm to the target data, and estimates the signal factor vector \mathbf{u}_k and the signal mixing matrix \mathbf{A} . In the second step, \mathbf{B} and \mathbf{A} are fixed at the values obtained in the first step. There is a different version of the PFA algorithm in which \mathbf{B} and \mathbf{A} are also updated in the second step. The details of the algorithm are given in [1].

The free energy is computed in exactly the same manner as in Eq. (5.100) except that $\bar{\mathbf{u}}_k$ is replaced by $\bar{\mathbf{z}}_k$. Thus, the free energy is expressed as

$$\mathcal{F} = \frac{K}{2} \log \frac{|\boldsymbol{\Lambda}|}{|\boldsymbol{\Gamma}|} - \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{y}_k + \frac{1}{2} \sum_{k=1}^K \bar{\mathbf{z}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{z}}_k + \frac{M}{2} \log \frac{|\boldsymbol{\alpha}|}{|\boldsymbol{\Psi}|}. \quad (5.131)$$

This free energy is used for monitoring the progress of the PFA-VBEM iteration.

Using the PFA algorithm, the estimate of the signal of interest, $\hat{\mathbf{y}}_k^S$, is given by.

$$\hat{\mathbf{y}}_k^S = E_{(A, \mathbf{u})}[A\mathbf{u}_k] = E_A[A]E_{\mathbf{u}}[\mathbf{u}_k] = \bar{\mathbf{A}}\bar{\mathbf{u}}_k. \quad (5.132)$$

The sample covariance matrix computed only using the signal of interest is given by

$$\bar{\mathbf{R}} = \frac{1}{K} \bar{\mathbf{A}} \mathbf{R}_{uu} \bar{\mathbf{A}}^T + \frac{1}{K} \boldsymbol{\Lambda}^{-1} \text{tr}(\mathbf{R}_{uu} \boldsymbol{\Psi}^{-1}), \quad (5.133)$$

where \mathbf{R}_{uu} is defined in Eq. (5.128). This $\bar{\mathbf{R}}$ can be used in source imaging algorithms such as the adaptive beamformers, and resultant images can be free from the influence of interferences.

5.5 Saketini: Source Localization Algorithm Based on the VBFA Model

5.5.1 Data Model

This section describes a virtual-sensor type source localization algorithm, called Saketini [2]. The Saketini algorithm is based on Bayesian factor analysis, and enables the estimation of source activity. In the Saketini algorithm, the sensor data $\mathbf{y}(t)$ is modeled using

$$\mathbf{y}(t) = \mathbf{L}(\mathbf{r})\mathbf{s}(\mathbf{r}, t) + \mathbf{A}\mathbf{u}(t) + \boldsymbol{\varepsilon}, \quad (5.134)$$

where $\mathbf{s}(\mathbf{r}, t)$ is the source vector defined in Eq. (2.3). An $M \times 3$ matrix, $\mathbf{L}(\mathbf{r})$, is the lead field matrix at \mathbf{r} , which is defined in Eq. (2.4), and $\boldsymbol{\varepsilon}$ is the additive sensor noise. In Eq. (5.134), the signal from the source activity at \mathbf{r} is represented by $\mathbf{L}(\mathbf{r})\mathbf{s}(\mathbf{r}, t)$. On the basis of the factor analysis model, the interference is modeled using $\mathbf{A}\mathbf{u}(t)$ where \mathbf{A} is an $M \times L$ factor mixing matrix, and $\mathbf{u}(t)$ is an $L \times 1$ factor vector. Here, $\mathbf{A}\mathbf{u}(t)$ represents all interference and source activities except the source activity at \mathbf{r} .

As in the preceding sections, $\mathbf{y}(t_k)$, $\mathbf{s}(\mathbf{r}, t_k)$, and $\mathbf{u}(t_k)$ are denoted \mathbf{y}_k , \mathbf{s}_k , and \mathbf{u}_k , respectively. $\mathbf{L}(\mathbf{r})$ is denoted \mathbf{L} for simplicity. Then, Eq. (5.134) is rewritten as

$$\mathbf{y}_k = \mathbf{L}\mathbf{s}_k + \mathbf{A}\mathbf{u}_k + \boldsymbol{\varepsilon}. \quad (5.135)$$

The goal of this algorithm is to estimate \mathbf{s}_k from the data \mathbf{y}_k . Once \mathbf{s}_k is obtained, the algorithm pointing location \mathbf{r} is scanned over the whole brain region to obtain the spatiotemporal reconstruction of the source activity over the whole brain. To formulate

the estimation problem based on the VBFA algorithm, Eq. (5.135) is rewritten as

$$\mathbf{y}_k = \mathbf{L}\mathbf{s}_k + \mathbf{A}\mathbf{u}_k + \boldsymbol{\varepsilon} = [\mathbf{L}, \mathbf{A}] \begin{bmatrix} \mathbf{s}_k \\ \mathbf{u}_k \end{bmatrix} + \boldsymbol{\varepsilon} = \mathbf{A}_s \mathbf{z}_k + \boldsymbol{\varepsilon}, \quad (5.136)$$

where $\mathbf{z}_k = [\mathbf{s}_k^T, \mathbf{u}_k^T]^T$, and $\mathbf{A}_s = [\mathbf{L}, \mathbf{A}]$. This equation indicates that the problem of estimating \mathbf{s}_k can be expressed by the factor analysis model using the augmented factor vector \mathbf{z}_k and the augmented mixing matrix \mathbf{A}_s . Equation (5.136) is, in principle, the same as Eq. (5.107). Therefore, the algorithm developed here is very similar to the PFA algorithm described in Sect. 5.4.

5.5.2 Probability Model

We assume the prior distribution for \mathbf{s}_k to be zero-mean Gaussian, such that

$$p(\mathbf{s}_k) = \mathcal{N}(\mathbf{s}_k | \mathbf{0}, \boldsymbol{\Phi}^{-1}), \quad (5.137)$$

where $\boldsymbol{\Phi}$ is a 3×3 (non-diagonal) precision matrix of this prior distribution. The prior distribution of the factor vector \mathbf{u}_k is assumed to be:

$$p(\mathbf{u}_k) = \mathcal{N}(\mathbf{u}_k | \mathbf{0}, \mathbf{I}). \quad (5.138)$$

Then, the prior distribution for \mathbf{z}_k is derived as

$$\begin{aligned} p(\mathbf{z}_k) &= p(\mathbf{s}_k)p(\mathbf{u}_k) = \mathcal{N}(\mathbf{s}_k | \mathbf{0}, \boldsymbol{\Phi}^{-1})\mathcal{N}(\mathbf{u}_k | \mathbf{0}, \mathbf{I}) \\ &= \left| \frac{\boldsymbol{\Phi}}{2\pi} \right|^{1/2} \exp \left[-\frac{1}{2} \mathbf{s}_k^T \boldsymbol{\Phi} \mathbf{s}_k \right] \left| \frac{\mathbf{I}}{2\pi} \right|^{1/2} \exp \left[-\frac{1}{2} \mathbf{u}_k^T \mathbf{u}_k \right] \\ &= \left| \frac{\tilde{\boldsymbol{\Phi}}}{2\pi} \right|^{1/2} \exp \left[-\frac{1}{2} \mathbf{z}_k^T \tilde{\boldsymbol{\Phi}} \mathbf{z}_k \right] = \mathcal{N}(\mathbf{z}_k | \mathbf{0}, \tilde{\boldsymbol{\Phi}}^{-1}), \end{aligned} \quad (5.139)$$

where

$$\tilde{\boldsymbol{\Phi}} = \begin{bmatrix} \boldsymbol{\Phi} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

Equation (5.139) indicates that the prior distribution of \mathbf{z}_k is the mean-zero Gaussian with the precision matrix of $\tilde{\boldsymbol{\Phi}}$. The sensor noise is also assumed to be zero-mean Gaussian with a diagonal precision matrix $\boldsymbol{\Lambda}$. Hence, we have

$$p(\mathbf{y}_k | \mathbf{s}_k, \mathbf{u}_k) = \mathcal{N}(\mathbf{y}_k | \mathbf{L}\mathbf{s}_k + \mathbf{A}\mathbf{u}_k, \boldsymbol{\Lambda}^{-1}) = \mathcal{N}(\mathbf{y}_k | \mathbf{A}_s \mathbf{z}_k, \boldsymbol{\Lambda}^{-1}). \quad (5.140)$$

We assume the same prior distribution for the mixing matrix \mathbf{A} , as shown in Eq. (5.42).

5.5.3 VBEM Algorithm

5.5.3.1 E-Step

The posterior $p(z|y)$ is Gaussian, and assumed to be,

$$p(z|y) = \prod_{k=1}^K p(z_k|y_k), \quad \text{and} \quad p(z_k|y_k) = \mathcal{N}(z_k|\bar{z}_k, \boldsymbol{\Gamma}^{-1}), \quad (5.141)$$

where \bar{z}_k and $\boldsymbol{\Gamma}$ are, respectively, the mean and precision of $p(z_k|y_k)$. Using similar arguments as for the E-step of the PFA algorithm, the estimate of $p(z_k|y_k)$, $\hat{p}(z_k|y_k)$, is given by

$$\log \hat{p}(z_k|y_k) = E_A[\log p(y_k|z_k, A)] + \log p(z_k), \quad (5.142)$$

where we omit terms that do not contain z_k . The precision $\boldsymbol{\Gamma}$ of the posterior is derived from the coefficient of $\frac{\partial}{\partial z_k} \hat{p}(z_k|y_k)$, and the mean \bar{z}_k is derived as the z_k that makes this derivative zero.

Defining $\bar{A}_s = E_A([L, A]) = [L, \bar{A}]$, we obtain

$$\begin{aligned} \frac{\partial}{\partial z_k} \log \hat{p}(z_k|y_k) &= E_A[A_s^T \boldsymbol{\Lambda} (y_k - A_s z_k)] - \tilde{\boldsymbol{\Phi}} z_k \\ &= E_A[A_s^T \boldsymbol{\Lambda} y_k] - E_A[A_s^T \boldsymbol{\Lambda} A_s] z_k - \tilde{\boldsymbol{\Phi}} z_k, \end{aligned} \quad (5.143)$$

where

$$E_A[A_s^T \boldsymbol{\Lambda} y_k] = \bar{A}_s^T \boldsymbol{\Lambda} y_k$$

and

$$\begin{aligned} E_A[A_s^T \boldsymbol{\Lambda} A_s] &= \begin{bmatrix} E_A[L^T \boldsymbol{\Lambda} L] & E_A[L^T \boldsymbol{\Lambda} A] \\ E_A[A^T \boldsymbol{\Lambda} L] & E_A[A^T \boldsymbol{\Lambda} A] \end{bmatrix} \\ &= \begin{bmatrix} L^T \boldsymbol{\Lambda} L & L^T \boldsymbol{\Lambda} \bar{A} \\ \bar{A}^T \boldsymbol{\Lambda} L & \bar{A}^T \boldsymbol{\Lambda} \bar{A} + M \boldsymbol{\Psi}^{-1} \end{bmatrix} = \bar{A}_s^T \boldsymbol{\Lambda} \bar{A}_s + M \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Psi}^{-1} \end{bmatrix}. \end{aligned} \quad (5.144)$$

Hence, we have

$$E_A[A_s^T \boldsymbol{\Lambda} A_s] z_k = (\bar{A}_s^T \boldsymbol{\Lambda} \bar{A}_s + M \boldsymbol{\Psi}_s^{-1}) z_k, \quad (5.145)$$

where

$$\boldsymbol{\Psi}_s^{-1} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Psi}^{-1} \end{bmatrix}. \quad (5.146)$$

Thus, we derive

$$\frac{\partial}{\partial \mathbf{z}_k} \log \hat{p}(\mathbf{z}_k | \mathbf{y}_k) = \bar{\mathbf{A}}_s^T \boldsymbol{\Lambda} \mathbf{y}_k - (\bar{\mathbf{A}}_s^T \boldsymbol{\Lambda} \bar{\mathbf{A}}_s + M \boldsymbol{\Psi}_s^{-1} + \tilde{\boldsymbol{\Phi}}) \mathbf{z}_k. \quad (5.147)$$

Consequently, the coefficient of \mathbf{z}_k gives

$$\boldsymbol{\Gamma} = \bar{\mathbf{A}}_s^T \boldsymbol{\Lambda} \bar{\mathbf{A}}_s + M \boldsymbol{\Psi}_s^{-1} + \tilde{\boldsymbol{\Phi}}. \quad (5.148)$$

The mean $\bar{\mathbf{z}}_k$ is obtained as

$$\bar{\mathbf{z}}_k = \begin{bmatrix} \bar{s}_k \\ \bar{\mathbf{u}}_k \end{bmatrix} = \boldsymbol{\Gamma}^{-1} \begin{bmatrix} \mathbf{L}^T \\ \bar{\mathbf{A}}^T \end{bmatrix} \boldsymbol{\Lambda} \mathbf{y}_k. \quad (5.149)$$

5.5.3.2 M-Step

The derivation for the M step is also very similar to the M step of the PFA-VBEM algorithm in Sect. 5.4.3.2. The estimation of $\hat{p}(A|\mathbf{y})$ is based on

$$\begin{aligned} \log \hat{p}(A|\mathbf{y}) &= -\frac{1}{2} E_{\mathbf{z}} \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k) \right] \\ &\quad - \frac{1}{2} \sum_{j=1}^M \mathbf{a}_j^T \lambda_j \boldsymbol{\alpha} \mathbf{a}_j. \end{aligned} \quad (5.150)$$

The derivative $\frac{\partial}{\partial \mathbf{A}} \log \hat{p}(A|\mathbf{y})$ is obtained as,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{A}} \log \hat{p}(A|\mathbf{y}) &= E_{\mathbf{z}} \left[\boldsymbol{\Lambda} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k) \mathbf{u}_k^T \right] - \boldsymbol{\Lambda} \boldsymbol{\Lambda} \boldsymbol{\alpha} \\ &= \boldsymbol{\Lambda} \mathbf{R}_{yu} - \boldsymbol{\Lambda} \mathbf{L} \mathbf{R}_{su} - \boldsymbol{\Lambda} \boldsymbol{\Lambda} (\mathbf{R}_{uu} + \boldsymbol{\alpha}). \end{aligned} \quad (5.151)$$

The precision, $\lambda_j \boldsymbol{\Psi}$, is equal to the coefficient of \mathbf{a}_j , and it is obtained as $\lambda_j (\mathbf{R}_{uu} + \boldsymbol{\alpha})$. Hence we have

$$\boldsymbol{\Psi} = (\mathbf{R}_{uu} + \boldsymbol{\alpha}). \quad (5.152)$$

Setting the derivative to zero gives $\bar{\mathbf{A}}$, i.e.,

$$\bar{\mathbf{A}} = (\mathbf{R}_{yu} - \mathbf{L} \mathbf{R}_{su})(\mathbf{R}_{uu} + \boldsymbol{\alpha})^{-1} = (\mathbf{R}_{yu} - \mathbf{L} \mathbf{R}_{su}) \boldsymbol{\Psi}^{-1}. \quad (5.153)$$

In the equations above, \mathbf{R}_{su} and \mathbf{R}_{uu} are obtained in the following manner. Defining

$$\boldsymbol{\Gamma}^{-1} = \begin{bmatrix} \tilde{\boldsymbol{\Gamma}}_{ss} & \tilde{\boldsymbol{\Gamma}}_{su} \\ \tilde{\boldsymbol{\Gamma}}_{us} & \tilde{\boldsymbol{\Gamma}}_{uu} \end{bmatrix}, \quad (5.154)$$

we have

$$\begin{aligned} \begin{bmatrix} \mathbf{R}_{ss} & \mathbf{R}_{su} \\ \mathbf{R}_{us} & \mathbf{R}_{uu} \end{bmatrix} &= E_z \left[\sum_{k=1}^K \begin{bmatrix} \mathbf{s}_k \\ \mathbf{u}_k \end{bmatrix} [\mathbf{s}_k^T, \mathbf{u}_k^T] \right] \\ &= \begin{bmatrix} \sum_{k=1}^K \bar{\mathbf{s}}_k \bar{\mathbf{s}}_k^T & \sum_{k=1}^K \bar{\mathbf{s}}_k \bar{\mathbf{u}}_k^T \\ \sum_{k=1}^K \bar{\mathbf{u}}_k \bar{\mathbf{s}}_k^T & \sum_{k=1}^K \bar{\mathbf{u}}_k \bar{\mathbf{u}}_k^T \end{bmatrix} + K \begin{bmatrix} \tilde{\boldsymbol{\Gamma}}_{ss} & \tilde{\boldsymbol{\Gamma}}_{su} \\ \tilde{\boldsymbol{\Gamma}}_{us} & \tilde{\boldsymbol{\Gamma}}_{uu} \end{bmatrix}. \end{aligned} \quad (5.155)$$

Therefore, we obtain:

$$\mathbf{R}_{ss} = \sum_{k=1}^K \bar{\mathbf{s}}_k \bar{\mathbf{s}}_k^T + K \tilde{\boldsymbol{\Gamma}}_{ss}, \quad (5.156)$$

$$\mathbf{R}_{su} = \sum_{k=1}^K \bar{\mathbf{s}}_k \bar{\mathbf{u}}_k^T + K \tilde{\boldsymbol{\Gamma}}_{su}, \quad (5.157)$$

$$\mathbf{R}_{uu} = \sum_{k=1}^K \bar{\mathbf{u}}_k \bar{\mathbf{u}}_k^T + K \tilde{\boldsymbol{\Gamma}}_{uu}. \quad (5.158)$$

Also, \mathbf{R}_{yu} is obtained as

$$\mathbf{R}_{yu} = \sum_{k=1}^K \mathbf{y}_k \bar{\mathbf{u}}_k^T. \quad (5.159)$$

5.5.3.3 Update Equations for Hyperparameters

The update equations for $\boldsymbol{\Phi}$, $\boldsymbol{\Lambda}$, and $\boldsymbol{\alpha}$ are derived by maximizing the free energy, which is rewritten as

$$\mathcal{F}(\boldsymbol{\Lambda}, \boldsymbol{\alpha}, \boldsymbol{\Phi}) = E_{(z, A)}[\log p(\mathbf{y}|z, \mathbf{A}_s) + \log p(z) + \log p(\mathbf{A})], \quad (5.160)$$

where the terms $-\log \hat{p}(z|\mathbf{y}) - \log \hat{p}(\mathbf{A}|\mathbf{y})$ are ignored because these terms contain none of these hyperparameters. To compute $\mathcal{F}(\boldsymbol{\Lambda}, \boldsymbol{\alpha}, \boldsymbol{\Phi})$, we use the relationships,

$$\begin{aligned}
\log p(\mathbf{y}|\mathbf{z}, \mathbf{A}_s) &= \sum_{k=1}^K \log p(\mathbf{y}_k|\mathbf{z}_k, \mathbf{A}_s) \\
&= \frac{K}{2} \log |\mathbf{\Lambda}| - \frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k)^T \mathbf{\Lambda} (\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k),
\end{aligned} \tag{5.161}$$

$$\begin{aligned}
\log p(\mathbf{z}) &= \sum_{k=1}^K \log p(z_k) = \frac{K}{2} \log |\tilde{\boldsymbol{\Phi}}| - \frac{1}{2} \sum_{k=1}^K \mathbf{z}_k^T \tilde{\boldsymbol{\Phi}} \mathbf{z}_k \\
&= \frac{K}{2} \log |\boldsymbol{\Phi}| - \frac{1}{2} \sum_{k=1}^K \mathbf{s}_k^T \boldsymbol{\Phi} \mathbf{s}_k - \frac{1}{2} \sum_{k=1}^K \mathbf{u}_k^T \mathbf{u}_k,
\end{aligned} \tag{5.162}$$

and

$$\log p(\mathbf{A}) = \frac{1}{2} \sum_{j=1}^M \log |\lambda_j \boldsymbol{\alpha}| - \frac{1}{2} \sum_{j=1}^M \mathbf{a}_j^T \lambda_j \boldsymbol{\alpha} \mathbf{a}_j. \tag{5.163}$$

Let us derive the update equation for $\boldsymbol{\Phi}$. In Eq. (5.160), the only term that contains $\boldsymbol{\Phi}$ is $\log p(\mathbf{z})$. Thus, we have

$$\begin{aligned}
\frac{1}{\partial \boldsymbol{\Phi}} \mathcal{F}(\mathbf{\Lambda}, \boldsymbol{\alpha}, \boldsymbol{\Phi}) &= \frac{1}{\partial \boldsymbol{\Phi}} \frac{K}{2} \log |\boldsymbol{\Phi}| - \frac{1}{\partial \boldsymbol{\Phi}} \frac{1}{2} E_{(\mathbf{z}, \mathbf{A})} \left[\sum_{k=1}^K \mathbf{s}_k^T \boldsymbol{\Phi} \mathbf{s}_k \right] \\
&= \frac{K}{2} \boldsymbol{\Phi}^{-1} - \frac{1}{2} E_s \left[\sum_{k=1}^K \mathbf{s}_k \mathbf{s}_k^T \right] = \frac{K}{2} \boldsymbol{\Phi}^{-1} - \frac{1}{2} \mathbf{R}_{ss},
\end{aligned} \tag{5.164}$$

and setting the right-hand side to zero gives the update equation,

$$\boldsymbol{\Phi}^{-1} = \frac{1}{K} \mathbf{R}_{ss}, \tag{5.165}$$

where \mathbf{R}_{ss} is obtained in Eq. (5.156). We can derive the update equation for $\boldsymbol{\alpha}$ in a similar manner. However, Since in Eq. (5.160), the only term that contains $\boldsymbol{\alpha}$ is $\log p(\mathbf{A})$, the update equation for $\boldsymbol{\alpha}$ is exactly the same as that in Eq. (5.79). The update equation for $\mathbf{\Lambda}$ is given by

$$\mathbf{\Lambda}^{-1} = \frac{1}{K} \text{diag}[\mathbf{R}_{yy} - \mathbf{R}_{ys} \mathbf{L}^T - \mathbf{L} \mathbf{R}_{sy} + \mathbf{L} \mathbf{R}_{ss} \mathbf{L}^T - \bar{\mathbf{A}} \boldsymbol{\Psi} \bar{\mathbf{A}}^T]. \tag{5.166}$$

Since the derivation of the equation above is lengthy, it is presented in Sect. 5.7.4.

5.5.4 Summary of the Saketini Algorithm

The Saketini algorithm is summarized as follows. It estimates the source time course at a pointing location, based on the data model in Eq. (5.134). It uses the VBEM algorithm in which the E-step updates the posterior precision $\boldsymbol{\Gamma}$ and the posterior mean $\bar{\mathbf{z}}_k$ by using Eqs. (5.148) and (5.149). The M-step updates $\boldsymbol{\Psi}$ and \boldsymbol{A} using Eqs. (5.152) and (5.153). The hyperparameter $\boldsymbol{\Phi}$ is updated using Eq. (5.165). The hyperparameter α is updated using (5.79). The noise precision \boldsymbol{A} is updated using Eq. (5.166). The free energy can be computed using the same equation as in Eq. (5.131) with the sole modification of adding $\log |\boldsymbol{\Phi}|$. That is, the free energy is expressed as

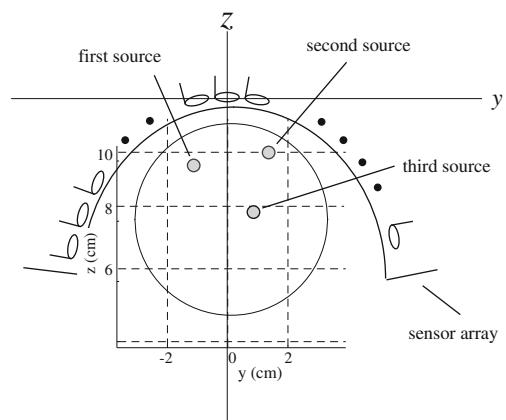
$$\mathcal{F} = \frac{K}{2} \log \frac{|\boldsymbol{A}||\boldsymbol{\Phi}|}{|\boldsymbol{\Gamma}|} - \frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{A} \mathbf{y}_k + \frac{1}{2} \sum_{k=1}^K \bar{\mathbf{z}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{z}}_k + \frac{M}{2} \log \frac{|\alpha|}{|\boldsymbol{\Psi}|}. \quad (5.167)$$

This free energy forms the lower bound of the marginal likelihood, and it is used for monitoring the progress of the VBEM iteration.

5.6 Numerical Examples

Computer simulation was performed to present typical results from the methods mentioned in this chapter. An alignment of the 275-sensor array from the Omega™ (VMS Medtech, Coquitlam, Canada) neuromagnetometer was used. The coordinate system and source-sensor configuration used in the computer simulation are depicted in Fig. 5.1. A vertical plane was assumed at the middle of the whole-head sensor array ($x = 0$ cm), and three sources were assumed to exist on this plane. The time courses shown in the upper three panels of Fig. 5.2 were assigned to the three sources.

Fig. 5.1 The coordinate system and source-sensor configuration used in the computer simulation. A vertical plane was assumed at the *middle* of the whole-head sensor array ($x = 0$ cm), and three sources were assumed to exist on this plane at $x = 0$ cm. The *large circle* shows the source space and the *small filled circles* show the locations of the three sources



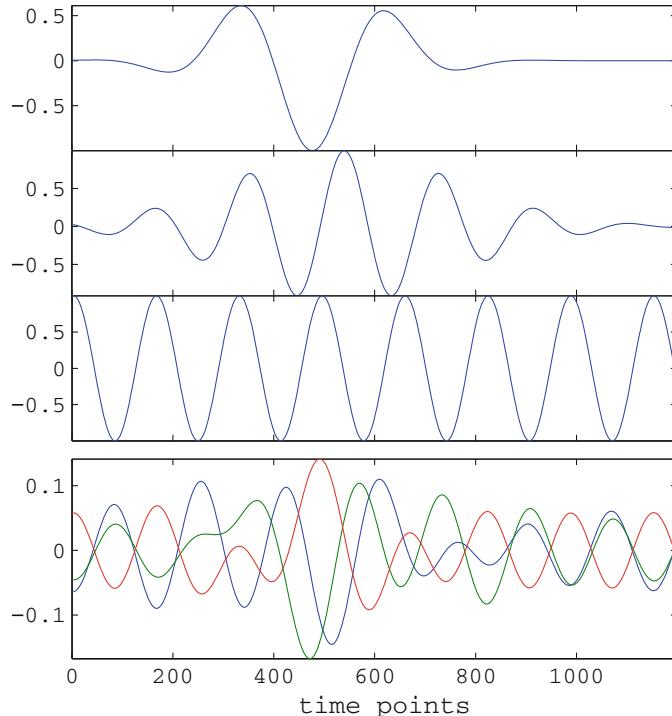


Fig. 5.2 The time courses of the three sources are shown in the *upper* three panels. The *bottom* panel shows the signal-only magnetic recordings (before noise is added). The selected three time courses from sensor #1, #101, and #201 are shown. The ordinates show relative values, and the abscissa shows the time points. A total of 1,200 time point data was generated

Simulated magnetic recordings were computed by projecting the time courses of the three sources onto the sensor space using the sensor lead field at the source locations. These recordings are shown for three selected sensor channels in the bottom panel of Fig. 5.2.

We first performed denoising experiments using the BFA and VBFA algorithms. We generated simulated MEG data by adding the sensor noise onto the computed magnetic recordings shown in the bottom panel of Fig. 5.2. Here, the noise generation was performed using a Gaussian random number, and the signal-to-noise ratio (SNR) was set to one. The noise-added simulated sensor data is shown in the top panel of Fig. 5.3 for the same selected sensor channels.

The BFA algorithm was applied to these simulated MEG data and the signal component was estimated by computing $A\bar{u}_k$. The results of this experiment are presented in Fig. 5.3. The results with the number of factors L set to 3 are shown in the second panel from the top, and the results with L set to 20 are in the third panel. Since we assume that three sources exist in this numerical experiment, the correct value of L is three.

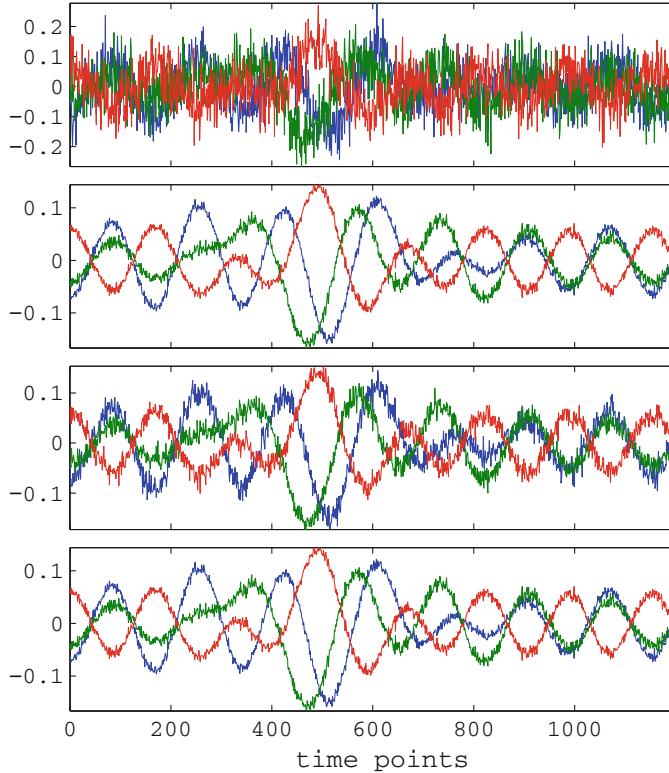


Fig. 5.3 Results of denoising experiments using BFA and VBFA. The *top panel* shows the sensor recordings before denoising. The *second panel* from the *top* shows the denoising results by BFA with the model order L set to three. The *third panel* shows the denoising results by BFA with L set to twenty. The *bottom panel* shows the denoising results by VBFA with L set to twenty. The time courses of the selected three sensor channels, (channel #1, channel #101, and channel #201), are shown. The ordinates show relative values, and abscissa shows the time points

Results in the second panel (results with the correct L) are better than the results in the third panel (results with a significantly overespecified L), and we can see that the overspecification of the number of factors, L , reduces the denoising capability of the algorithm. The VBFA algorithm was then applied, and the results are shown in the bottom panel of Fig. 5.3. Here, although the number of factors L was set to 20, results close to those obtained using the BFA algorithm with the correct value of L were obtained. The results demonstrate the fact that the VBFA algorithm incorporates the model order determination and is insensitive to the overspecification of the number of factors.

We next performed experiments on the interference removal using the PFA algorithm. In this experiment, the simulated MEG recordings were generated with 2,400 time points where the first half period does not contain the signal activity, which is contained in the second half period. The interference data was computed by

generating fifty interference sources in the source space. The interference sources had Gaussian random time courses, and the interference magnetic recordings were generated by projecting these time courses onto the sensor space using the sensor lead field at the interference source locations. The simulated MEG recordings were generated by adding the interference (and a small amount of sensor noise) onto the same simulated signal-only sensor data in the previous denoising experiment.

In the PFA application, the first half period was used as the control period and the mixing matrix of the interference \mathbf{B} and the diagonal noise precision matrix \mathbf{A} were estimated using this control period. We then estimated the signal mixing matrix \mathbf{A} and the signal factor activity \mathbf{u}_k using the second half period. The signal component was estimated using $\hat{\mathbf{y}}_k^S$ in Eq.(5.132).

Results are shown in Fig. 5.4. In this figure, the top panel shows the control period (the first half period) of the original interference-added sensor data, and the middle panel shows the target period (the second half period) of the original interference-added sensor data. The bottom panel shows the interference-removed sensor data, $\hat{\mathbf{y}}_k^S$. Here, the number of factors for the interference was set to 100, and that for the signal of interest was set to 10. Note that the true number of independent interference activities is 50, and that for the signal activity is 3. The results in Fig. 5.4 demonstrate that

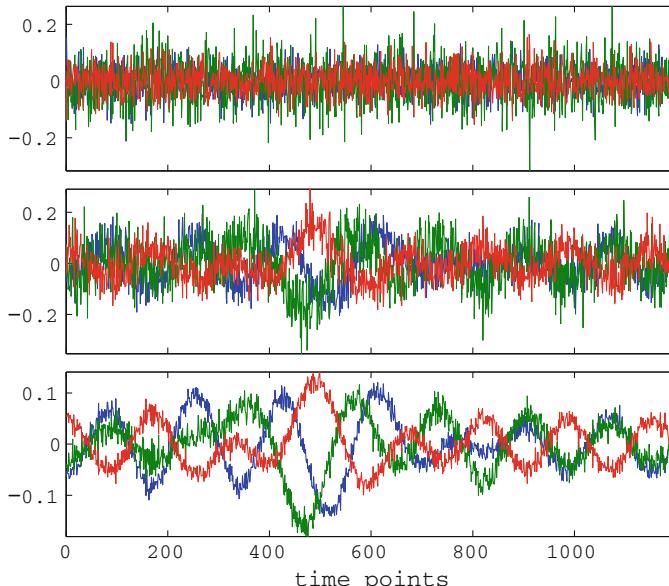


Fig. 5.4 Results of interference-removal experiments using the PFA algorithm. The *top panel* shows the control period (the first 1,200 time points) of the interference-added simulated sensor data. The *middle panel* shows the target period (the second 1,200 time points) of the interference-added simulated sensor data. The *bottom panel* shows the results of interference removal experiments. Here, results of $\hat{\mathbf{y}}_k^S$ in Eq.(5.132) are plotted. The selected three sensor time courses (channel #1, channel #101, and channel #201) are shown in these three panels

the PFA algorithm successfully removes the interference, even though the number of factors for the signal and interference components are significantly overspecified.

We next performed a source reconstruction experiment. We computed the interference-removed covariance matrix $\tilde{\mathbf{R}}$ in Eq. (5.133), and used it with the adaptive beamformer algorithm. The results are shown in Fig. 5.5. In Fig. 5.5a, the results of source reconstruction, obtained using the signal-only data (data before adding the interferences), are shown. The results of source reconstruction obtained with the interference-added data are shown in Fig. 5.5b. We can observe a certain amount of distortion in Fig. 5.5b. The results of reconstruction using $\tilde{\mathbf{R}}$ are shown in Fig. 5.5c. The distortion is significantly reduced, demonstrating the effectiveness of the PFA algorithm.

Finally, a source reconstruction experiment using the Saketini algorithm was performed. Simulated MEG data was generated in which the signal-to-noise ratio (SNR) was set equal to two. We computed the power map using $\text{tr}(\Phi^{-1})$ with the number of factors L set to 20. The results are shown in Fig. 5.6. The reconstructed source power map on the plane of $x = 0$ is shown in Fig. 5.6a. The reconstructed time courses

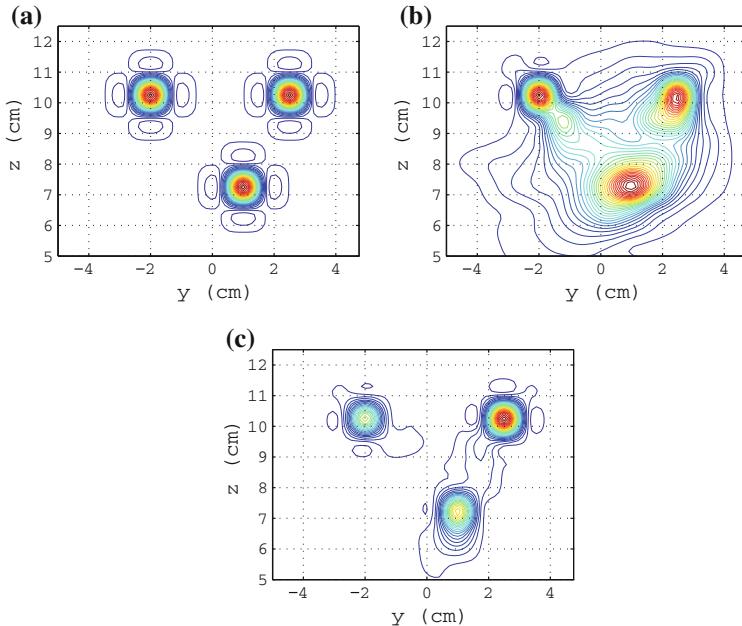


Fig. 5.5 Results of source reconstruction experiments using the interference-removed covariance matrix $\tilde{\mathbf{R}}$ in Eq. (5.133). The adaptive beamformer algorithm was used for source reconstruction, and the reconstructed source distribution at $x = 0$ cm is shown. **a** The results obtained with a sample covariance matrix computed from signal-only simulated recordings. **b** The results obtained with a sample covariance matrix computed from the interference-added sensor data. **c** The results of the interference-removal experiment, obtained with $\tilde{\mathbf{R}}$ in Eq. (5.133). The (y, z) coordinates of the three sources were set to $(-2.0, 10.2)$ cm, $(2.5, 10.2)$ cm, and $(1.0, 7.2)$ cm in this experiment

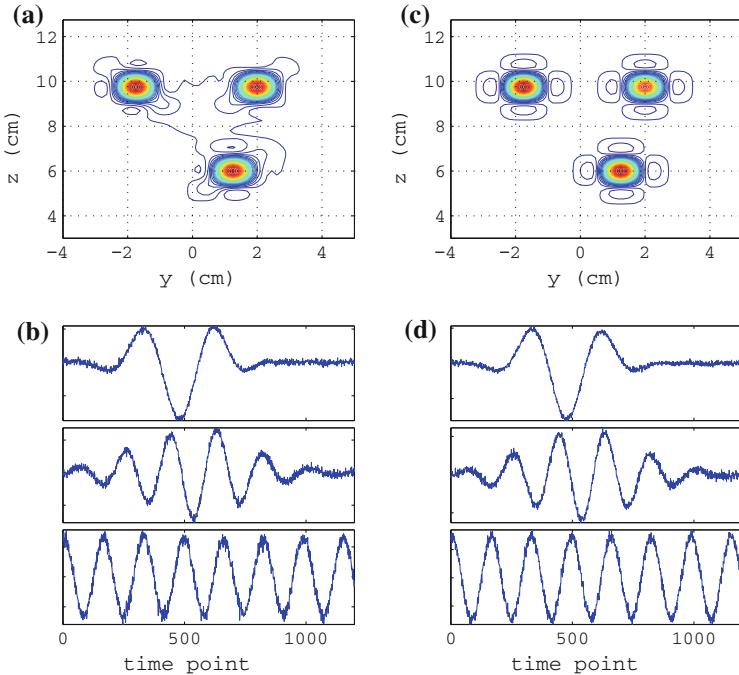


Fig. 5.6 Results of source reconstruction experiments using the Saketini algorithm. **a** The reconstructed source power map on the plane at $x = 0\text{ cm}$. **b** The reconstructed time courses of the three sources: The time courses of the first, second, and third sources are shown in the *top*, *middle*, and *bottom panels*, respectively. **c** The reconstructed source power map on the plane at $x = 0\text{ cm}$, obtained using the Champagne algorithm. **d** The reconstructed time courses of the three sources, obtained using the Champagne algorithm. The locations (y, z) of the first, second, and third sources were set to $(-1.8, 9.8)$, $(2.0, 9.8)$, and $(1.3, 6.0)$, respectively

of the voxels closest to the three sources are shown in Fig. 5.6b. The source power map and source time courses obtained using the Champagne algorithm are shown in Fig. 5.6c, d for comparison. The results show that both algorithms are capable of providing accurate spatiotemporal reconstruction of the three sources.

5.7 Appendix to This Chapter

5.7.1 Proof of Eq. (5.84)

We present the proof of Eq. (5.84), which is rewritten as

$$E_A[A\Psi A^T] = \bar{A}\Psi\bar{A}^T + L\Lambda^{-1} = R_{yu}\bar{A}^T + L\Lambda^{-1}.$$

First, we have

$$[\mathbf{A}\Psi\mathbf{A}^T]_{i,j} = \sum_{k,\ell} A_{i,k} [\Psi]_{k,\ell} A_{j,\ell} = \sum_{k,\ell} [\Psi]_{k,\ell} A_{i,k} A_{j,\ell}.$$

Equation (5.44) expresses the posterior distribution of the mixing matrix \mathbf{A} where the two elements, $A_{i,k}$ and $A_{j,\ell}$, are independent when $i \neq j$ and when $i = j$, the covariance of these elements is equal to $[\lambda_j^{-1}\Psi^{-1}]_{k,\ell}$. Thus, we have

$$E_A[A_{i,k} A_{j,\ell}] = \bar{A}_{i,k} \bar{A}_{j,\ell} + \delta_{i,j} \frac{1}{\lambda_j} [\Psi^{-1}]_{k,\ell}. \quad (5.168)$$

Therefore, we get

$$E_A[(\mathbf{A}\Psi\mathbf{A}^T)_{i,j}] = \sum_{k,\ell} [\Psi]_{k,\ell} \bar{A}_{i,k} \bar{A}_{j,\ell} + \sum_{k,\ell} [\Psi]_{k,\ell} \delta_{i,j} \frac{1}{\lambda_j} [\Psi^{-1}]_{k,\ell}. \quad (5.169)$$

The first term on the right-hand side of the equation above is rewritten as

$$\sum_{k,\ell} [\Psi]_{k,\ell} \bar{A}_{i,k} \bar{A}_{j,\ell} = [\bar{\mathbf{A}}\Psi\bar{\mathbf{A}}^T]_{i,j}.$$

Considering

$$\sum_{k,\ell} [\Psi]_{k,\ell} [\Psi^{-1}]_{k,\ell} = \text{tr}(\Psi\Psi^{-1}) = \text{tr}(\mathbf{I}) = L, \quad (5.170)$$

the second term in Eq. (5.169) is rewritten as

$$\sum_{k,\ell} [\Psi]_{k,\ell} \delta_{i,j} \frac{1}{\lambda_j} [\Psi^{-1}]_{k,\ell} = L \delta_{i,j} \frac{1}{\lambda_j} = L [\Lambda^{-1}]_{i,j}, \quad (5.171)$$

because Λ is diagonal. Consequently, we get the following relationship:

$$E_A[\mathbf{A}\Psi\mathbf{A}^T] = \bar{\mathbf{A}}\Psi\bar{\mathbf{A}}^T + L\Lambda^{-1}. \quad (5.172)$$

Using Eqs. (5.68) and (5.70), we obtain the relationship

$$\bar{\mathbf{A}} = \mathbf{R}_{yu} \Psi^{-1} \quad \text{or} \quad \bar{\mathbf{A}}\Psi = \mathbf{R}_{yu}.$$

Thus, we finally derive Eq. (5.84), such that

$$E_A[\mathbf{A}\Psi\mathbf{A}^T] = \mathbf{R}_{yu} \bar{\mathbf{A}}^T + L\Lambda^{-1}. \quad (5.173)$$

5.7.2 Proof of Eq. (5.94)

We provide a proof of Eq. (5.94), which is repeated here for convenience.

$$\begin{aligned} E_{(\mathbf{A}, \mathbf{u})} & \left[-\frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k) - \frac{1}{2} \sum_{k=1}^K \mathbf{u}_k \mathbf{u}_k^T \right] \\ & = -\frac{1}{2} \sum_{j=1}^K \mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{y}_k + \frac{1}{2} \sum_{j=1}^K \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{u}}_k. \end{aligned}$$

The left-hand side of the equation above is changed to:

$$\begin{aligned} E_{(\mathbf{A}, \mathbf{u})} & \left[-\frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{A}\mathbf{u}_k) - \frac{1}{2} \sum_{k=1}^K \mathbf{u}_k \mathbf{u}_k^T \right] \\ & = -\frac{1}{2} \sum_{k=1}^K \mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{y}_k \\ & \quad - \frac{1}{2} \sum_{k=1}^K E_{(\mathbf{A}, \mathbf{u})} \left[-\mathbf{u}_k^T \mathbf{A}^T \boldsymbol{\Lambda} \mathbf{y}_k - \mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{A} \mathbf{u}_k + \mathbf{u}_k^T \mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} \mathbf{u}_k + \mathbf{u}_k^T \mathbf{u}_k \right]. \end{aligned} \tag{5.174}$$

On the right-hand side of Eq. (5.174), the following relationship:

$$\begin{aligned} E_{(\mathbf{A}, \mathbf{u})} \left[\mathbf{u}_k^T \mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} \mathbf{u}_k + \mathbf{u}_k^T \mathbf{u}_k \right] & = E_{\mathbf{u}} \left[\mathbf{u}_k^T \left[E_{\mathbf{A}} \left(\mathbf{A}^T \boldsymbol{\Lambda} \mathbf{A} \right) + \mathbf{I} \right] \mathbf{u}_k \right] \\ & = E_{\mathbf{u}} \left[\mathbf{u}_k^T \left[\bar{\mathbf{A}}^T \boldsymbol{\Lambda} \bar{\mathbf{A}} + M \boldsymbol{\Psi}^{-1} + \mathbf{I} \right] \mathbf{u}_k \right] \\ & = E_{\mathbf{u}} \left[\mathbf{u}_k^T \boldsymbol{\Gamma} \mathbf{u}_k \right] = E_{\mathbf{u}} \left[\text{tr} \left(\mathbf{u}_k \mathbf{u}_k^T \boldsymbol{\Gamma} \right) \right] \\ & = \text{tr} \left[E_{\mathbf{u}} \left(\mathbf{u}_k \mathbf{u}_k^T \right) \boldsymbol{\Gamma} \right] = \text{tr} \left[(\bar{\mathbf{u}}_k \bar{\mathbf{u}}_k^T + \boldsymbol{\Gamma}^{-1}) \boldsymbol{\Gamma} \right] \\ & = \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{u}}_k \end{aligned} \tag{5.175}$$

holds, where constant terms are omitted. In the equation above, we use Eq. (5.56). Also, the relationship,

$$\begin{aligned} E_{(\mathbf{A}, \mathbf{u})} \left[\mathbf{y}_k^T \boldsymbol{\Lambda} \mathbf{A} \mathbf{u}_k + \mathbf{u}_k^T \mathbf{A}^T \boldsymbol{\Lambda} \mathbf{y}_k \right] & = \mathbf{y}_k^T \boldsymbol{\Lambda} \bar{\mathbf{A}} \bar{\mathbf{u}}_k + \bar{\mathbf{u}}_k^T \bar{\mathbf{A}}^T \boldsymbol{\Lambda} \mathbf{y}_k \\ & = \mathbf{y}_k^T \boldsymbol{\Lambda} \bar{\mathbf{A}} \boldsymbol{\Gamma}^{-1} \boldsymbol{\Gamma} \bar{\mathbf{u}}_k + \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \boldsymbol{\Gamma}^{-1} \bar{\mathbf{A}}^T \boldsymbol{\Lambda} \mathbf{y}_k \\ & = 2 \bar{\mathbf{u}}_k^T \boldsymbol{\Gamma} \bar{\mathbf{u}}_k \end{aligned} \tag{5.176}$$

holds. Substituting Eqs. (5.175) and (5.176) into (5.174), we get Eq. (5.94).

5.7.3 Proof of Eq. (5.103)

Next, the proof of Eq. (5.103) is presented, which is

$$E_A [\mathbf{A} \mathbf{R}_{uu} \mathbf{A}^T] = \bar{\mathbf{A}} \mathbf{R}_{uu} \bar{\mathbf{A}}^T + \boldsymbol{\Lambda}^{-1} \text{tr}(\mathbf{R}_{uu} \boldsymbol{\Psi}).$$

The proof is quite similar to the one provided in Sect. 5.7.1. First, we have

$$[\mathbf{A} \mathbf{R}_{uu} \mathbf{A}^T]_{i,j} = \sum_{k,\ell} A_{i,k} [\mathbf{R}_{uu}]_{k,\ell} A_{j,\ell} = \sum_{k,\ell} [\mathbf{R}_{uu}]_{k,\ell} A_{i,k} A_{j,\ell}.$$

Using Eq. (5.168), we get

$$E_A [[\mathbf{A} \mathbf{R}_{uu} \mathbf{A}^T]_{i,j}] = \sum_{k,\ell} [\mathbf{R}_{uu}]_{k,\ell} \bar{A}_{i,k} \bar{A}_{j,\ell} + \sum_{k,\ell} [\mathbf{R}_{uu}]_{k,\ell} \delta_{i,j} \frac{1}{\lambda_j} [\boldsymbol{\Psi}^{-1}]_{k,\ell}. \quad (5.177)$$

The first term on the right-hand side of the equation above is rewritten as

$$\sum_{k,\ell} [\mathbf{R}_{uu}]_{k,\ell} \bar{A}_{i,k} \bar{A}_{j,\ell} = [\bar{\mathbf{A}} \mathbf{R}_{uu} \bar{\mathbf{A}}^T]_{i,j}.$$

Considering

$$\sum_{k,\ell} [\mathbf{R}_{uu}]_{k,\ell} [\boldsymbol{\Psi}^{-1}]_{k,\ell} = \text{tr}(\mathbf{R}_{uu} \boldsymbol{\Psi}^{-1}), \quad (5.178)$$

the second term in Eq. (5.177) is rewritten as

$$\sum_{k,\ell} [\mathbf{R}_{uu}]_{k,\ell} \delta_{i,j} \frac{1}{\lambda_j} [\boldsymbol{\Psi}^{-1}]_{k,\ell} = [\boldsymbol{\Lambda}^{-1}]_{i,j} \text{tr}(\mathbf{R}_{uu} \boldsymbol{\Psi}^{-1}). \quad (5.179)$$

Consequently, we get the following relationship:

$$E_A [\mathbf{A} \mathbf{R}_{uu} \mathbf{A}^T] = \bar{\mathbf{A}} \mathbf{R}_{uu} \bar{\mathbf{A}}^T + \boldsymbol{\Lambda}^{-1} \text{tr}(\mathbf{R}_{uu} \boldsymbol{\Psi}^{-1}). \quad (5.180)$$

5.7.4 Proof of Eq. (5.166)

Next, the proof of Eq. (5.166) is presented. To derive this equation, we should consider $\log p(\mathbf{y}|\mathbf{z}, \mathbf{A})$ and $\log p(\mathbf{A})$, because $\boldsymbol{\Lambda}$ is contained in these two terms. We have

$$\begin{aligned}
& \frac{\partial}{\partial \Lambda} \log p(\mathbf{y}|\mathbf{z}, \Lambda) \\
&= \frac{\partial}{\partial \Lambda} \left[\frac{K}{2} \log |\Lambda| - \frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k)^T \Lambda (\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k) \right] \\
&= \frac{K}{2} \Lambda^{-1} - \frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k)(\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k)^T. \tag{5.181}
\end{aligned}$$

According to Eq. (5.81), we have

$$\frac{\partial}{\partial \Lambda} \log p(\Lambda) = \frac{L}{2} \Lambda^{-1} - \frac{1}{2} \text{diag} [\mathbf{A}\boldsymbol{\alpha}\mathbf{A}^T]. \tag{5.182}$$

We can thus compute the derivative of $\mathcal{F}(\Lambda, \boldsymbol{\alpha}, \Phi)$ with respect to Λ , which is

$$\begin{aligned}
\frac{\partial \mathcal{F}}{\partial \Lambda} &= E_{(\Lambda, \mathbf{z})} \left[\frac{K}{2} \Lambda^{-1} - \frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k)(\mathbf{y}_k - \mathbf{L}\mathbf{s}_k - \mathbf{A}\mathbf{u}_k)^T \right. \\
&\quad \left. + \frac{L}{2} \Lambda^{-1} - \frac{1}{2} \mathbf{A}\boldsymbol{\alpha}\mathbf{A}^T \right] \\
&= \frac{K}{2} \Lambda^{-1} - \frac{1}{2} [\mathbf{R}_{yy} - \mathbf{R}_{yu} \bar{\mathbf{A}}^T - \bar{\mathbf{A}} \mathbf{R}_{uy} - \mathbf{R}_{ys} \mathbf{L}^T - \mathbf{L} \mathbf{R}_{sy} \\
&\quad + \mathbf{L} \mathbf{R}_{su} \bar{\mathbf{A}}^T + \bar{\mathbf{A}} \mathbf{R}_{us} \mathbf{L}^T + \mathbf{L} \mathbf{R}_{ss} \mathbf{L}^T] \\
&\quad - \frac{1}{2} E_{\mathbf{A}}[\mathbf{A}\mathbf{R}_{uu}\mathbf{A}^T] + \frac{L}{2} \Lambda^{-1} - \frac{1}{2} E_{\mathbf{A}}[\mathbf{A}\boldsymbol{\alpha}\mathbf{A}^T]. \tag{5.183}
\end{aligned}$$

According to Eqs. (5.83) and (5.84), we get

$$\begin{aligned}
E_{\mathbf{A}}[\mathbf{A}\mathbf{R}_{uu}\mathbf{A}^T + \mathbf{A}\boldsymbol{\alpha}\mathbf{A}^T] &= E_{\mathbf{A}}[\mathbf{A}(\mathbf{R}_{uu} + \boldsymbol{\alpha})\mathbf{A}^T] \\
&= E_{\mathbf{A}}[\mathbf{A}\boldsymbol{\Psi}\mathbf{A}^T] = \bar{\mathbf{A}}\boldsymbol{\Psi}\bar{\mathbf{A}}^T + L\Lambda^{-1}. \tag{5.184}
\end{aligned}$$

Using the equation above and the relationship $\bar{\mathbf{A}}\boldsymbol{\Psi} = (\mathbf{R}_{yu} - \mathbf{L}\mathbf{R}_{su})$, and setting the right-hand side of Eq. (5.183) equal to zero, we derive the update equation for Λ , such that,

$$\Lambda^{-1} = \frac{1}{K} \text{diag}[\mathbf{R}_{yy} - \mathbf{R}_{ys} \mathbf{L}^T - \mathbf{L} \mathbf{R}_{sy} + \mathbf{L} \mathbf{R}_{ss} \mathbf{L}^T - \bar{\mathbf{A}}\boldsymbol{\Psi}\bar{\mathbf{A}}^T].$$

The equation above is equal to Eq. (5.166).

References

1. S.S. Nagarajan, H.T. Attias, K.E. Hild, K. Sekihara, A probabilistic algorithm for robust interference suppression in bioelectromagnetic sensor data. *Stat. Med.* **26**, 3886–3910 (2007)
2. J.M. Zumer, H.T. Attias, K. Sekihara, S.S. Nagarajan, A probabilistic algorithm integrating source localization and noise suppression for MEG and EEG data. *NeuroImage* **37**, 102–115 (2007)
3. H. Attias, Inferring parameters and structure of latent variable models by variational bayes, in *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence*, (Morgan Kaufmann Publishers Inc, 1999), pp. 21–30
4. A.P. Dempster, N.M. Laird, D.B. Rubin, Maximum likelihood from incomplete data via the EM algorithm. *J. R. Stat. Soc. Ser. B (Methodological)* **39**, 1–38 (1977)
5. R.H. Shumway, D.S. Stoffer, An approach to time series smoothing and forecasting using the EM algorithm. *J. Time Ser. Anal.* **3**, 253–263 (1982)
6. H.T. Attias, A variational Bayesian framework for graphical models, in *Advances in Neural Information Processing* (MIT Press, Cambridge, 2000), pp. 209–215

Chapter 6

A Unified Bayesian Framework for MEG/EEG Source Imaging

6.1 Introduction

Magnetoencephalography (MEG) and related electroencephalography (EEG) use an array of sensors to take electromagnetic field (or voltage) measurements from on or near the scalp surface with excellent temporal resolution. In both MEG and EEG, the observed field can in many cases be explained by synchronous, compact current sources located within the brain. Although useful for research and clinical purposes, accurately determining the spatial distribution of these unknown sources is a challenging inverse problem. The relevant estimation problem can be posed as follows: The measured electromagnetic signal is $\mathbf{y} \in \mathbb{R}^{d_y \times d_t}$, where d_y equals the number of sensors and d_t is the number of time points at which measurements are made. Each unknown source $s_r \in \mathbb{R}^{d_c \times d_t}$ is a d_c -dimensional neural current dipole, at d_t timepoints, projecting from the r th (discretized) voxel or candidate location distributed throughout the brain. These candidate locations can be obtained by segmenting a structural MR scan of a human subject and tessellating the brain volume with a set of vertices, \mathbf{y} and each s_r are related by the likelihood model

$$\mathbf{y} = \sum_{r=1}^{d_s} \mathbf{L}_r s_r + \boldsymbol{\varepsilon}, \quad (6.1)$$

where d_s is the number of voxels under consideration and $\mathbf{L}_r \in \mathbb{R}^{d_y \times d_c}$ is the so-called lead-field matrix for the r th voxel. The k th column of \mathbf{L}_r represents the signal vector that would be observed at the scalp given a unit current source/dipole at the r th vertex with orientation in the k th direction. It is common to assume $d_c = 2$ (for MEG with a single spherical shell model) or $d_c = 3$ (for EEG), which allows flexible source orientations to be estimated in 2D or 3D space. Multiple methods based on the physical properties of the brain and Maxwell's equations are available for the computation of each L_i , as described in Appendix A. Finally, $\boldsymbol{\varepsilon}$ is a noise-plus-interference term where we assume, for simplicity, that the columns

are drawn independently from $\mathcal{N}(0, \Sigma_\epsilon)$. However, temporal correlations can easily be incorporated if desired using a simple transformation outlined in [1] or using the spatio-temporal framework introduced in [2]. In this chapter, we will mostly assume that Σ_ϵ is known; however, robust procedures for its estimation can be found in Chap. 5, and can naturally be incorporated into the proposed model. However, it should be noted that joint estimation of brain source activity and Σ_ϵ potentially leads to identifiability issues, and remains an unsolved problem to date.

To obtain reasonable spatial resolution, the number of candidate source locations will necessarily be much larger than the number of sensors ($d_s \gg d_y$). The salient inverse problem then becomes the ill-posed estimation of regions with significant brain activity, which are reflected by voxels i such that $|s_i| > 0$; we refer to these as *active* dipoles or sources. The severe underdetermine nature of this MEG (or related EEG) source localization problem (since the mapping from source activity configuration $s \triangleq [s_1^T, \dots, s_{d_s}^T]^T$ to sensor measurement y is many to one), requires the incorporation of prior assumptions when choosing an appropriate solution out of an infinite set of candidates.

Bayesian approaches are useful in this capacity because they allow these assumptions to be explicitly quantified using postulated prior distributions. However, the means by which these priors are chosen, as well as the estimation and inference procedures that are subsequently adopted to affect localization, have led to a daunting array of algorithms with seemingly very different properties and assumptions.

While seemingly quite different in many respects, we present a generalized framework that encompasses all of these methods and points to intimate connections between algorithms. The underlying motivation here is to leverage analytical tools and ideas from machine learning, Bayesian inference, and convex analysis that have not as of yet been fully exploited in the context of MEG/EEG source localization. Specifically, here we address how a simple Gaussian scale mixture prior with flexible covariance components underlie and generalize all of the above. This process demonstrates a number of surprising similarities or out-right equivalences between what might otherwise appear to be very different methodologies. Theoretical properties related to convergence, global and local minima, and localization bias are analyzed and fast algorithms are derived that improve upon existing methods. This perspective leads to explicit connections between many established algorithms and suggests natural extensions for handling unknown dipole orientations, extended source configurations, correlated sources, temporal smoothness, and computational expediency. Specific imaging methods elucidated under this paradigm include weighted minimum ℓ_2 -norm, FOCUSS [3], MCE [4], VESTAL [5], sLORETA [6], ReML [7] and covariance component estimation, beamforming, variational Bayes, the Laplace approximation, and automatic relevance determination (ARD) [8, 9]. Perhaps surprisingly, all of these methods can be formulated as particular cases of covariance component estimation using different concave regularization terms and optimization rules, making general theoretical analyses and algorithmic extensions/improvements particularly relevant. By providing a unifying theoretical perspective and comprehensive analyses, neuroelectromagnetic imaging practitioners will be better able to

assess the relative strengths of many Bayesian strategies with respect to particular applications; it will also help ensure that different methods are used to their full potential and not underutilized.

6.2 Bayesian Modeling Framework

In a Bayesian framework all prior assumptions are embedded in the distribution $p(s)$. If under a given experimental or clinical paradigm this $p(s)$ were somehow known exactly, then the posterior distribution $p(s|y)$ can be computed via Bayes rule:

$$p(s|y) = \frac{p(y|s)p(s)}{p(y)}. \quad (6.2)$$

This distribution contains all possible information about the unknown s conditioned on the observed data y . Two fundamental problems prevent using $p(s|y)$ for source localization. First, for most priors $p(s)$, the distribution $p(y)$ given by:

$$p(y) = \int p(y|s)p(s)ds. \quad (6.3)$$

cannot be computed. Because this quantity, which is sometimes referred to as the model evidence or marginal likelihood, is required to compute posterior moments and is also sometimes used to facilitate model selection, this deficiency can be very problematic. Of course if only a point estimate for s is desired, then this normalizing distribution may not be needed. For example, a popular estimator involves finding the value of s that maximizes the posterior distribution, often called the maximum a posteriori or MAP estimate, and is invariant to $p(y)$. However MAP estimates may be unrepresentative of posterior mass and are unfortunately intractable to compute for most $p(s)$ given reasonable computational resources. Secondly, we do not actually know the prior $p(s)$ and so some appropriate distribution must be assumed, perhaps based on neurophysiological constraints or computational considerations. In fact, it is this choice, whether implicitly or explicitly, that differentiates a wide variety of localization methods at a very high level.

Such a prior is often considered to be fixed and known, as in the case of minimum ℓ_2 -norm approaches [2], minimum current estimation (MCE), FOCUSS, sLORETA, and minimum variance beamformers. Alternatively, a number of empirical Bayesian approaches have been proposed that attempt a form of model selection by using the data to guide the search for an appropriate $p(s)$. In this scenario, candidate priors are distinguished by a set of flexible hyperparameters γ that must be estimated via a variety of data-driven iterative procedures including hierarchical covariance component models, automatic relevance determination (ARD), and several related variational Bayesian methods (for a more comprehensive list of citations see Wipf et al. [10]).

6.3 Bayesian Modeling Using General Gaussian Scale Mixtures and Arbitrary Covariance Components

In this section, we present a general-purpose Bayesian framework for source localization and discuss a central distinction between fixed-prior MAP estimation schemes and empirical Bayesian approaches that adopt a flexible, parameterized prior. While often derived using different assumptions and methodology, they can be related via a simple hierarchical structure based on general Gaussian scale mixture distributions with arbitrary covariance components. Numerous special cases of this model have been considered previously in the context of MEG and EEG source localization and related problems as will be discussed in subsequent sections.

6.3.1 The Generative Model

To begin we invoke the noise model from (6.1), which fully defines the assumed data likelihood

$$p(\mathbf{y}|\mathbf{s}) \propto \exp\left(-\frac{1}{2} \left\| \mathbf{y} - \sum_{i=1}^{d_s} \mathbf{L}_i \mathbf{s}_i \right\|_{\Sigma_\epsilon^{-1}}^2\right), \quad (6.4)$$

where $\|\mathbf{X}\|_W^2$ denotes the weighted matrix norm $\sqrt{\text{trace}[\mathbf{X}^T \mathbf{W} \mathbf{X}]}$.

While the unknown noise covariance can also be parameterized and seamlessly estimated from the data via the proposed paradigm, for simplicity we assume that Σ_ϵ is known, estimated from the data using a variational Bayesian factor analysis (VBFA) model as discussed in Sect. 5.3 and that it is fixed. Next we adopt the following source prior for \mathbf{s} :

$$p(\mathbf{s}|\gamma) \propto \exp\left(-\frac{1}{2} \text{trace}[\mathbf{s}^T \Sigma_s^{-1} \mathbf{s}]\right), \quad \Sigma_s = \sum_{i=1}^{d_\gamma} \gamma_i C_i. \quad (6.5)$$

This is equivalent to applying independently, at each time point, a zero-mean Gaussian distribution with covariance Σ_s to each column of \mathbf{s} . Here $\gamma \triangleq [\gamma_1, \dots, \gamma_{d_\gamma}]$ is a vector of d_γ nonnegative hyperparameters that control the relative contribution of each covariance basis matrix C_i . While the hyperparameters are unknown, the set of components $C \triangleq \{C_i : i = 1, \dots, d_\gamma\}$ is assumed to be fixed and known. Such a formulation is extremely flexible however, because a rich variety of candidate covariance bases can be proposed as will be discussed in more detail later. Moreover, this structure has been advocated by a number of others in the context of neuroelectromagnetic source imaging [7, 11]. Finally, we assume a hyperprior on γ of the form

$$p(\gamma) = \prod_{i=1}^{d_\gamma} \frac{1}{2} \exp[-f_i(\gamma_i)] \quad (6.6)$$

where each $f_i(\cdot)$ is an unspecified function that is assumed to be known. The implicit prior on \mathbf{s} , obtained by integrating out (marginalizing) the unknown γ , is known as a Gaussian scale mixture.

$$p(\mathbf{s}) = \int p(\mathbf{s}|\gamma)p(\gamma)d\gamma. \quad (6.7)$$

6.3.2 Estimation and Inference

Estimation and inference using the proposed model can be carried out in multiple ways depending how the unknown quantities \mathbf{s} and γ are handled. This leads to a natural partitioning of a variety of inverse methods. We briefly summarize three possibilities before discussing the details and close interrelationships.

6.3.2.1 Hyperparameter MAP or Empirical Bayesian

The first option is if γ were somehow known (and assuming Σ_s is known as well), then the conditional distribution $p(\mathbf{s}|\mathbf{y}, \gamma) \propto p(\mathbf{y}|\mathbf{s})p(\mathbf{s}|\gamma)$ is a fully specified Gaussian distribution with mean and covariance given by

$$\begin{aligned} E_{p(\mathbf{s}|\mathbf{y}, \gamma)}[\mathbf{s}] &= \gamma \mathbf{L}^T \left(\Sigma_\epsilon + \mathbf{L} \Sigma_s \mathbf{L}^T \right)^{-1} \mathbf{B} \\ \text{Cov}_{p(s_j|\mathbf{y}, \gamma)}[s_j] &= \Sigma_s - \Sigma_s \mathbf{L}^T \left(\Sigma_\epsilon + \mathbf{L} \Sigma_s \mathbf{L}^T \right)^{-1} \mathbf{L} \Sigma_s, \quad \forall j, \end{aligned} \quad (6.8)$$

where $\mathbf{L} \triangleq [\mathbf{L}_1, \dots, \mathbf{L}_{d_s}]$ and s_j denotes the j th column of \mathbf{s} (i.e., the sources at the j th time point) and individual columns are uncorrelated.

It is then common to use the simple estimator $\hat{\mathbf{s}} = E_{p(\mathbf{s}|\mathbf{y}, \gamma)}[\mathbf{s}]$ for the unknown sources. However, since γ is actually not known, a suitable approximation $\hat{\gamma} \approx \gamma$ must first be found. One principled way to accomplish this is to integrate out the sources \mathbf{s} and then solve

$$\hat{\gamma} = \arg \max_{\gamma} \int p(\mathbf{y}|\mathbf{s})p(\mathbf{s}|\gamma)p(\gamma)d\mathbf{s} \quad (6.9)$$

This treatment is sometimes referred to as empirical Bayes because the γ -dependent prior on \mathbf{s} , $p(\mathbf{s}|\gamma)$, is empirically learned from the data, often using expectation-maximization (EM) algorithms which treat \mathbf{s} as hidden data. Additionally, the process of marginalization provides a natural regularizing mechanism that can shrink many elements of γ , to exactly zero, in effect pruning the associated covariance component from the model, with only the relevant components remaining. Consequently, estimation under this model is sometimes called automatic relevance determination (ARD). This procedure can also be leveraged to obtain a rigorous lower bound on

$\log p(\mathbf{y})$. While knowing $p(s|\mathbf{y})$ is useful for source estimation given a particular model, access to $p(\mathbf{y})$ (or equivalently $\log p(\mathbf{y})$) can assist model selection.

γ -MAP obtains a point estimate for the unknown γ by first integrating out the unknown sources s producing the hyperparameter likelihood equation

$$p(\mathbf{y}|\gamma) = \int p(\mathbf{y}|s)p(s|\gamma)ds \propto \exp\left(-\frac{1}{2}\mathbf{y}^T \Sigma_y^{-1} \mathbf{y}\right), \quad (6.10)$$

where

$$\Sigma_y = \Sigma_\epsilon + \mathbf{L} \Sigma_s \mathbf{L}^T. \quad (6.11)$$

To estimate γ we then solve:

$$\hat{\gamma} = \arg \max_{\gamma} p(\gamma|\mathbf{y}) = \arg \max_{\gamma} p(\mathbf{y}|\gamma)p(\gamma), \quad (6.12)$$

which is equivalent to minimizing the cost function

$$\mathcal{L}(\gamma) \triangleq -2 \log p(\mathbf{y}|\gamma)p(\gamma) \equiv \text{trace}[C_y \Sigma_y^{-1}] + \log |\Sigma_y| + \frac{1}{n} \left[\sum_{i=1}^n f_i(\gamma_i) \right], \quad (6.13)$$

where $C_y \triangleq n^{-1} \mathbf{y} \mathbf{y}^T$ is the empirical covariance. This cost function is composed of three parts. The first is a data fit term based on the dissimilarity between the empirical covariance and the model covariance; in general, this factor encourages γ to be large. The second term provides the primary regularizing or sparsifying effect, penalizing a measure of the volume formed by the model covariance. Since the volume of any high dimensional space is more effectively reduced by collapsing individual dimensions as close to zero as possible (as opposed to reducing all dimensions isometrically), this penalty term promotes a model covariance that is maximally degenerate (or non-spherical), which pushes elements of γ to exactly zero.

Finally, the third term follows directly from the hyperprior, which we have thus far assumed to be arbitrary. This term can be useful for incorporating specific prior information, perhaps from fMRI data or some other imaging modality. It can also be used to expedite hyperparameter pruning or conversely, to soften the pruning process or prevent pruning altogether. One popular choice is the inverse-Gamma hyperprior [12] given by

$$p_i(\gamma_i^{-1}) = \text{Gam}(a_i)^{-1} b_i^{a_i} \gamma_i^{1-a_i} \exp\left(-\frac{b_i}{\gamma_i}\right) \quad (6.14)$$

for the i th hyper parameter, where $\text{Gam}(x) = \int_0^{\inf} z^{x-1} \exp(-z) dz$ is the standard Gamma function and a_i, b_i are the shape parameters. In the limit as $a_i, b_i \rightarrow 0$ then the prior converges to a non-informative Jeffreys prior, when optimization is performed on $\log \gamma_i$ space as is customary in some applications, the effective prior is flat and can therefore be ignored. In contrast, $a_i, b_i \rightarrow \infty$, all hyperparameters

are constrained to have equal value and essentially no learning (or pruning) occurs. Consequently, the standard weighted minimum ℓ_2 -norm solution can be seen as a special case.

We will often concern ourselves with flat hyperpriors when considering the γ -MAP option. In this case, the third term in Eq. (6.13) vanishes and the only regularization will come from the $\log |\cdot|$ term. In this context, the optimization problem with respect to the unknown hyperparameters is sometimes referred to as type-II maximum likelihood. It is also equivalent to the restricted maximum likelihood (ReML) cost function, discussed by Friston et al. [7] Regardless of how γ is optimized, once some $\hat{\gamma}$ is obtained, we compute $\hat{\Sigma}_s$ which fully specifies our assumed empirical prior on s . To the extent that the “learned” prior $p(s|\hat{\gamma})$ is realistic, this posterior quantifies regions of significant current density and point estimates for the unknown sources can be obtained by evaluating the posterior mean.

6.3.2.2 Optimization of γ -MAP

The primary objective of this section is to minimize Eq. (6.13) with respect to γ . For simplicity, we will first present updates with $f_i(\gamma_i) = 0$ (i.e. a flat hyper prior). We then address natural adaptations to the more general cases (e.g., not just conjugate priors). Of course one option is to treat the problem as a general nonlinear optimization task and perform gradient descent or some other generic procedure. In contrast, here we will focus on methods specifically tailored for minimizing Eq. (6.13) using principled methodology. We begin with methods based directly on the EM algorithm and then diverge to alternatives that draw on convex analysis to achieve faster convergence. One approach to minimizing Eq. (6.13) is the restricted likelihood method (ReML) which utilizes what amounts to EM-based updates treating s as hidden data.

For the E-step, the mean and covariance of s are computed given some estimate of the hyperparameters $\hat{\gamma}$. For the M-step, we then must update $\hat{\gamma}$ using these moments as the true values. Unfortunately, the optimal value of $\hat{\gamma}$ cannot be obtained in closed form for arbitrary covariance component sets, so a second-order Fisher scoring procedure is adopted to approximate the desired solution. While effective for estimating small numbers of hyperparameters, this approach requires inverting a $d_\gamma \times d_\gamma$ Fisher information matrix, which is not computationally feasible for large d_γ . Moreover, unlike exact EM implementations, there is no guarantee that such a Fisher scoring method will decrease the likelihood function Eq. (6.13) at each iteration.

Consequently, here we present alternative optimization procedures that apply to the arbitrary covariance model discussed above, and naturally guarantee that $\gamma_i \geq 0$ for all i . All of these methods rely on reparameterizing the generative model such that the implicit M-step can be solved in closed form. First, we note that $\mathcal{L}(\gamma)$ only depends on the data y through the $d_y \times d_y$ sample correlation matrix C_y . Therefore, to reduce the computational burden, we replace y with a matrix $\tilde{y} \in \mathbb{R}^{d_y \times \text{rank}(y)}$ such that $yy^T = \tilde{y}\tilde{y}^T$. This removes any per-iteration dependency on n , which can potentially be large, without altering that actual cost function. It also implies that, for purposes of computing γ , the number of columns of s is reduced to match $\text{rank}(y)$.

Next we introduce the decomposition

$$\mathbf{s} = \sum_{i=1}^{d_\gamma} A_i \tilde{\mathbf{s}}_i = A \tilde{\mathbf{s}}, \quad (6.15)$$

where each A_i is selected such that $A_i A_i^T = C_i$ and $A \triangleq [A_1 \dots A_{d_\gamma}]$, $\mathbf{s} \triangleq [\mathbf{s}_1^T \dots \mathbf{s}_{d_\gamma}^T]^T$. Letting $\tilde{\mathbf{L}} \triangleq [\tilde{\mathbf{L}}_1, \dots, \tilde{\mathbf{L}}_{d_\gamma}] = \mathbf{L}[A_1, \dots A_{d_\gamma}]$, this allows us to re-express the original hierarchical Bayesian model as

$$p(\tilde{\mathbf{y}}|\tilde{\mathbf{s}}) \propto \exp\left(-\frac{1}{2}\|\mathbf{y} - \tilde{\mathbf{L}}\tilde{\mathbf{s}}\|_{\Sigma_\epsilon^{-1}}^2\right) \quad (6.16)$$

$$p(\tilde{\mathbf{s}}_i|\gamma_i) \propto \exp\left(-\frac{1}{2\gamma_i}\|\tilde{\mathbf{s}}_i\|_{\mathcal{F}}^2\right), \quad \forall i = 1, \dots, d_\gamma, \quad (6.17)$$

where $\|\mathbf{X}\|_{\mathcal{F}}$ is the standard Frobenius norm $\sqrt{\text{trace}[\mathbf{X}\mathbf{X}^T]}$. The hyperprior remains unaltered. It is easily verified by the rules for transformation of random variables that Eq. (6.16) and the original model are consistent. It also follows that

$$\Sigma_y = \Sigma_\epsilon + \mathbf{L}\left(\sum_{i=1}^{d_\gamma} \gamma_i C_i\right)\mathbf{L}^T = \Sigma_\epsilon + \sum_{i=1}^{d_\gamma} \gamma_i \tilde{\mathbf{L}}_i \tilde{\mathbf{L}}_i^T = \Sigma_\epsilon + \tilde{\mathbf{L}} \Sigma_{\tilde{\mathbf{s}}} \tilde{\mathbf{L}}^T \quad (6.18)$$

where $\Sigma_{\tilde{\mathbf{s}}}$ is the diagonal, γ -dependent prior covariance of the pseudo sources.

There are three ways to derive update rules for γ : EM iteration, Fixed-point (MacKay) updates, and Convexity based updates. Although the EM iteration update rules are the most straightforward to derive, they have empirically slow convergence rates. In contrast, the fixed-point updates have fast convergence rates but no convergence guarantees. In contrast, the convexity based updates have both fast convergence properties as well as guaranteed convergence properties. Details of derivations and properties of these update rules can be found elsewhere, but the three types of update rules are listed here for completeness.

1. EM-Updates

$$\gamma_i^{(k+1)} \rightarrow \frac{1}{nr_i} \|\gamma_i^{(k)} \tilde{\mathbf{L}}_i^T (\Sigma_y^{(k)})^{-1} \tilde{\mathbf{y}}\|_{\mathcal{F}}^2 + \frac{1}{r_i} \text{trace}\left[\gamma_i^{(k)} I - \gamma_i^{(k)} \tilde{\mathbf{L}}_i^T (\Sigma_y^{(k)})^{-1} \tilde{\mathbf{L}}_i \gamma_i^{(k)}\right] \quad (6.19)$$

2. MacKay update:

$$\gamma_i^{(k+1)} \rightarrow \frac{1}{n} \|\gamma_i^{(k)} \tilde{\mathbf{L}}_i^T (\Sigma_y^{(k)})^{-1} \tilde{\mathbf{y}}\|_{\mathcal{F}}^2 (\text{trace}[\gamma_i^{(k)} \tilde{\mathbf{L}}_i^T (\Sigma_y^{(k)})^{-1} \tilde{\mathbf{L}}_i])^{-1} \quad (6.20)$$

3. Convexity-based update

$$\gamma_i^{(k+1)} \rightarrow \frac{\gamma_i^{(k)}}{n} \|\tilde{\mathbf{L}}_i (\Sigma_y^{(k)})^{-1} \tilde{\mathbf{y}}\|_{\mathcal{F}} (\text{trace}[\tilde{\mathbf{L}}_i^T (\Sigma_y^{(k)})^{-1} \tilde{\mathbf{L}}_i])^{-1/2} \quad (6.21)$$

6.3.2.3 Analysis of γ -MAP

Previously, we have claimed that the γ -MAP process naturally forces excessive/irrelevant hyperparameters to converge to zero, thereby reducing model complexity. Note that, somewhat counterintuitively, this occurs even when a flat hyperprior is assumed. While this observation has been verified empirically by ourselves and others in various application settings, there has been relatively little corroborating theoretical evidence, largely because of the difficulty in analyzing the potentially multimodal, non-convex γ -MAP cost-function. We can then show that: Every local minimum of the generalized γ -MAP cost function, is achieved at a solution with utmost $\text{rank}(\mathbf{y})d_y \leq d_y^2$ non-zero hyper parameters if $f_i(\gamma_i)$ is concave and non-decreasing for all i , including flat hyper priors. Therefore, we can be confident that the pruning mechanism of γ -MAP is not merely an empirical phenomena. Nor is it dependent on a particular sparse hyperprior, the result holds when a flat (uniform) hyperprior is assumed.

The number of observation vectors n also plays an important role in shaping γ -MAP solutions. Increasing n has two primary benefits: (i) it facilitates convergence to the global minimum (as opposed to getting stuck in a suboptimal extrema) and (ii), it improves the quality of this minimum by mitigating the effects of noise. Finally, a third benefit to using $n > 1$ is that it leads to temporal smoothing of estimated time courses (i.e., rows of $\hat{\mathbf{s}}$). This occurs because the selected covariance components do not change across time, as would be the case if a separate set of hyperparameters were estimated at each time point. For purposes of model selection, a rigorous bound on $\log p(\mathbf{y})$ can be derived using principles from convex analysis that have been successfully applied in general-purpose probabilistic graphical models (see Wipf and Nagarajan [13]).

6.3.3 Source MAP or Penalized Likelihood Methods

The second option is to integrate out the unknown γ , we can treat $p(s)$ as the effective prior and attempt to compute a MAP estimate of s via

$$\hat{s} = \arg \max_s \int p(\mathbf{y}|s)p(s|\gamma)p(\gamma)d\gamma = \arg \max_s p(\mathbf{y}|s)p(s) \quad (6.22)$$

While it may not be immediately transparent, solving s -MAP also leads to a shrinking and pruning of superfluous covariance components. In short, this occurs because the hierarchical model upon which it is based leads to a convenient, iterative EM algorithm-based implementation, which treats the hyperparameters γ as hidden data and computes their expectation for the E-step. Over the course of learning, this expectation collapses to zero for many of the irrelevant hyperparameters, removing them from the model in much the same way as γ -MAP.

In a general setting Eq. (6.22) can be a difficult optimization problem and furthermore, the nature of the underlying cost function is not immediately transparent. Consequently, we advocate an indirect alternative utilizing the pseudo-source decomposition given by \tilde{s} described previously, which leads to an efficient EM implementation and a readily interpretable cost function. It also demonstrates that both FOCUSS and MCE can be viewed as EM algorithms that are readily generalized to handle more complex spatio-temporal constraints. Explicitly, we will minimize

$$\begin{aligned}\mathcal{L}(\tilde{s}) &\triangleq -2 \log \int p(\tilde{\mathbf{y}}|\tilde{s}) p(\tilde{s}|\gamma) p(\gamma) d\gamma \\ &= -2 \log p(\tilde{\mathbf{y}}|\tilde{s}) p(\tilde{s}) \\ &\equiv \|\tilde{\mathbf{y}} - \tilde{\mathbf{L}}\tilde{s}\|_{\Sigma_\epsilon^{-1}}^2 + \sum_{i=1}^{d_\gamma} g_i(\|\tilde{s}\|_{\mathcal{F}}^2),\end{aligned}\quad (6.23)$$

where $g_i(\cdot)$ is defined as

$$g_i(\|\tilde{s}\|_{\mathcal{F}}^2) \triangleq -2 \log \int p(\tilde{s}_i|\gamma_i) p_i(\gamma_i) d\gamma_i. \quad (6.24)$$

For many choices of the hyperprior, the associated $g_i(\cdot)$ may not be available in closed form. Moreover, it is often more convenient and transparent to directly assume the form of $g_i(\cdot)$ rather than infer its value from some postulated hyperprior. Virtually any non-decreasing, concave function $g_i(\cdot)$ of interest can be generated by the proposed hierarchical model. In other words, there will always exist some $p_i(\gamma_i)$, possibly improper, such that the stated Gaussian mixture representation will produce any desired concave $g_i(\cdot)$. For example, a generalized version of MCE and FOCUSS can be produced from the selection $g_i(z) = c_i z^{p/2}$, which is concave and amenable to a Gaussian scale-mixture representation for any $p \in (0, 2]$ and constant $c_i > 0$.

Presumably, there are a variety of ways to optimize Eq. (6.23). One particularly straightforward and convenient method exploits the hierarchical structure inherent in the assumed Bayesian model. This leads to simple and efficient EM-based update rules. It also demonstrates that the canonical FOCUSS iterations are equivalent to principled EM updates. Likewise, regularized MCE solutions can also be obtained in the same manner.

Ultimately, we would like to estimate each \tilde{s}_i , which in turn gives us the true sources S . If we knew the values of the hyperparameters γ this would be straightforward; however, these are of course unknown. Consequently, in the EM framework, γ is treated as hidden data whose distribution (or relevant expectation) is computed during the E-step. The M-step then computes the MAP estimate of \tilde{s} assuming that γ equals the appropriate expectation. For the $(k+1)$ th E-step, the expected value of each γ_i^{-1} under the distribution $p(\gamma|\mathbf{y}, \tilde{s}^{(k)})$ is required (see the M-step below) and can be computed analytically assuming $g_i(\cdot)$ is differentiable, regardless of the underlying form of $p(\gamma)$. Assuming $g_i(z) \propto z^{p/2}$, it can be shown that

$$\gamma_i^{(k+1)} \triangleq E_{p(\gamma_i | \tilde{\mathbf{y}}, \tilde{\mathbf{s}}_i^{(k)})} [\gamma_i^{-1}]^{-1} = \left[\frac{1}{r_i n} \|\tilde{\mathbf{s}}_i^{(k)}\|_{\mathcal{F}}^2 \right]^{\frac{2-p}{2}} \quad (6.25)$$

The associated M-step is then readily computed using

$$\tilde{\mathbf{s}}_i^{(k+1)} \rightarrow \gamma_i \tilde{\mathbf{L}}_i^T (\Sigma_{\epsilon} + \sum_{i=1}^{d_{\gamma}} \gamma_i \tilde{\mathbf{L}}_i \tilde{\mathbf{L}}_i^T)^{-1} \tilde{\mathbf{y}}. \quad (6.26)$$

$\mathbf{s}^{(k+1)} = A \tilde{\mathbf{s}}^{(k+1)}$ is then the $(k+1)$ th estimate of the source activity. Each iteration of this procedure decreases the cost function of Eq. (6.23) and converges to some fixed point is guaranteed. From a computational stand point, the s -MAP updates are of the same per-iteration complexity as the γ -MAP updates. Roughly speaking, the EM iterations above can be viewed as coordinate descent over a particular auxiliary cost function dependent on both s and γ . There exists a formal duality between s -MAP and γ -MAP procedures that is beyond the scope of this chapter, but details of which are elaborate in [10, 14].

Given $n = 1$, and $A_i = \mathbf{e}_i$, where each \mathbf{e}_i is a standard indexing vector of zeros with a “1” for the i th element, we get,

$$\gamma_i^{(k+1)} \rightarrow [\gamma_i^k \mathbf{L}_i^T (\lambda I + \mathbf{L} \text{ diag}[\gamma^k] \mathbf{L}^T)^{-1} \mathbf{y}]^{(2-p)} \quad (6.27)$$

where \mathbf{L}_i is the i th column of L . We then recover the exact FOCUSS updates when $p \rightarrow 0$ and a FOCUSS-like update for MCE when $p = 1$. Note however, that while previous applications of MCE and FOCUSS to electromagnetic imaging using $n = 1$ require a separate iterative solution to be computed at each time point in isolation, here the entire s can be computed at once with $n > 1$ for about the same computational cost as a single FOCUSS run.

The nature of the EM updates for s -MAP, where an estimate of γ is obtained via the E-step, suggest that this approach is indirectly performing some form of covariance component estimation. But if this is actually the case, it remains unclear exactly what cost function these covariance component estimates are minimizing. This is unlike the case of γ -MAP where it is more explicit. The fundamental difference between s -MAP and γ -MAP lies in the regularization mechanism of the covariance components. Unlike γ -MAP the penalty term in Eq. (6.23), is a separable summation that depends on the value of p to affect hyperparameter pruning; and importantly there is no volume-based penalty. For $p < 1$ the penalty is concave in γ and hence, every local minimum of Eq. (6.23) is achieved at a solution with at most $\text{rank}(\mathbf{y}) d_y \leq d_y^2$ non-zero hyper parameters. Rather than promoting sparsity at the level of individual source elements at a given voxel and time (as occurs with standard MCE and FOCUSS when $n = 1$, $C_i = \mathbf{e}_i$, here sparsity is encouraged at the level of the pseudo-sources, $\tilde{\mathbf{s}}$. The function $g_i(\cdot)$ operates on the Frobenius norm of each $\tilde{\mathbf{S}}_i$ and favors solutions with many $\|\tilde{\mathbf{s}}_i\|_{\mathcal{F}} = 0$ for many indices of i . Notably though, by virtue of the Frobenius norm over time, within a non-zero $\tilde{\mathbf{s}}_i$ temporally smooth (non-sparse)

solutions are favored. This is why, one obtains continuous source estimates over time, as opposed to application of FOCUSS or MCE at each individual time-point.

Even though sparsity bounds for *s-MAP* and γ -*MAP* are somewhat similar, the actual performance in practical situations is quite distinct. This is because of global minimum convergence properties of *s-MAP* that are impacted by the choice of $g_i(\cdot)$ or p . For instance, although setting $p = 1$ leads to a convex cost function devoid of non-global minima, and the update rules for generalized-MCE will converge to the a global minimum, the resultant estimate can have problems recovering the true source estimate, especially for conditions related to MEG and EEG source reconstructions. This is because lead-fields L required for MEG and EEG are highly correlated, and violate the restricted isometric properties (RIP) that are required for accurate performance of ℓ_1 -norm procedures. More details on this point can be found in [14]. These theoretical restrictions essentially render the conditions for MCE performance to reconstruction of 1–2 dipoles at best, with no localization bias guarantees. The problem with MCE is not the existence of local minima, rather, it is that the global minimum may be unrepresentative of the true source distribution even for simple dipolar source configurations. In this situation, and especially when lead field columns are highly correlated, the MCE solution may fail to find sufficiently sparse source representations consistent with the assumption of a few equivalent current dipoles, and this issue will also persist with more complex covariance components and source configurations. Nevertheless, the unimodal, convex nature of the generalized MCE like procedures are its attractive advantage.

Furthermore, if $p < 1$ (which implies that $g_i(z^2)$ is concave in z) then more pseudo sources will be pruned at any global solution, which often implies that those which remain may be more suitable than the MCE estimate. Certainly this is true when estimating dipolar sources, but it likely holds in more general situations as well. However, local minima can be an unfortunate menace with $p < 1$. For instance when $p \rightarrow 0$, we get the the canonical FOCUSS algorithm, and it has a combinatorial number of local minima satisfying

$$d_s - 1C_{dy} + 1 \leq \# \text{ of FOCUSS Local Minima} \leq_{d_s} C_{dy}. \quad (6.28)$$

which is a huge number for practical lead-field matrices, and this property largely explains the sensitivity of FOCUSS to initialization and noise. While the FOCUSS cost function can be shown to have zero localization bias at the global solution, because of the tendency to become stuck at local optima, in practice a bias can be observed when recovering even a single dipolar source. Other selections of p between zero and one can lead to a similar fate. In the general case, a natural trade-off exists with *s-MAP* procedures: greater sparsity of solutions at the global minimum the less possibility that this minimum is biased, but the higher the chance of suboptimal convergence to a biased local minimum, and the optimal balance could be application dependent. Nevertheless, *s-MAP* is capable of successfully handling large numbers of diverse covariance components, and therefore simultaneous constraints on the source space.

6.3.4 Variational Bayesian Approximation

From the perspective of a Bayesian purist however, the pursuit of MAP estimates for unknown quantities of interest, whether parameters s or hyperparameters γ , can be misleading since these estimates discount uncertainty and may not reflect regions of significant probability mass, unlike (for example) the posterior mean. Variational Bayesian methods, which have successfully been applied to a wide variety of hierarchical Bayesian models in the machine learning literature offer an alternative to s -MAP and γ -MAP. Therefore, a third possibility involves finding formal approximations to $p(s|y)$ as well as the marginal $p(y)$ using an intermediary approximation for $p(y)$. However, because of the intractable integrations involved in obtaining either distribution, practical implementation requires additional assumptions leading to different types of approximation strategies. The principle idea here is that all unknown quantities should either be marginalized (integrated out) when possible or approximated with tractable distributions that reflect underlying uncertainty and have computable posterior moments. Practically, we would like to account for ambiguity regarding γ when estimating $p(s|y)$, and potentially, we would like a good approximation for $p(y)$, or a bound on the model evidence $\log p(y)$ for application to model selection. The only meaningful difference between VB and γ -MAP, at least in the context of the proposed generative model, involves approximations to the model evidence $\log p(y)$, with VB and γ -MAP giving different estimates.

In this section, we discuss two types of variational approximations germane to the source localization problem: the mean field approximation (VB-MF), and a fixed-form, Laplace approximation (VB-LA). It turns out that both are related to γ -MAP but with important distinctions. A mean-field approximation makes the simplifying assumption that the joint distribution over unknowns s and γ and factorizes, meaning $p(s, \gamma|y) \approx \hat{p}(s|y)\hat{p}(\gamma|y)$ where $\hat{p}(s|y)$ and $\hat{p}(\gamma|y)$ are chosen to minimize the Kullback-Leibler divergence between the factorized and full posterior. This is accomplished via an iterative process akin to EM, effectively using two E-steps (one for s and one for γ). It also produces a rigorous lower bound on $\log p(y)$ similar to γ -MAP. A second possibility applies a second-order Laplace approximation to the posterior on the hyperparameters (after marginalizing over the sources s), which is then iteratively matched to the true posterior; the result can then be used to approximate $p(s|y)$ and $\log p(y)$. Both of these VB methods lead to posterior approximations $\hat{p}(s|y) = p(s|y, \gamma = \hat{\gamma})$, $\hat{\gamma}$ is equivalently computed via γ -MAP. Consequently, VB has the same level of component pruning as γ -MAP.

6.3.4.1 Mean-Field Approximation

The basic strategy here is to replace intractable posterior distributions with approximate ones that, while greatly simplified and amenable to simple inference procedures, still retain important characteristics of the full model. In the context of our presumed model structure, both the posterior source distribution $p(s|y)$, which is maximized

with s -MAP, as well as the hyperparameter posterior $(\gamma|y)$, which is maximized via γ -MAP, are quite complex and can only be expressed up to some unknown scaling factor (the integration required for normalization is intractable). Likewise, the joint posterior $p(s, \gamma|y)$ over all unknowns is likewise intractable and complex. VB attempts to simplify this situation by finding an approximate joint posterior that factorizes as

$$p(s, \gamma|y) \approx \hat{p}(s, \gamma|y) = \hat{p}(s|y)\hat{p}(\gamma|y), \quad (6.29)$$

where $\hat{p}(s|y)$ and $\hat{p}(\gamma|y)$ are amenable to closed-form computation of posterior quantities such as means and variances (unlike the full posteriors upon which our model is built). This is possible because the enforced factorization, often called the mean-field approximation reflecting its origins in statistical physics, simplifies things significantly. The cost function optimized to find this approximate distribution is

$$\hat{p}(s|y), \hat{p}(\gamma|y) = \underset{q(s), q(\gamma)}{\operatorname{argmin}} \text{KL}[q(s)q(\gamma)||p(s, \gamma|y)], \quad (6.30)$$

where $q(s)$ and $q(\gamma)$ are arbitrary probability distributions and $\text{KL}[\cdot||\cdot]$ indicates the Kullback-Leibler divergence measure.

Recall that γ -MAP iterations effectively compute an approximate distribution for s (E-step) and then a point estimate for γ (M-step); s -MAP does the exact opposite. In contrast, here an approximating distribution is required for both parameters s and hyperparameters γ . While it is often convenient that conjugate hyperpriors must be employed such that Eq. (6.30) is solvable, in fact this problem can be solved by coordinate descent over $q(s)$ and $q(\gamma)$ for virtually any hyperprior. It can be shown that γ -MAP and mean-field approximations can be equivalent in terms of the cost function being optimized and the source activity estimates obtained, given certain choice of hyperpriors [10]. Nevertheless offer a whole class of algorithms within this framework with different covariance component sets, and possible hyperpriors selected, and how the optimization is performed. The main advantage of VB is that strict lower bounds on $\log p(y)$ automatically fall out of the VB framework, given by:

$$\begin{aligned} \log p(y) &\geq F \triangleq \log p(y) - \text{KL}[q(s)q(\gamma)||p(s, \gamma|y)] \\ &= \int \hat{p}(s|y)\hat{p}(\gamma|y) \log \frac{p(y, s, \gamma)}{\hat{p}(s|y)\hat{p}(\gamma|y)} d\gamma, \end{aligned} \quad (6.31)$$

where the inequality follows by the non-negativity of the Kullback-Leibler divergence. The quantity F is sometimes referred to the variational free energy. Evaluation of F requires the full distribution $\hat{p}(\gamma|y)$ and therefore necessitates using conjugate priors or further approximations.

6.3.4.2 Laplace Approximation

The Laplace approximation has been advocated to finding a tractable posterior distribution on the hyperparameters and then using this $\hat{p}(\gamma|y)$ to find approximations

to $p(s|y)$ and $\log p(y)$. To facilitate this process, the hyperparameters are re-parameterized via the transformation $\lambda_i \triangleq \log \gamma_i$, $\forall i$, which now allows them to have positive or negative values. The Laplace approximation then involves the assumption that the posterior distribution of these new hyperparameters $\boldsymbol{\lambda} \triangleq [\lambda_1, \dots, \lambda_{d_\gamma}]^T$ satisfies $p(\boldsymbol{\lambda}|y) \approx \mathcal{N}(\boldsymbol{\lambda}|\mu_\lambda, \Sigma_\lambda)$. There are a variety of ways μ_λ and Σ_λ can be chosen to form the best approximation. For instance, a MAP estimate of $\boldsymbol{\lambda}$ can be computed by maximizing $\log p(y, \boldsymbol{\lambda})$ using a second-order Fisher scoring procedure [1], although this method not guaranteed to increase the cost function and requires a very expensive $O(d_\gamma^3)$ inverse computation. Once the mode $\hat{\boldsymbol{\lambda}}$ is obtained, setting $\mu_\lambda = \hat{\boldsymbol{\lambda}}$ and match the second-order statistics of the true and approximate distributions at this mode using $\Sigma_\lambda = -[\ell''(\mu_\lambda)]^{-1}$ where $\ell(\mu_\lambda) \triangleq \log p(y, \boldsymbol{\lambda} = \hat{\boldsymbol{\lambda}})$ and $\ell''(\mu_\lambda)$ is the corresponding Hessian matrix evaluated at μ_λ . Furthermore, because of the non-negativity of the KL divergence, it always holds that

$$\begin{aligned} \log p(y) &\geq F \triangleq \log p(y) - \text{KL}[\mathcal{N}(\boldsymbol{\lambda}|\mu_\lambda, \Sigma_\lambda)||p(\boldsymbol{\lambda}|y)] \\ &= \int \mathcal{N}(\boldsymbol{\lambda}|\mu_\lambda, \Sigma_\lambda) \log \frac{p(y, \boldsymbol{\lambda})}{\mathcal{N}(\boldsymbol{\lambda}|\mu_\lambda, \Sigma_\lambda)} d\boldsymbol{\lambda}. \end{aligned} \quad (6.32)$$

Unfortunately however, this bound cannot be computed in closed form, so the computable approximation is used as a surrogate instead [1].

In summary, the approximate distributions $\hat{p}(s|y)$ obtained from both variational methods can be shown to be equivalent to a γ -MAP estimate given the appropriate equalizing hyperprior. Consequently, the optimization procedures for γ -MAP can be adapted for VB. Furthermore, the sparsifying properties of the underlying γ -MAP cost functions also apply to VB. Whereas VB-MF gives more principled bounds on the likelihood, the Laplace approximation methods do not. However, some limitations of the Laplace's approximation method can potentially be ameliorated by combining the Laplace and mean-field approximations as outline in Appendix D of [13].

6.4 Selection of Covariance Components C

While details and assumptions may differ γ -MAP, s -MAP, and VB can all be leveraged to obtain a point estimate of S expressed in the Tikhonov regularized form.

$$\hat{S} = \hat{\Sigma}_S \mathbf{L}^T (\Sigma_\epsilon + \mathbf{L} \hat{\Sigma}_S \mathbf{L}^T)^{-1} \mathbf{y} \quad \text{with} \quad \hat{\Sigma}_S = \sum_i \hat{\gamma}_i C_i. \quad (6.33)$$

Importantly, since each method intrinsically sets $\hat{\gamma}_i = 0$ for superfluous covariance components, we can in principle allow the cardinality of C to be very large, and rely on a data-learning process to select the most appropriate subset. Specific choices for C lead to a variety of established algorithms as discussed next.

In the simplest case, the single-component assumption $\Sigma_S = \gamma_1 C_1 = \gamma_1 I$, where I is the Identity matrix, leads to a weighted minimum- ℓ_2 -norm algorithm. More

interesting covariance component terms have been used to effect spatial smoothness, depth bias compensation, and candidate locations of likely activity. With regard to the latter, it has been suggested that prior information about a source location can be coded by including a second term C_2 i.e. $C = C_1C_2$, with all zeros except a patch of 1's along the diagonal signifying a location of probable source activity, perhaps based on fMRI data. For $s\text{-MAP}$, $\gamma\text{-MAP}$, VB , we will obtain a source estimate representing a trade-off (modulated by the relative values of the associated $\hat{\gamma}_1$ and $\hat{\gamma}_2$) between honoring the prior information imposed by C_2 and the smoothness implied by C_1 . The limitation of this proposal is that we generally do not know, a priori, the regions where activity is occurring with both high spatial and temporal resolution. Therefore, we cannot reliably known how to choose an appropriate location-prior term C_2 in many situations. A potential solution to this dilemma is to try out many different (or even all possible) combinations of location priors. For example, if we assume the underlying source currents are formed from a collection of dipolar point sources located at each vertex of the leadfield grid, then we may choose $C = \mathbf{e}_i\mathbf{e}_i^T$ $i = 1, \dots, d_S$, where each \mathbf{e}_i is a standard indexing vector of zeros with a “1” for the i th element (and so $C_i = \mathbf{e}_i\mathbf{e}_i^T$ encodes a prior preference for a dipolar source at location i). This specification for the prior involves the counterintuitive addition of an unknown hyperparameter for every candidate source location which, on casual analysis may seem prone to severe overfitting. As suggested previously however, $s\text{-MAP}$, $\gamma\text{-MAP}$, VB all possess an intrinsic, sparsity-based regularization mechanism. This ameliorates the overfitting problem substantially and effectively reduces the space of possible active source locations by choosing a small relevant subset of active dipolar locations. In general, the methodology is quite flexible and other prior specifications can be included as well, such as temporal and spectral constraints.

In summary, there are two senses with which to understand the notion of covariance component selection. First, there is the selection of which components to include in the model before any estimation takes place, i.e., the choice of C . Second, there is the selection that occurs within C as a natural byproduct of many hyperparameters being driven to zero during the learning process. Such components are necessarily pruned by the model; those that remain have therefore been ‘selected’ in some sense. Here we have argued that in many cases the later, data-driven selection can be used to ease the burden of often ad hoc user-specified selections.

6.5 Discussion

The efficacy of modern Bayesian techniques for quantifying uncertainty and explicitly accounting for prior assumptions make them attractive candidates for source localization. However, it is not always transparent how these methods relate, nor how they can be extended to handle more challenging problems, nor which ones should be expected to perform best in various situations relevant to MEG/EEG source imaging. Starting from a hierarchical Bayesian model constructed using Gaussian scale mixtures with flexible covariance components, we analyze and, where pos-

sible, extend three broad classes of Bayesian inference methods: γ -MAP, which involves integrating out the unknown sources and optimizing the hyperparameters, and s -MAP, which integrates out the hyperparameters and directly optimizes over the sources, and variational approximation methods, which attempt to account for uncertainty in all unknowns. Together, these three encompass a surprisingly wide range of existing source reconstruction approaches, which makes general theoretical analyses and algorithmic extensions/improvements pertaining to the particularly relevant. Thus far, we have attempted to relate and extend three large classes of Bayesian inverse methods, all of which turn out to be performing covariance component estimation/pruning using different sparsity promoting regularization procedures. We now provide some summary points related to connections to existing methods.

1. s -MAP, γ -MAP, and VB can be viewed as procedures for learning a source covariance model using a set of predetermined symmetric, positive semi-definite covariance components. The number of components in this set, each of which acts as a constraint on the source space, can be extremely large, potentially much larger than the number of sensors. However, a natural pruning mechanism effectively discards components that are unsupported by the data. This occurs because of an intrinsic sparsity preference in the Gaussian scale mixture model, which is manifested in an explicit sparsity-inducing regularization term. Consequently, it is not crucial that the user/analyst manually determine an optimal set of components a priori; many components can be included initially allowing the learning process to remove superfluous ones.
2. The wide variety of Bayesian source localization methods that fall under this framework can be differentiated by the following factors: (1) Selection of covariance component regularization term; (2) Choice of initial covariance component set C; (3) Optimization method/ Update rules; and (4) Approximation to $\log p(\mathbf{y})$; this determines whether we are ultimately performing s -MAP, γ -MAP, or VB.
3. Covariance component possibilities include geodesic neural basis functions for estimating distributed sources [34], spatial smoothing factors [24], indicator matrices to couple dipole components or learn flexible orientations [36], fMRI-based factors [31], and temporal and spectral constraints [21].
4. With large numbers of covariance components, s -MAP, γ -MAP, and VB provably remove or prune a certain number of components which are not necessary for representing the observed data.
5. In principle, the noise-plus-interference covariance can be jointly estimated as well, competing with all the other components to model the data. However, identifiability issues can be a concern here and so we consider it wiser to estimate Σ_ϵ via other means (e.g., using VBFA applied to prestimulus data as described in Chap. 5).
6. The latent structure inherent to the Gaussian scale-mixture model leads to an efficient, principled family of update rules for s -MAP, γ -MAP, and VB. This facilitates the estimation of complex covariance structures modulated by very large numbers of hyperparameters (e.g., 10^5+) with relatively little difficulty.

7. Previous focal source imaging techniques such as FOCUSS and MCE display undesirable discontinuities across time as well significant biases in estimating dipole orientations. Consequently, various heuristics have been proposed to address these deficiencies. However, the general spatiotemporal framework of *s-MAP*, γ -*MAP*, and VB handles both of these concerns in a robust, principled fashion by the nature of their underlying cost function. The standard weighted minimum norm can be seen as a limiting case of γ *MAP*.
8. As described in other chapters here, adaptive beamformers are spatial filters that pass source signals in particular focused locations while suppressing interference from elsewhere. The widely-used minimum variance adaptive beamformer (MVAB) creates such filters using a sample covariance estimate; however, the quality of this estimate deteriorates when the sources are correlated or the number of samples n is small. The simpler γ -*MAP* strategy can also be used to enhance beamforming in a way that is particularly robust to source correlations and limited data [15]. Specifically, the estimated γ -*MAP* data covariance matrix $\hat{\Sigma}_y = \Sigma_\epsilon + \sum_i \hat{\gamma}_i \mathbf{L} C_i \mathbf{L}^T$ can be used to replace the problematic sample covariance $C_y = \mathbf{y} \mathbf{y}^T$. This substitution has the natural ability to remove the undesirable effects of correlations or limited data. When n becomes large and assuming uncorrelated sources, this method reduces to the exact MVAB. Additionally, the method can potentially enhance a variety of traditional signal processing methods that rely on robust sample covariance estimates.
9. It can be shown that sLORETA is equivalent to performing a single iteration of a particular γ -*MAP* optimization procedure. Consequently, the latter can be viewed as an iterative refinement of sLORETA. This is exactly analogous to the view of FOCUSS as an iterative refinement of a weighted minimum ℓ_2 -norm estimate.
10. γ -*MAP* and VB have theoretically zero localization bias estimating perfectly uncorrelated dipoles given the appropriate hyperprior and initial set of covariance component.
11. The role of the hyperprior $p(\gamma)$ is heavily dependent on the estimation algorithm being performed. In the *s-MAP* framework, the hyperprior functions through its role in creating the concave regularization function $g_i(\cdot)$. In practice, it is much more transparent to formulate a model directly based on a desired $g_i(\cdot)$ as opposed to working with some supposedly plausible hyperprior $p(\gamma)$ and then inferring the what the associated $g_i(\cdot)$ would be. In contrast, with γ -*MAP* and VB the opposite is true. Choosing a model based on the desirability of some $g_i(\cdot)$ can lead to a model with an underlying hyperprior $p(\gamma)$ that performs poorly. Both VB and γ -*MAP* give rigorous bounds on the model evidence $\log p(\mathbf{y})$.

In summary, we hope that these ideas help to bring an insightful perspective to Bayesian source imaging methods, reduce confusion about how different techniques relate, expand the range of feasible applications of these methods. We have observed a number of surprising similarities or out-right equivalences between what might otherwise appear to be very different methodologies. Additionally, there are numerous promising directions for future research, including time-frequency extensions,

alternative covariance component parameterizations, and robust interference suppression that inspire our own current body of work in this area, and we hope will also inspire other researchers!

References

1. K. Friston, L. Harrison, J. Daunizeau, S. Kiebel, C. Phillips, N. Trujillo-Barreto, R. Henson, G. Flandin, J. Mattout, Multiple sparse priors for the M/EEG inverse problem. *NeuroImage* **39**(3), 1104–1120 (2008)
2. A. Bolstad, B.V. Veen, R. Nowak, Space-time event sparse penalization for magneto-/electroencephalography. *NeuroImage* **46**(4), 1066–1081 (2009)
3. I.F. Gorodnitsky, J.S. George, B.D. Rao, Neuromagnetic source imaging with FOCUSS: a recursive weighted minimum norm algorithm. *Electroencephalogr. Clin. Neurophysiol.* **95**, 231–251 (1995)
4. K. Uutela, M. Hämäläinen, E. Somersalo, Visualization of magnetoencephalographic data using minimum current estimate. *NeuroImage* **10**, 173–180 (1999)
5. M.-X. Huang, A.M. Dale, T. Song, E. Halgren, D.L. Harrington, I. Podgorny, J.M. Canive, S. Lewis, R.R. Lee, Vector-based spatial-temporal minimum l1-norm solution for MEG. *NeuroImage* **31**(3), 1025–1037 (2006)
6. R.D. Pascual-Marqui, Standardized low resolution brain electromagnetic tomography (sLORETA): technical details. *Methods Find. Exp. Clin. Pharmacol.* **24**, 5–12 (2002)
7. K.J. Friston, W. Penny, C. Phillips, S. Kiebel, G. Hinton, J. Ashburner, Classical and Bayesian inference in neuroimaging: theory. *NeuroImage* **16**(2), 465–483 (2002)
8. D.J.C. MacKay, Bayesian interpolation. *Neural Comput.* **4**, 415–447 (1992)
9. R.R. Ramírez, S. Makeig, Neuroelectromagnetic source imaging using multiscale geodesic neural bases and sparse Bayesian learning, in *12th Annual Meeting of the Organization for Human Brain Mapping* (Florence, Italy, 2006)
10. D.P. Wipf, J.P. Owen, H.T. Attias, K. Sekihara, S.S. Nagarajan, Robust Bayesian estimation of the location, orientation, and time course of multiple correlated neural sources using MEG. *NeuroImage* **49**, 641–655 (2010)
11. J. Mattout, C. Phillips, W.D. Penny, M.D. Rugg, K.J. Friston, MEG source localization under multiple constraints: an extended Bayesian framework. *NeuroImage* **30**(3), 753–767 (2006)
12. M.-A. Sato, T. Yoshioka, S. Kajihara, K. Toyama, N. Goda, K. Doya, M. Kawato, Hierarchical Bayesian estimation for MEG inverse problem. *NeuroImage* **23**(3), 806–826 (2004)
13. D. Wipf, S. Nagarajan, A unified Bayesian framework for MEG/EEG source imaging. *NeuroImage* **44**(3), 947–966 (2009)
14. D.P. Wipf, B.D. Rao, S. Nagarajan, Latent variable Bayesian models for promoting sparsity. *IEEE Trans. Inf. Theory* **57**(9), 6236–6255 (2011)
15. D. Wipf, S. Nagarajan, Beamforming using the relevance vector machine, in *Proceedings of the 24th International Conference on Machine Learning* (ACM, 2007), pp. 1023–1030

Chapter 7

Source-Space Connectivity Analysis Using Imaginary Coherence

7.1 Introduction

There has been tremendous interest in estimating the functional connectivity of neuronal activities across different brain regions using electromagnetic brain imaging. Functional connectivity analysis has traditionally been implemented in the sensor space, but lately, a number of studies have begun to use source-space analysis, in which voxel time courses are first estimated by an inverse algorithm, and brain interactions are then analyzed using those estimated voxel time courses [1–3]. Although a certain degree of inaccuracy exists in the source estimation process, the source-space analysis has the potential of providing more accurate information regarding which brain regions are functionally coupled.

Source-space connectivity analysis computes a metric of brain interaction called a connectivity metric, using the voxel time courses. Among existing connectivity metrics, a representative metric is the coherence [2–5]. However, in the source-space coherence analysis, a serious problem arises from spurious coherence caused by the leakage of an inverse algorithm. Such leakages are more or less inevitable in all types of inverse algorithms [6]. One representative ramification of this spurious coherence is an artifactual large peak around the seed voxel, called seed blur, in the resulting coherence image. To remove such spurious coherence, the use of the imaginary part of coherence, which is called the imaginary coherence, has been proven to be effective [7]. It should be mentioned that the use of imaginary coherence was originally proposed by Nolte et al. [8] to remove the spurious coherence caused by the volume conduction in EEG sensor-space analysis.

This chapter reviews the source-space imaginary coherence analysis and related methods. We first provide a detailed theoretical analysis on how the use of imaginary coherence leads to the removal of the spurious coherence caused by the algorithm leakage. We then discuss several related methods of imaginary coherence, including corrected imaginary coherence, canonical coherence, and envelope-to-envelope correlation/coherence. We present numerical examples that confirm our arguments.

7.2 Source-Space Coherence Imaging

In the source-space analysis, the first step estimates voxel time courses using an inverse algorithm. Since the source is a three-dimensional vector with three (x , y , and z) components,¹ most source reconstruction algorithms produce component-wise, multiple time courses at each voxel. Accordingly, we should compute a “representative” single time course by projecting the multiple time courses onto the direction of the source orientation.

However, in practical applications, the source orientation is generally unknown, and it must be estimated from the data. One quick and easy way to estimate the source orientation is to use the direction that maximizes the reconstructed voxel power. The reconstructed source-time courses at the j th voxel is expressed as

$$\widehat{\mathbf{S}}_j = \begin{bmatrix} \widehat{s}_x(\mathbf{r}_j, t_1) & \widehat{s}_x(\mathbf{r}_j, t_2) & \cdots & \widehat{s}_x(\mathbf{r}_j, t_K) \\ \widehat{s}_y(\mathbf{r}_j, t_1) & \widehat{s}_y(\mathbf{r}_j, t_2) & \cdots & \widehat{s}_y(\mathbf{r}_j, t_K) \\ \widehat{s}_z(\mathbf{r}_j, t_1) & \widehat{s}_z(\mathbf{r}_j, t_2) & \cdots & \widehat{s}_z(\mathbf{r}_j, t_K) \end{bmatrix}, \quad (7.1)$$

where we assume that the data is collected at time points, t_1, t_2, \dots, t_K , and \mathbf{r}_j indicates the location of the j th voxel. Denoting the orientation of the source at the j th voxel as $\boldsymbol{\eta}_j$, the estimate of $\boldsymbol{\eta}_j$, $\widehat{\boldsymbol{\eta}}_j$, is obtained using the following maximization:

$$\widehat{\boldsymbol{\eta}}_j = \underset{\boldsymbol{\eta}_j}{\operatorname{argmax}} \boldsymbol{\eta}_j^T (\widehat{\mathbf{S}}_j \widehat{\mathbf{S}}_j^T) \boldsymbol{\eta}_j.$$

The optimum estimate is obtained as the eigenvector corresponding to the maximum eigenvalue of a matrix $\widehat{\mathbf{S}}_j \widehat{\mathbf{S}}_j^T$, i.e.,

$$\widehat{\boldsymbol{\eta}}_j = \vartheta_{\max}\{\widehat{\mathbf{S}}_j \widehat{\mathbf{S}}_j^T\}, \quad (7.2)$$

where the notation $\vartheta_{\max}\{\cdot\}$, defined in Sect.C.9, indicates the eigenvector corresponding to the maximum eigenvalue of a matrix between the parentheses. The representative time course at the j th voxel, $u_j(t)$, is obtained using

$$[u_j(t_1), u_j(t_2), \dots, u_j(t_K)] = \widehat{\boldsymbol{\eta}}_j^T \widehat{\mathbf{S}}_j. \quad (7.3)$$

Once a representative voxel time course is obtained at each voxel, the next step is to compute a voxel-pairwise coherence. This step involves first setting a reference point, called the seed point, and computing the coherence between the time course from the seed point and that from another voxel’s location, referred to as the target location. By scanning through all target locations in a brain, a three-dimensional

¹ Although the source has two components when the homogeneous spherical conductor model is used, for arguments in this chapter, we assume that the source vector has three x , y , and z components.

mapping of source coherence—a source coherence image with respect to the seed location—can be obtained. We can further scan not only the target location but also the seed location to obtain a six-dimensional coherence image, called the voxel-to-voxel coherence matrix.

Let us define the spectra of the seed and the target voxels as $\sigma_S(f)$ and $\sigma_T(f)$, respectively. The coherence $\phi(f)$ is obtained by computing the correlation of these spectra,

$$\phi(f) = \frac{\langle \sigma_T(f) \sigma_S^*(f) \rangle}{\sqrt{\langle |\sigma_T(f)|^2 \rangle \langle |\sigma_S(f)|^2 \rangle}}, \quad (7.4)$$

where the superscript * indicates the complex conjugate, and the brackets $\langle \cdot \rangle$ indicate the ensemble average. In practical applications, this ensemble average is computed by averaging across multiple trials. When only a single continuous data set is measured, the single data set is divided into many trials and coherence is obtained by averaging across these trials.

It is apparent in Eq. (7.4) that if the seed and target spectra contain common components that do not result from true brain interactions, then the coherence may contain spurious components. In source coherence imaging, the leakage of the imaging algorithm is a major source of such spurious coherence.

7.3 Real and Imaginary Parts of Coherence

We here take a look at the nature of real and imaginary parts of coherence before proceeding with the arguments of the leakage influence on the coherence imaging. Let us define the time course from the seed voxel as $u_S(t)$ and the time course from the target voxel as $u_T(t)$. The cross correlation of the seed and target time courses, $R(\tau)$, is then defined such that

$$R(\tau) = \langle u_T(t + \tau) u_S(t) \rangle = \int_{-\infty}^{\infty} u_T(t + \tau) u_S(t) dt. \quad (7.5)$$

The cross correlation is related to the cross spectrum density through

$$R(\tau) = \int_{-\infty}^{\infty} \Psi(f) e^{-i2\pi f \tau} df, \quad (7.6)$$

where the cross spectrum is

$$\Psi(f) = \langle \sigma_T(f) \sigma_S^*(f) \rangle. \quad (7.7)$$

According to Eq. (7.6), we have

$$R(0) = \int_{-\infty}^{\infty} \Psi(f) df = \int_{-\infty}^{\infty} \Re[\Psi(f)] df = \int_{-\infty}^{\infty} \Re[\langle \sigma_T(f) \sigma_S^*(f) \rangle] df, \quad (7.8)$$

where $\Re[\cdot]$ indicates the real part of the complex number in the square brackets. The right-hand side of the above equation holds because the real part of $\Psi(f)$ is an even function and the imaginary part is an odd function. (This is because $R(\tau)$ is a real-valued function.) Therefore, for a narrow-band signal,

$$R(0) = \int_{-\infty}^{\infty} \Re[\langle \sigma_T(f) \sigma_S^*(f) \rangle] df \approx \Re[\langle \sigma_T(f) \sigma_S^*(f) \rangle]. \quad (7.9)$$

Here, $R(0)$, which represents the zero-time-lag correlation, is related to the real part of the cross spectrum, and thus, the real part of coherence represents the instantaneous interaction.

Applying Parseval's theorem to the Fourier transform relationship in Eq. (7.6), we get

$$\int_{-\infty}^{\infty} R(\tau)^2 d\tau = \int_{-\infty}^{\infty} |\Psi(f)|^2 df = \int_{-\infty}^{\infty} |\langle \sigma_T(f) \sigma_S^*(f) \rangle|^2 df. \quad (7.10)$$

Thus, for a narrow-band signal, the following relationship holds:

$$\int_{-\infty}^{\infty} R(\tau)^2 d\tau \approx |\langle \sigma_T(f) \sigma_S^*(f) \rangle|^2 = \Re[\langle \sigma_T(f) \sigma_S^*(f) \rangle] + \Im[\langle \sigma_T(f) \sigma_S^*(f) \rangle], \quad (7.11)$$

where $\Im[\cdot]$ indicate the imaginary part of the complex number in the square brackets. Using Eq. (7.9), we have

$$\Im[\langle \sigma_T(f) \sigma_S^*(f) \rangle] \approx \int_{-\infty}^{\infty} R(\tau)^2 d\tau - R(0)^2 = \sum_{\tau \neq 0} R(\tau)^2. \quad (7.12)$$

The imaginary part of the cross spectrum is equal to the sum of nonzero-lag correlations. Therefore, we can see that the imaginary part of coherence represents the non-instantaneous interaction.

7.4 Effects of the Leakage

7.4.1 Leakage Effects on the Magnitude Coherence

We analyze the effects of leakage of imaging algorithms in source coherence imaging. When the leakage is taken into account, the estimated source time courses at the seed and the target voxels, $\hat{u}_S(t)$ and $\hat{u}_T(t)$, can be expressed as

$$\hat{u}_S(t) = u_S(t) + d_1 u_T(t) + c_S(t), \quad (7.13)$$

$$\hat{u}_T(t) = u_T(t) + d_2 u_S(t) + c_T(t). \quad (7.14)$$

where $u_S(t)$ and $u_T(t)$ are the true source time courses of the seed and the target locations. In the above equations, $d_1 u_T(t)$ indicates the leakage of the target signal in the estimated seed signal, and $d_2 u_S(t)$ indicates the leakage of the seed signal in the estimated target signal. The real-valued constants d_1 and d_2 express the relative amount of these leakages.

In the equations above, $c_S(t)$ and $c_T(t)$ express the interference terms, which may include the leakage from the third source, contributions from external disturbances, and that of the sensor noise. The influence of these interference terms can be considered separately from the influence of the leakage terms [7]. In the arguments here, ignoring $c_S(t)$ and $c_T(t)$, the estimated spectra at the seed and target voxels, $\hat{\sigma}_S(f)$ and $\hat{\sigma}_T(f)$, are expressed as

$$\hat{\sigma}_S = \sigma_S + d_1 \sigma_T, \quad (7.15)$$

$$\hat{\sigma}_T = \sigma_T + d_2 \sigma_S. \quad (7.16)$$

where we omit the explicit notation of (f) for simplicity.

The magnitude coherence between the seed and the target voxels is expressed as

$$|\hat{\phi}| = \left| \frac{\langle \hat{\sigma}_T \hat{\sigma}_S^* \rangle}{\sqrt{\langle |\hat{\sigma}_T|^2 \rangle \langle |\hat{\sigma}_S|^2 \rangle}} \right|. \quad (7.17)$$

Using Eqs. (7.15) and (7.16), we have

$$\langle \hat{\sigma}_T \hat{\sigma}_S^* \rangle = \langle \sigma_T \sigma_S^* \rangle + d_1 \langle |\sigma_T|^2 \rangle + d_2 \langle |\sigma_S|^2 \rangle + d_1 d_2 \langle \sigma_S \sigma_T^* \rangle, \quad (7.18)$$

$$\langle |\hat{\sigma}_T|^2 \rangle = \langle |\sigma_S|^2 \rangle + d_1^2 \langle |\sigma_T|^2 \rangle + 2d_1 \Re(\langle \sigma_T \sigma_S^* \rangle), \quad (7.19)$$

$$\langle |\hat{\sigma}_S|^2 \rangle = \langle |\sigma_T|^2 \rangle + d_2^2 \langle |\sigma_S|^2 \rangle + 2d_2 \Re(\langle \sigma_T \sigma_S^* \rangle). \quad (7.20)$$

The equations above show that even when there is no true source interaction, i.e., even when $\langle \sigma_T \sigma_S^* \rangle = 0$ and $\langle \sigma_S \sigma_T^* \rangle = 0$, $|\hat{\phi}|$ has a nonzero value, which is equal to

$$|\hat{\phi}| = \frac{|d_1\langle|\sigma_T|^2\rangle + d_2\langle|\sigma_S|^2\rangle|}{\sqrt{(\langle|\sigma_S|^2\rangle + d_1^2\langle|\sigma_T|^2\rangle)(\langle|\sigma_T|^2\rangle + d_2^2\langle|\sigma_S|^2\rangle)}}. \quad (7.21)$$

7.4.2 Leakage Effects on the Imaginary Coherence

We next analyze the effects of the algorithm leakage on the imaginary coherence. Using Eq.(7.18) and the relationship

$$\langle\sigma_T\sigma_S^*\rangle + \langle\sigma_S\sigma_T^*\rangle = 2\Re(\langle\sigma_S\sigma_T^*\rangle),$$

the cross spectrum $\langle\hat{\sigma}_T\hat{\sigma}_S^*\rangle$ can be expressed as

$$\langle\hat{\sigma}_T\hat{\sigma}_S^*\rangle = (1 - d_1d_2)\langle\sigma_T\sigma_S^*\rangle + d_1\langle|\sigma_T|^2\rangle + d_2\langle|\sigma_S|^2\rangle + 2d_1d_2\Re(\langle\sigma_T\sigma_S^*\rangle). \quad (7.22)$$

By taking the imaginary part of Eq.(7.22), we can derive

$$\Im(\langle\hat{\sigma}_T\hat{\sigma}_S^*\rangle) = (1 - d_1d_2)\Im(\langle\sigma_T\sigma_S^*\rangle). \quad (7.23)$$

We can then obtain the imaginary part of the estimated coherence $\Im(\hat{\phi})$ as

$$\Im(\hat{\phi}) = \frac{\Im(\langle\hat{\sigma}_T\hat{\sigma}_S^*\rangle)}{\sqrt{\langle|\hat{\sigma}_T|^2\rangle\langle|\hat{\sigma}_S|^2\rangle}} = \frac{(1 - d_1d_2)\Im(\langle\sigma_T\sigma_S^*\rangle)}{\sqrt{\langle|\hat{\sigma}_T|^2\rangle\langle|\hat{\sigma}_S|^2\rangle}} = \Lambda \Im(\phi), \quad (7.24)$$

where $\Im(\phi)$ indicates the true value of the imaginary coherence. Using Eqs.(7.19) and (7.20), Λ is obtained as

$$\Lambda = \frac{(1 - d_1d_2)}{\sqrt{\tau_1\tau_2}} \quad (7.25)$$

where

$$\tau_1 = 1 + d_1^2 \frac{\langle|\sigma_T|^2\rangle}{\langle|\sigma_S|^2\rangle} + 2d_1 \frac{\Re(\langle\sigma_T\sigma_S^*\rangle)}{\langle|\sigma_S|^2\rangle}, \quad (7.26)$$

and

$$\tau_2 = 1 + d_2^2 \frac{\langle|\sigma_S|^2\rangle}{\langle|\sigma_T|^2\rangle} + 2d_2 \frac{\Re(\langle\sigma_T\sigma_S^*\rangle)}{\langle|\sigma_T|^2\rangle}. \quad (7.27)$$

Equation (7.24) shows that when $\Im(\phi) = 0$, we have $\Im(\hat{\phi}) = 0$, indicating that no spurious imaginary coherence has been caused. However, Eq.(7.24) also indicates that the value of $\Im(\hat{\phi})$ differs from the true value $\Im(\phi)$, i.e., the intensity of the estimated imaginary coherence is biased and the bias is represented by Λ in Eq.(7.25).

Let us calculate the intensity bias factor Λ assuming a simple scenario in which $d_1 = d_2 = d$ and $\langle |\sigma_S|^2 \rangle = \langle |\sigma_T|^2 \rangle$. Under these assumptions, the bias factor is simplified to

$$\frac{1 + |d|}{1 - |d|} \geq \Lambda = \frac{1 - d^2}{1 + d^2 + 2\Re(\phi)d} \geq \frac{1 - |d|}{1 + |d|}. \quad (7.28)$$

The equation above shows how the amount of leakage $|d|$ affects the bias factor Λ . It shows that if $|d|$ is as small as $|d| < 0.1$, the intensity bias is less than 10 %. However, when $|d|$ is as large as 0.4, we have $2.3 \geq \Lambda \geq 0.4$ and the intensity value of the imaginary coherence may have a very large bias. In the following, we introduce a metric referred to as the corrected imaginary coherence, which avoids this intensity bias.

7.5 Corrected Imaginary Coherence

7.5.1 Modification of Imaginary Coherence

The imaginary coherence can avoid spurious results caused by the algorithm leakage. However, its intensity value is affected by the algorithm leakage. In this section, we introduce a modified form of the imaginary coherence whose intensity values are unaffected by the algorithm leakage. Let us define the corrected imaginary coherence as

$$\xi = \frac{\Im(\phi)}{\sqrt{1 - \Re(\phi)^2}}. \quad (7.29)$$

Let us express the estimate of ξ , $\hat{\xi}$, such that

$$\hat{\xi} = \frac{\Im(\hat{\phi})}{\sqrt{1 - \Re(\hat{\phi})^2}} = \frac{\Im((\hat{\sigma}_T \hat{\sigma}_S^*))}{\sqrt{\langle |\hat{\sigma}_T|^2 \rangle \langle |\hat{\sigma}_S|^2 \rangle - \Re((\hat{\sigma}_T \hat{\sigma}_S^*))}}.$$

Then, considering

$$\Re((\hat{\sigma}_T \hat{\sigma}_S^*)) = d_1^2 \langle |\sigma_T|^2 \rangle + d_2^2 \langle |\sigma_S|^2 \rangle + (1 + d_1 d_2) \Re(\langle \sigma_T \sigma_S^* \rangle), \quad (7.30)$$

and using Eqs. (7.19) and (7.20), we derive

$$\langle |\hat{\sigma}_T|^2 \rangle \langle |\hat{\sigma}_S|^2 \rangle - \Re((\hat{\sigma}_T \hat{\sigma}_S^*)) = (1 - d_1 d_2)^2 \left(\langle |\sigma_T|^2 \rangle \langle |\sigma_S|^2 \rangle - \Re(\langle \sigma_T \sigma_S^* \rangle) \right). \quad (7.31)$$

Combining this with Eq. (7.23), we finally obtain

$$\begin{aligned}
\widehat{\xi} &= \frac{\Im(\widehat{\phi})}{\sqrt{1 - \Re(\widehat{\phi})^2}} \\
&= \frac{\Im(\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle)}{\sqrt{\langle |\widehat{\sigma}_T|^2 \rangle \langle |\widehat{\sigma}_S|^2 \rangle - \Re(\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle)}} \\
&= \frac{(1 - d_1 d_2) \Im(\langle \sigma_T \sigma_S^* \rangle)}{(1 - d_1 d_2) \sqrt{\langle |\sigma_T|^2 \rangle \langle |\sigma_S|^2 \rangle - \Re(\langle \sigma_T \sigma_S^* \rangle)}} \\
&= \frac{\Im(\phi)}{\sqrt{1 - \Re(\phi)^2}}.
\end{aligned} \tag{7.32}$$

The above equation indicates that the corrected imaginary coherence computed using voxel spectra is exactly equal to the true corrected imaginary coherence. This implies that the corrected imaginary coherence is unaffected by the algorithm leakage.

In the above arguments, the corrected imaginary coherence is introduced somewhat in an ad hoc manner. In the following, we derive the corrected imaginary coherence in two different manners: factorization of the mutual coherence and regression of the target signal with the seed signal. These derivations provide some insights into the nature of corrected imaginary coherence.

7.5.2 Factorization of Mutual Information

Here we show that the corrected imaginary coherence is derived from the factorization of the mutual information into the instantaneous and non-instantaneous components. We assume that $\widehat{\sigma}_T$ and $\widehat{\sigma}_S$ are Gaussian distributed, circular complex random variables. Concise explanations on the complex Gaussian random variable are found in Sect. C.2 in the Appendix.

To derive the mutual information in the frequency domain, we define the entropy of voxel spectra $\widehat{\sigma}_T$ and $\widehat{\sigma}_S$. To do so, we first define real-valued 2×1 vectors $\zeta_T = [\Re(\widehat{\sigma}_T), \Im(\widehat{\sigma}_T)]^T$ and $\zeta_S = [\Re(\widehat{\sigma}_S), \Im(\widehat{\sigma}_S)]^T$. Using these vectors, the entropy is defined for $\widehat{\sigma}_T$ and $\widehat{\sigma}_S$ such that,

$$\mathcal{H}(\widehat{\sigma}_T) = - \int p(\zeta_T) \log p(\zeta_T) d\zeta_T, \tag{7.33}$$

$$\mathcal{H}(\widehat{\sigma}_S) = - \int p(\zeta_S) \log p(\zeta_S) d\zeta_S. \tag{7.34}$$

The entropy is a metric for uncertainty. $\mathcal{H}(\widehat{\sigma}_T)$ represents the uncertainty when $\widehat{\sigma}_T$ is unknown, and $\mathcal{H}(\widehat{\sigma}_S)$ represents the uncertainty when $\widehat{\sigma}_S$ is unknown. The joint entropy is defined as

$$\mathcal{H}(\widehat{\sigma}_T, \widehat{\sigma}_S) = - \iint p(\zeta_T, \zeta_S) \log p(\zeta_T, \zeta_S) d\zeta_T d\zeta_S. \quad (7.35)$$

The mutual information between $\widehat{\sigma}_T$ and $\widehat{\sigma}_S$ is then defined as

$$\mathcal{I}(\widehat{\sigma}_T, \widehat{\sigma}_S) = \mathcal{H}(\widehat{\sigma}_T) + \mathcal{H}(\widehat{\sigma}_S) - \mathcal{H}(\widehat{\sigma}_T, \widehat{\sigma}_S). \quad (7.36)$$

When $\widehat{\sigma}_T$ and $\widehat{\sigma}_S$ are independent, we have $\mathcal{I}(\widehat{\sigma}_T, \widehat{\sigma}_S) = 0$.

Under the assumption that $\widehat{\sigma}_T$ and $\widehat{\sigma}_S$ are complex Gaussian distributed, the entropy is expressed as

$$\mathcal{H}(\widehat{\sigma}_T) = \log \left\langle |\widehat{\sigma}_T|^2 \right\rangle, \quad (7.37)$$

$$\mathcal{H}(\widehat{\sigma}_S) = \log \left\langle |\widehat{\sigma}_S|^2 \right\rangle. \quad (7.38)$$

The joint entropy is given by

$$\begin{aligned} \mathcal{H}(\widehat{\sigma}_T, \widehat{\sigma}_S) &= \log \left| \left\langle \begin{bmatrix} \widehat{\sigma}_T \\ \widehat{\sigma}_S \end{bmatrix} \widehat{\phi} [\widehat{\sigma}_T^* \widehat{\sigma}_S^*] \right\rangle \right| \\ &= \log \left| \frac{\langle |\widehat{\sigma}_T|^2 \rangle \langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle}{\langle \widehat{\sigma}_S \widehat{\sigma}_T^* \rangle \langle |\widehat{\sigma}_S|^2 \rangle} \right| \\ &= \log \left(\langle |\widehat{\sigma}_T|^2 \rangle \langle |\widehat{\sigma}_S|^2 \rangle - |\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle|^2 \right) \end{aligned} \quad (7.39)$$

Therefore, the mutual information between $\widehat{\sigma}_T$ and $\widehat{\sigma}_S$ is obtained as

$$\begin{aligned} \mathcal{I}(\widehat{\sigma}_T, \widehat{\sigma}_S) &= \mathcal{H}(\widehat{\sigma}_T) + \mathcal{H}(\widehat{\sigma}_S) - \mathcal{H}(\widehat{\sigma}_T, \widehat{\sigma}_S) \\ &= \log \frac{\langle |\widehat{\sigma}_T|^2 \rangle \langle |\widehat{\sigma}_S|^2 \rangle}{\langle |\widehat{\sigma}_T|^2 \rangle \langle |\widehat{\sigma}_S|^2 \rangle - |\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle|^2}. \end{aligned} \quad (7.40)$$

It is easy to see that the mutual information $\mathcal{I}(\widehat{\sigma}_T, \widehat{\sigma}_S)$ is related to the magnitude coherence $|\widehat{\phi}|$ such that

$$\mathcal{I}(\widehat{\sigma}_T, \widehat{\sigma}_S) = -\log(1 - |\widehat{\phi}|^2), \quad (7.41)$$

where

$$|\widehat{\phi}|^2 = \frac{|\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle|^2}{\langle |\widehat{\sigma}_T|^2 \rangle \langle |\widehat{\sigma}_S|^2 \rangle}. \quad (7.42)$$

Let us consider the factorization of the mutual information using the real and imaginary parts of coherence. Using $|\widehat{\phi}|^2 = \Re(\widehat{\phi})^2 + \Im(\widehat{\phi})^2$, we have

$$\begin{aligned}
\mathcal{I}(\widehat{\sigma}_T, \widehat{\sigma}_S) &= -\log(1 - |\widehat{\phi}|^2) = -\log(1 - \Re(\widehat{\phi})^2 - \Im(\widehat{\phi})^2) \\
&= -\log \left[\left(1 - \Re(\widehat{\phi})^2\right) \left[1 - \frac{\Im(\widehat{\phi})^2}{1 - \Re(\widehat{\phi})^2}\right] \right] \\
&= -\log(1 - \Re(\widehat{\phi})^2) - \log \left[1 - \frac{\Im(\widehat{\phi})^2}{1 - \Re(\widehat{\phi})^2} \right]. \tag{7.43}
\end{aligned}$$

On the right-hand side, the first term,

$$\mathcal{I}_R(\widehat{\sigma}_T, \widehat{\sigma}_S) = -\log(1 - \Re(\widehat{\phi})^2)$$

represents a component of the mutual information corresponding to the real part of coherence. This $\mathcal{I}_R(\widehat{\sigma}_T, \widehat{\sigma}_S)$ can be interpreted as the instantaneous component, which corresponds to the zero-lag correlation between the target and the seed time courses, as discussed in Sect. 7.3.

The second term,

$$\mathcal{I}_I(\widehat{\sigma}_T, \widehat{\sigma}_S) = -\log \left[\left(1 - \frac{\Im(\widehat{\phi})^2}{1 - \Re(\widehat{\phi})^2}\right) \right]$$

represents a component of the mutual information corresponding to the imaginary part of coherence. It is interpreted as the non-instantaneous component, which corresponds to the nonzero-lag correlation between $u_T(t)$ and $u_S(t)$. It is easy to see

$$\mathcal{I}_I(\widehat{\sigma}_T, \widehat{\sigma}_S) = -\log \left[\left(1 - \frac{\Im(\widehat{\phi})^2}{1 - \Re(\widehat{\phi})^2}\right) \right] = -\log(1 - \widehat{\xi}^2), \tag{7.44}$$

where $\widehat{\xi}$ is the corrected imaginary coherence. The equation above indicates that the corrected imaginary coherence can be interpreted as a coherence-domain expression of the non-instantaneous component of the mutual information. Arguments similar to those in this section are found in [9].

7.5.3 Residual Coherence

We next show that the corrected imaginary coherence can be derived using regression analysis. Let us regress the target spectrum $\widehat{\sigma}_T$ using the seed spectrum $\widehat{\sigma}_S$, such that

$$\widehat{\sigma}_T = \alpha \widehat{\sigma}_S + v, \tag{7.45}$$

where α is a real-valued constant, and v is a residual signal of this regression. The value for α is determined using the least-squares fit:

$$\alpha = \operatorname{argmin}_{\alpha} \langle |\widehat{\sigma}_T - \alpha \widehat{\sigma}_S|^2 \rangle, \quad (7.46)$$

and is derived as

$$\alpha = \frac{\Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]}{\langle |\widehat{\sigma}_S|^2 \rangle}. \quad (7.47)$$

Thus, the residual signal is expressed as

$$v = \widehat{\sigma}_T - \alpha \widehat{\sigma}_S = \widehat{\sigma}_T - \frac{\Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]}{\langle |\widehat{\sigma}_S|^2 \rangle} \widehat{\sigma}_S. \quad (7.48)$$

We then define the residual coherence, $\widehat{\phi}_R$, such that

$$\widehat{\phi}_R = \frac{\langle v \widehat{\sigma}_S^* \rangle}{\sqrt{\langle |v|^2 \rangle \langle |\widehat{\sigma}_S|^2 \rangle}}. \quad (7.49)$$

Let us show that this residual coherence is equal to the corrected imaginary coherence. First, the cross spectrum between the residual and seed signals is expressed as

$$\begin{aligned} \langle v \widehat{\sigma}_S^* \rangle &= \langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle - \frac{\Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]}{\langle |\widehat{\sigma}_S|^2 \rangle} \langle |\widehat{\sigma}_S|^2 \rangle \\ &= \langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle - \Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle] \\ &= i \Im[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]. \end{aligned} \quad (7.50)$$

The equation above shows that the cross spectrum between v and $\widehat{\sigma}_S$ is equal to the imaginary part of the cross spectrum between $\widehat{\sigma}_T$ and $\widehat{\sigma}_S$. Using Eq. (7.48), we express $\langle |v|^2 \rangle$ such that

$$\begin{aligned} \langle |v|^2 \rangle &= \left\langle \left[\widehat{\sigma}_T - \frac{\Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]}{\langle |\widehat{\sigma}_S|^2 \rangle} \widehat{\sigma}_S \right] \left[\widehat{\sigma}_T - \frac{\Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]}{\langle |\widehat{\sigma}_S|^2 \rangle} \widehat{\sigma}_S \right]^* \right\rangle \\ &= \langle |\widehat{\sigma}_T|^2 \rangle - \frac{\Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]}{\langle |\widehat{\sigma}_S|^2 \rangle} (\widehat{\sigma}_S \widehat{\sigma}_T^* + \widehat{\sigma}_T \widehat{\sigma}_S^*) + \frac{\Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]^2}{\langle |\widehat{\sigma}_S|^2 \rangle^2} \langle |\widehat{\sigma}_S|^2 \rangle \\ &= \langle |\widehat{\sigma}_T|^2 \rangle - \frac{\Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]^2}{\langle |\widehat{\sigma}_S|^2 \rangle}. \end{aligned} \quad (7.51)$$

Therefore, we get

$$\langle |v|^2 \rangle \langle |\widehat{\sigma}_S|^2 \rangle = \langle |\widehat{\sigma}_T|^2 \rangle \langle |\widehat{\sigma}_S|^2 \rangle - \Re[\langle \widehat{\sigma}_T \widehat{\sigma}_S^* \rangle]^2. \quad (7.52)$$

Substituting Eqs. (7.52) and (7.50) into Eq. (7.49), and neglecting the constant i , the residual coherence is expressed as

$$\hat{\phi}_R = \frac{\Im[(\hat{\sigma}_T \hat{\sigma}_S^*)]}{\sqrt{(\langle |\hat{\sigma}_T|^2 \rangle \langle |\hat{\sigma}_S|^2 \rangle - \Re[\langle \hat{\sigma}_T \hat{\sigma}_S^* \rangle]^2}} = \frac{\Im(\hat{\phi})}{\sqrt{1 - \Re(\hat{\phi})^2}} = \hat{\xi}. \quad (7.53)$$

The above equation shows that the residual coherence is exactly equal to the corrected imaginary coherence $\hat{\xi}$. That is, the corrected imaginary coherence is derived as the coherence between the seed signal and the residual signal, which is obtained by regressing out the seed signal from the target signal.

7.5.4 Phase Dependence of the Corrected Imaginary Coherences

The coherence ϕ is complex valued, and it is expressed as the amplitude and phase, such that

$$\phi = |\phi| e^{i\theta}, \quad (7.54)$$

where θ is the phase of the coherence, and the notation (f) is omitted for simplicity. When the coherence is caused by true brain interactions, we naturally interpret that the amplitude $|\phi|$, which is equal to the magnitude coherence, represents the strength of the brain interactions.

The imaginary coherence is expressed as

$$\Im(\phi) = |\phi| \sin \theta. \quad (7.55)$$

The value of the imaginary coherence depends on the phase θ . This is one weak point of the imaginary coherence, because even if there is a strong brain interaction, the imaginary coherence becomes very small when the phase θ has a value close to one of multiples of π .

Let us see how the corrected imaginary coherence ξ depends on the phase. Using Eq. (7.54), we have

$$\xi = \frac{|\phi| \sin \theta}{\sqrt{1 - |\phi|^2 \cos^2 \theta}}. \quad (7.56)$$

The plots of ξ with respect to the phase are shown in Fig. 7.1 with four values of the magnitude coherence $|\phi|$. This figure shows that the corrected imaginary coherence becomes asymptotically independent of the phase with the limit of $|\phi| \rightarrow 1$. However, the corrected imaginary coherence has almost the same phase dependence as the imaginary coherence, when the magnitude coherence $|\phi|$ is less than 0.6. When $|\phi|$ becomes closer to 1, the corrected imaginary coherence ξ becomes significantly greater than the imaginary coherence $\Im(\phi)$.

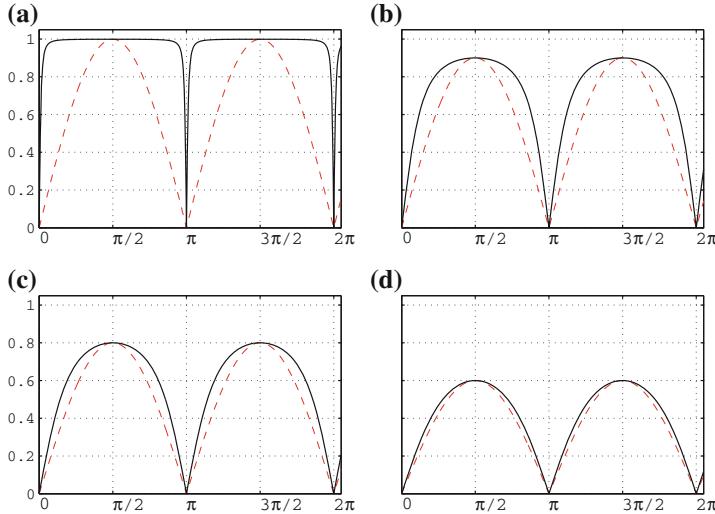


Fig. 7.1 Plots of the imaginary coherence and corrected imaginary coherence with respect to the phase for four values of magnitude coherence. **a** $|\phi| = 0.99$. **b** $|\phi| = 0.9$. **c** $|\phi| = 0.8$. **d** $|\phi| = 0.6$. The *solid line* shows the corrected imaginary coherence and the *broken line* shows the imaginary coherence

7.6 Canonical Coherence

7.6.1 Canonical Magnitude Coherence

So far, we have discussed methods to compute voxel-based coherence. However, when the number of voxels is large and our target is to analyze all-voxel-to-all-voxel connections, the interpretation and visualization of analysis results may not be easy because the number of voxel connections becomes huge. One way to reduce this difficulty is to compute coherence between regions determined based on the neurophysiology and/or neuroanatomy of a brain. However, in this case, since a brain region contains multiple voxels, we must compute the coherence between groups of multiple time courses. We here extend the theory of canonical correlation described in Sect. C.3 in the Appendix to compute the canonical coherence, which should be effective for region-based connectivity analysis. The explanation here follows those in [10, 11].

In this section, random variables \mathbf{x} and \mathbf{y} are complex-valued column vectors in the frequency domain, i.e.,

$$\mathbf{x}(f) = \begin{bmatrix} x_1(f) \\ x_2(f) \\ \vdots \\ x_p(f) \end{bmatrix} \quad \text{and} \quad \mathbf{y}(f) = \begin{bmatrix} y_1(f) \\ y_2(f) \\ \vdots \\ y_q(f) \end{bmatrix}. \quad (7.57)$$

In the following arguments, the notation (f) is omitted for simplicity. Using the same arguments as those for the canonical correlation in Sect.C.3, the random vectors \mathbf{x} and \mathbf{y} are, respectively, projected in the directions of \mathbf{a} and \mathbf{b} where \mathbf{a} and \mathbf{b} are complex-valued $p \times 1$ and $q \times 1$ column vectors. That is, defining $\hat{\mathbf{x}} = \mathbf{a}^H \mathbf{x}$ and $\hat{\mathbf{y}} = \mathbf{b}^H \mathbf{y}$, the coherence between $\hat{\mathbf{x}}$ and $\hat{\mathbf{y}}$, ϕ , is given by

$$\begin{aligned} \phi &= \frac{\langle \hat{\mathbf{x}} \hat{\mathbf{y}}^* \rangle}{\sqrt{\langle \hat{\mathbf{x}} \hat{\mathbf{x}}^* \rangle} \sqrt{\langle \hat{\mathbf{y}} \hat{\mathbf{y}}^* \rangle}} = \frac{\mathbf{a}^H \langle \mathbf{x} \mathbf{y}^H \rangle \mathbf{b}}{\sqrt{[\mathbf{a}^H \langle \mathbf{x} \mathbf{x}^H \rangle \mathbf{a}] [\mathbf{b}^H \langle \mathbf{y} \mathbf{y}^H \rangle \mathbf{b}]}} \\ &= \frac{\mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b}}{\sqrt{[\mathbf{a}^H \boldsymbol{\Sigma}_{xx} \mathbf{a}] [\mathbf{b}^H \boldsymbol{\Sigma}_{yy} \mathbf{b}]}}, \end{aligned} \quad (7.58)$$

where the superscript H indicates the Hermitian transpose. Here the cross-spectral matrices are defined as $\boldsymbol{\Sigma}_{xy} = \langle \mathbf{x} \mathbf{y}^H \rangle$, $\boldsymbol{\Sigma}_{xx} = \langle \mathbf{x} \mathbf{x}^H \rangle$, and $\boldsymbol{\Sigma}_{yy} = \langle \mathbf{y} \mathbf{y}^H \rangle$. We derive a formula to compute the canonical magnitude coherence. Using Eq.(7.58), the magnitude coherence between $\hat{\mathbf{x}}$ and $\hat{\mathbf{y}}$ is expressed as

$$|\phi|^2 = \frac{|\mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b}|^2}{[\mathbf{a}^H \boldsymbol{\Sigma}_{xx} \mathbf{a}] [\mathbf{b}^H \boldsymbol{\Sigma}_{yy} \mathbf{b}]}.$$
(7.59)

The canonical squared magnitude coherence $|\psi|^2$ is defined as the maximum of $|\phi|^2$ (with respect to \mathbf{a} and \mathbf{b}), which is obtained by solving the optimization problem

$$|\psi|^2 = \max_{\mathbf{a}, \mathbf{b}} |\mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b}|^2 \quad \text{subject to } \mathbf{a}^H \boldsymbol{\Sigma}_{xx} \mathbf{a} = 1 \quad \text{and} \quad \mathbf{b}^H \boldsymbol{\Sigma}_{yy} \mathbf{b} = 1.$$
(7.60)

This constrained optimization problem can be solved in the following manner. Since $|\mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b}|^2$ is a scalar, we have the relationship,

$$|\mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b}|^2 = (\mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b})^H (\mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b}) = \mathbf{b}^H \boldsymbol{\Sigma}_{xy}^H \mathbf{a} \mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b}.$$
(7.61)

We first fix \mathbf{a} , and solve the maximization problem with respect to \mathbf{b} . The maximization problem is

$$\max_{\mathbf{b}} \mathbf{b}^H \boldsymbol{\Sigma}_{xy}^H \mathbf{a} \mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b} \quad \text{subject to} \quad \mathbf{b}^H \boldsymbol{\Sigma}_{yy} \mathbf{b} = 1.$$
(7.62)

Using the same derivation mentioned in Sect. C.9, the maximum value is equal to the largest eigenvalue of the matrix

$$\boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^H \mathbf{a} \mathbf{a}^H \boldsymbol{\Sigma}_{xy}.$$

The matrix above is a rank-one matrix and it only has a single nonzero eigenvalue. According to Sect. C.8 (Property No. 10), this eigenvalue is equal to

$$\mathbf{a}^H \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^H \mathbf{a}, \quad (7.63)$$

which is a scalar.

Thus, the optimization

$$|\psi|^2 = \max_{\mathbf{a}, \mathbf{b}} \mathbf{b}^H \boldsymbol{\Sigma}_{xy}^H \mathbf{a} \mathbf{a}^H \boldsymbol{\Sigma}_{xy} \mathbf{b} \quad \text{subject to } \mathbf{a}^H \boldsymbol{\Sigma}_{xx} \mathbf{a} = 1 \quad \text{and} \quad \mathbf{b}^H \boldsymbol{\Sigma}_{yy} \mathbf{b} = 1, \quad (7.64)$$

is now rewritten as

$$|\psi|^2 = \max_{\mathbf{a}} \mathbf{a}^H \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^H \mathbf{a} \quad \text{subject to} \quad \mathbf{a}^H \boldsymbol{\Sigma}_{xx} \mathbf{a} = 1. \quad (7.65)$$

Using again the derivation described in Sect. C.9, the solution of this maximization is obtained as the maximum eigenvalue of the matrix

$$\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^H.$$

That is, denoting the eigenvalues of this matrix as γ_j where $j = 1, \dots, d$ and $d = \min\{p, q\}$, the canonical squared magnitude coherence is derived as

$$|\psi|^2 = \mathcal{S}_{\max}\{\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^H\} = \gamma_1, \quad (7.66)$$

where the notation $\mathcal{S}_{\max}\{\cdot\}$ indicates the maximum eigenvalue of a matrix between the parentheses, as is defined in Sect. C.9.

This canonical squared magnitude coherence is considered the best overall magnitude coherence measure between the two sets of multiple spectra x_1, \dots, x_p and y_1, \dots, y_q , and it is equal to the maximum eigenvalue γ_1 in Eq. (7.66). However, other eigenvalues may have information complementary to γ_1 , and therefore, a metric that uses all the eigenvalues may be preferable. Let us assume the random vectors \mathbf{x} and \mathbf{y} to be complex Gaussian. According to Eq. (C.52), we can then define the mutual information between \mathbf{x} and \mathbf{y} such that

$$\mathcal{I}(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^d \log \frac{1}{1 - \gamma_j}, \quad (7.67)$$

which is a metric using all the eigenvalues. In the arguments in Sect. 7.5.2, the relationship between the coherence and the mutual information is given in Eq. (7.41). This relationship can lead to an alternative definition of canonical magnitude coherence based on the mutual information, such that

$$|\tilde{\psi}|^2 = 1 - \exp[-\mathcal{I}(\mathbf{x}, \mathbf{y})] = 1 - \prod_{j=1}^d (1 - \gamma_j). \quad (7.68)$$

This metric uses all the eigenvalues, and since its value is normalized between 0 and 1, the interpretation is intuitively easy, compared to the original mutual information $\mathcal{I}(\mathbf{x}, \mathbf{y})$.

Also, the mutual information $\mathcal{I}(\mathbf{x}, \mathbf{y})$ in Eq. (7.67) depends on d , which is the sizes of the vectors \mathbf{x} and \mathbf{y} . Such property may not be appropriate for a brain interaction metric, and the metric independent of the sizes of vectors can be defined such that,

$$\tilde{\mathcal{I}}(\mathbf{x}, \mathbf{y}) = \frac{1}{d} \sum_{j=1}^d \log \frac{1}{1 - \gamma_j}. \quad (7.69)$$

The alternative definition of canonical magnitude coherence in this case is expressed as

$$|\check{\psi}|^2 = 1 - \exp[-\tilde{\mathcal{I}}(\mathbf{x}, \mathbf{y})] = 1 - \left[\prod_{j=1}^d (1 - \gamma_j) \right]^{1/d}. \quad (7.70)$$

The above metric may be more effective than the one in Eq. (7.68).

7.6.2 Canonical Imaginary Coherence

We next derive a formula to compute the canonical imaginary coherence. Following the arguments of Evald et al. [12] we define

$$\boldsymbol{\Sigma}_{xx} = \Re(\boldsymbol{\Sigma}_{xx}) + i\Im(\boldsymbol{\Sigma}_{xx}) = \boldsymbol{\Gamma}_{xx} + i\boldsymbol{\Upsilon}_{xx}, \quad (7.71)$$

$$\boldsymbol{\Sigma}_{yy} = \Re(\boldsymbol{\Sigma}_{yy}) + i\Im(\boldsymbol{\Sigma}_{yy}) = \boldsymbol{\Gamma}_{yy} + i\boldsymbol{\Upsilon}_{yy}, \quad (7.72)$$

$$\boldsymbol{\Sigma}_{xy} = \Re(\boldsymbol{\Sigma}_{xy}) + i\Im(\boldsymbol{\Sigma}_{xy}) = \boldsymbol{\Gamma}_{xy} + i\boldsymbol{\Upsilon}_{xy}. \quad (7.73)$$

We use real-valued \mathbf{a} and \mathbf{b} and express the imaginary part of the coherence between $\mathbf{a}^T \mathbf{x}$ and $\mathbf{b}^T \mathbf{y}$:

$$\Im(\phi) = \frac{\mathbf{a}^T \Im(\boldsymbol{\Sigma}_{xy}) \mathbf{b}}{\sqrt{[\mathbf{a}^T \boldsymbol{\Sigma}_{xx} \mathbf{a}] [\mathbf{b}^T \boldsymbol{\Sigma}_{yy} \mathbf{b}]}} = \frac{\mathbf{a}^T \boldsymbol{\Upsilon}_{xy} \mathbf{b}}{\sqrt{[\mathbf{a}^T \boldsymbol{\Sigma}_{xx} \mathbf{a}] [\mathbf{b}^T \boldsymbol{\Sigma}_{yy} \mathbf{b}]}}. \quad (7.74)$$

The canonical imaginary coherence ψ_I is defined as the maximum of $\Im(\phi)$, which is obtained using the maximization,

$$\psi_I = \max_{\mathbf{a}, \mathbf{b}} \mathbf{a}^T \boldsymbol{\Upsilon}_{xy} \mathbf{b}, \quad \text{subject to } \mathbf{a}^T \boldsymbol{\Sigma}_{xx} \mathbf{a} = 1 \quad \text{and} \quad \mathbf{b}^T \boldsymbol{\Sigma}_{yy} \mathbf{b} = 1. \quad (7.75)$$

The formulation above is exactly the same as that for the canonical correlation in Eq.(C.31). Thus, as in Sect.C.3.1, this maximization problem might seem to be rewritten as the eigenvalue problem,

$$\left(\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Upsilon}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Upsilon}_{xy}^T \right) \mathbf{a} = \psi_I^2 \mathbf{a}, \quad (7.76)$$

and ψ_I^2 is obtained as its maximum eigenvalue. However, in this eigenproblem, the eigenvector \mathbf{a} is generally complex-valued, conflicting with the assumption that the vectors \mathbf{a} and \mathbf{b} are real-valued.

We instead use real-valued $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ obtained such that

$$\boldsymbol{\alpha} = \boldsymbol{\Gamma}_{xx}^{1/2} \mathbf{a}, \quad (7.77)$$

$$\boldsymbol{\beta} = \boldsymbol{\Gamma}_{yy}^{1/2} \mathbf{b}. \quad (7.78)$$

Then, using $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$, we have

$$\begin{aligned} \mathbf{a}^T \boldsymbol{\Sigma}_{xx} \mathbf{a} &= \boldsymbol{\alpha}^T \boldsymbol{\Gamma}_{xx}^{-1/2} \boldsymbol{\Sigma}_{xx} \boldsymbol{\Gamma}_{xx}^{-1/2} \boldsymbol{\alpha} \\ &= \boldsymbol{\alpha}^T \boldsymbol{\Gamma}_{xx}^{-1/2} (\boldsymbol{\Gamma}_{xx} + i \boldsymbol{\Upsilon}_{xx}) \boldsymbol{\Gamma}_{xx}^{-1/2} \boldsymbol{\alpha} \\ &= \|\boldsymbol{\alpha}\|^2 + i \boldsymbol{\alpha}^T \boldsymbol{\Gamma}_{xx}^{-1/2} \boldsymbol{\Upsilon}_{xx} \boldsymbol{\Gamma}_{xx}^{-1/2} \boldsymbol{\alpha} = \|\boldsymbol{\alpha}\|^2 \end{aligned} \quad (7.79)$$

and

$$\begin{aligned} \mathbf{b}^T \boldsymbol{\Sigma}_{yy} \mathbf{b} &= \boldsymbol{\beta}^T \boldsymbol{\Gamma}_{yy}^{-1/2} \boldsymbol{\Sigma}_{yy} \boldsymbol{\Gamma}_{yy}^{-1/2} \boldsymbol{\beta} \\ &= \boldsymbol{\beta}^T \boldsymbol{\Gamma}_{yy}^{-1/2} (\boldsymbol{\Gamma}_{yy} + i \boldsymbol{\Upsilon}_{yy}) \boldsymbol{\Gamma}_{yy}^{-1/2} \boldsymbol{\beta} \\ &= \|\boldsymbol{\beta}\|^2 + i \boldsymbol{\beta}^T \boldsymbol{\Gamma}_{yy}^{-1/2} \boldsymbol{\Upsilon}_{yy} \boldsymbol{\Gamma}_{yy}^{-1/2} \boldsymbol{\beta} = \|\boldsymbol{\beta}\|^2. \end{aligned} \quad (7.80)$$

In the equations above, we use the fact that the matrices $\boldsymbol{\Gamma}_{xx}^{-1/2} \boldsymbol{\Upsilon}_{xx} \boldsymbol{\Gamma}_{xx}^{-1/2}$ and $\boldsymbol{\Gamma}_{yy}^{-1/2} \boldsymbol{\Upsilon}_{yy} \boldsymbol{\Gamma}_{yy}^{-1/2}$ are skew-symmetric.²

Using $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$, the optimization in Eq.(7.75) can be rewritten as

$$\psi_I = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \boldsymbol{\alpha}^T \boldsymbol{\Gamma}_{xx}^{-1/2} \boldsymbol{\Upsilon}_{xy} \boldsymbol{\Gamma}_{yy}^{-1/2} \boldsymbol{\beta} \quad \text{subject to } \boldsymbol{\alpha}^T \boldsymbol{\alpha} = 1 \quad \text{and} \quad \boldsymbol{\beta}^T \boldsymbol{\beta} = 1. \quad (7.81)$$

² If $\mathbf{A}^T = -\mathbf{A}$ holds, matrix \mathbf{A} is skew-symmetric. Since $\boldsymbol{\alpha}^T \mathbf{A} \boldsymbol{\alpha}$ is a scalar, if \mathbf{A} is skew-symmetric, $\boldsymbol{\alpha}^T \mathbf{A} \boldsymbol{\alpha} = 0$.

Following the derivation in Sect.C.3.1, the solution of the above optimization problem is expressed as the eigenvalue problem,

$$\left(\boldsymbol{\Pi}_{xy} \boldsymbol{\Pi}_{xy}^T \right) \boldsymbol{\alpha} = \psi_I^2 \boldsymbol{\alpha}, \quad (7.82)$$

where $\boldsymbol{\Pi}_{xy} = \boldsymbol{\Gamma}_{xx}^{-1/2} \boldsymbol{\Upsilon}_{xy} \boldsymbol{\Gamma}_{yy}^{-1/2}$. In the above equation, since the matrix, $\boldsymbol{\Pi}_{xy} \boldsymbol{\Pi}_{xy}^T$ is a real symmetric matrix, the eigenvector $\boldsymbol{\alpha}$ is real-valued. The vector $\boldsymbol{\beta}$ is obtained using

$$\boldsymbol{\beta} = \frac{1}{\psi_I} \boldsymbol{\Upsilon}_{xy}^T \boldsymbol{\alpha},$$

which is also real-valued.

The matrices $\boldsymbol{\Pi}_{xy} \boldsymbol{\Pi}_{xy}^T$ and $\boldsymbol{\Gamma}_{xx}^{-1} \boldsymbol{\Upsilon}_{xy} \boldsymbol{\Gamma}_{yy}^{-1} \boldsymbol{\Upsilon}_{xy}^T$ have the same eigenvalues, according to Sect.C.8 (Property No. 9). Let us define the eigenvalues of the matrix $\boldsymbol{\Gamma}_{xx}^{-1} \boldsymbol{\Upsilon}_{xy} \boldsymbol{\Gamma}_{yy}^{-1} \boldsymbol{\Upsilon}_{xy}^T$ as ζ_j ($j = 1, \dots, d$). The canonical imaginary coherence ψ_I^2 is derived as

$$\psi_I^2 = \mathcal{S}_{\max} \{ \boldsymbol{\Gamma}_{xx}^{-1} \boldsymbol{\Upsilon}_{xy} \boldsymbol{\Gamma}_{yy}^{-1} \boldsymbol{\Upsilon}_{xy}^T \} = \zeta_1. \quad (7.83)$$

Using the same arguments in the preceding section, the imaginary-coherence-based mutual information is given by

$$\mathcal{I}_s(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^d \log \frac{1}{1 - \zeta_j}. \quad (7.84)$$

The alternative definition of canonical imaginary coherence, based on $\mathcal{I}_s(\mathbf{x}, \mathbf{y})$, is obtained as

$$\tilde{\psi}_I^2 = 1 - \exp[-\mathcal{I}_s(\mathbf{x}, \mathbf{y})] = 1 - \prod_{j=1}^d (1 - \zeta_j), \quad (7.85)$$

which uses all the eigenvalues ζ_1, \dots, ζ_d . When we can assume $\zeta_j \ll 1$ (where $j = 1, \dots, d$), we have

$$\tilde{\psi}_I^2 \approx \sum_{j=1}^d \zeta_j. \quad (7.86)$$

The connectivity metric above has been proposed and called the multivariate interaction measure (MIM) in [12].

The mutual information independent of the sizes of vectors is defined such that,

$$\tilde{\mathcal{I}}_s(\mathbf{x}, \mathbf{y}) = \frac{1}{d} \sum_{j=1}^d \log \frac{1}{1 - \zeta_j}. \quad (7.87)$$

The alternative canonical imaginary coherence, based on $\tilde{\mathcal{I}}_s(\mathbf{x}, \mathbf{y})$, is obtained as

$$\check{\psi}_I^2 = 1 - \exp[-\mathcal{I}_s(\mathbf{x}, \mathbf{y})] = 1 - \left[\prod_{j=1}^d (1 - \zeta_j) \right]^{1/d}. \quad (7.88)$$

When we can assume $\zeta_j \ll 1$ (where $j = 1, \dots, d$), we have

$$\check{\psi}_I^2 \approx \frac{1}{d} \sum_{j=1}^d \zeta_j. \quad (7.89)$$

The connectivity metric above is called the global interaction measure (GIM) [12].

7.6.3 Canonical Residual Coherence

We can compute the canonical residual coherence, which is a multivariate version of the residual coherence described in Sect. 7.5.3. Let us assume that \mathbf{y} is the target spectra, \mathbf{x} is the seed spectra, and consider the regression

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{v}. \quad (7.90)$$

The real-valued regression coefficient matrix \mathbf{A} is obtained by the least-squares fit. The optimum \mathbf{A} is obtained as

$$\widehat{\mathbf{A}} = \Re \left(\boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \right). \quad (7.91)$$

Thus, the residual signal \mathbf{v} is expressed as

$$\mathbf{v} = \mathbf{y} - \widehat{\mathbf{A}}\mathbf{x} = \mathbf{y} - \Re \left(\boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \right) \mathbf{x}. \quad (7.92)$$

Let us define $\boldsymbol{\Sigma}_{vv}$ and $\boldsymbol{\Sigma}_{vx}$ such that $\boldsymbol{\Sigma}_{vv} = \langle \mathbf{v}\mathbf{v}^H \rangle$ and $\boldsymbol{\Sigma}_{vx} = \langle \mathbf{v}\mathbf{x}^H \rangle$. They are obtained as

$$\boldsymbol{\Sigma}_{vv} = \langle (\mathbf{y} - \widehat{\mathbf{A}}\mathbf{x})(\mathbf{y} - \widehat{\mathbf{A}}\mathbf{x})^H \rangle = \boldsymbol{\Sigma}_{yy} - \widehat{\mathbf{A}}\boldsymbol{\Sigma}_{xy} - \boldsymbol{\Sigma}_{yx}\widehat{\mathbf{A}}^T + \widehat{\mathbf{A}}\boldsymbol{\Sigma}_{xx}\widehat{\mathbf{A}}^T, \quad (7.93)$$

and

$$\boldsymbol{\Sigma}_{vx} = \langle (\mathbf{y} - \widehat{\mathbf{A}}\mathbf{x})\mathbf{x}^H \rangle = \boldsymbol{\Sigma}_{yx} - \widehat{\mathbf{A}}\boldsymbol{\Sigma}_{xx}. \quad (7.94)$$

We just follow the arguments in Sect. 7.6.1, and derive an expression for the canonical magnitude coherence between \mathbf{v} and \mathbf{x} . Using complex-valued \mathbf{a} and \mathbf{b} , we define $\widehat{v} = \mathbf{a}^H \mathbf{v}$ and $\widehat{x} = \mathbf{b}^H \mathbf{x}$, the magnitude coherence between \widehat{v} and \widehat{x} , $|\phi_R|^2$,

is given by

$$|\phi_R|^2 = \frac{|\langle \hat{v} \hat{x}^* \rangle|^2}{\sqrt{\langle \hat{v} \hat{v}^* \rangle} \sqrt{\langle \hat{x} \hat{x}^* \rangle}} = \frac{|\mathbf{a}^H \boldsymbol{\Sigma}_{vx} \mathbf{b}|^2}{[\mathbf{a}^H \boldsymbol{\Sigma}_{vv} \mathbf{a}] [\mathbf{b}^H \boldsymbol{\Sigma}_{xx} \mathbf{b}]}.$$
 (7.95)

The canonical squared residual coherence $|\psi_R|^2$ is defined as the maximum of $|\phi_R|^2$, which is obtained by solving the optimization problem,

$$|\psi_R|^2 = \max_{\mathbf{a}, \mathbf{b}} |\mathbf{a}^H \boldsymbol{\Sigma}_{vx} \mathbf{b}|^2, \text{ subject to } \mathbf{a}^H \boldsymbol{\Sigma}_{vv} \mathbf{a} = 1 \text{ and } \mathbf{b}^H \boldsymbol{\Sigma}_{xx} \mathbf{b} = 1.$$
 (7.96)

This maximization problem is exactly the same as that in Eq. (7.60), and the solution of this maximization is known to be the maximum eigenvalue of the matrix $\boldsymbol{\Sigma}_{vv}^{-1} \boldsymbol{\Sigma}_{vx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{vx}^H$. That is, the canonical residual coherence is derived as

$$|\psi_R|^2 = \mathcal{S}_{\max}\{\boldsymbol{\Sigma}_{vv}^{-1} \boldsymbol{\Sigma}_{vx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{vx}^H\},$$
 (7.97)

where $\boldsymbol{\Sigma}_{vv}$ is derived using Eq. (7.93), and $\boldsymbol{\Sigma}_{vx}$ is derived using Eq. (7.94).

7.6.4 Computing Coherence When Each Voxel has Multiple Time Courses

When the source vector has x , y , and z components, most source imaging algorithms generate three time courses corresponding to the x , y , and z components at each voxel. A usual way is to obtain a single representative time course at each voxel and compute coherence using such representative time courses, as described in Sect. 7.2. Here, we describe an alternative method to compute voxel coherence when each voxel has multiple time courses [13].

The method is a straightforward application of the canonical coherence. Let us assume that the voxel time courses are expressed as in Eq. (7.1). The voxel spectra at the target and seed voxels are respectively denoted using 3×1 vectors \mathbf{x} and \mathbf{y} , which are given by

$$\mathbf{x}(f) = \begin{bmatrix} x_1(f) \\ x_2(f) \\ x_3(f) \end{bmatrix} \quad \text{and} \quad \mathbf{y}(f) = \begin{bmatrix} y_1(f) \\ y_2(f) \\ y_3(f) \end{bmatrix},$$
 (7.98)

where $x_1(f)$, $x_2(f)$, and $x_3(f)$ are the spectra obtained from the source time courses in the x , y and z directions at the target voxel, and $y_1(f)$, $y_2(f)$, and $y_3(f)$ are the spectra obtained from the source time courses in the x , y and z directions at the seed voxel. Using these \mathbf{x} and \mathbf{y} , we compute canonical magnitude coherence $|\psi|$ in Eq. (7.66) or $|\tilde{\psi}|$ in Eq. (7.68), and the canonical imaginary coherence ψ_I in Eq. (7.83) or $\tilde{\psi}_I$ in Eq. (7.85). This method is equivalent to determining the source

orientations at the seed and target voxels by simultaneously maximizing the canonical (magnitude/imaginary) coherence between these voxels.

7.7 Envelope Correlation and Related Connectivity Metrics

7.7.1 Envelope Correlation

Coherence is a metric that measures the phase relationship between two spectra, and naturally it is sensitive to phase jitters. This property becomes problematic when we try to estimate brain interactions at higher frequencies, such as gamma and high-gamma activities. This is because, at high frequencies, a small time jitter could cause a large phase jitter, and coherence may not be able to detect connectivity relationships. That is, the coherence may not be appropriate as a connectivity metric at high frequencies. The envelope-to-envelope correlation [14] is considered an appropriate metric for such cases.

To compute the envelope correlation, we first convert the seed and target time courses into their analytic signals, such that

$$\mathcal{A}[u(t)] = u(t) + \frac{i}{\pi} \int \frac{u(t')}{t - t'} dt'. \quad (7.99)$$

On the left-hand side of the equation above, $\mathcal{A}[\cdot]$ indicates an operator that creates an analytic signal of the real-valued time signal in the parentheses. Let us define analytic signals from the seed- and target-voxel time courses $u_S(t)$ and $u_T(t)$, such that

$$\mathcal{A}[u_S(t)] = A_S(t)e^{i\theta_S(t)}, \quad (7.100)$$

$$\mathcal{A}[u_T(t)] = A_T(t)e^{i\theta_T(t)}, \quad (7.101)$$

where $A_S(t)$ and $A_T(t)$ are the amplitudes of the seed and target analytic signals, and $\theta_S(t)$ and $\theta_T(t)$ are their instantaneous phases. The envelope correlation is the correlation between the amplitudes, $A_S(t)$ and $A_T(t)$, and is computed such that

$$\Theta = \frac{\sum_{j=1}^K A_T(t_j)A_S(t_j)}{\sqrt{\left[\sum_{j=1}^K A_T(t_j)^2\right]\left[\sum_{j=1}^K A_S(t_j)^2\right]}}. \quad (7.102)$$

It is obvious in the above equation that common interferences such as the algorithm leakage cause spurious correlation, and the seed blur should exist in an image of envelope correlation.

7.7.2 Residual Envelope Correlation

Residual envelope correlation, which is free from the problem of seed blur, is proposed in [15]. To compute the residual envelope correlation, we first compute the residual time course $u_R(t)$ using

$$u_R(t) = \int_{-\infty}^{\infty} v(f) e^{i2\pi ft} df. \quad (7.103)$$

Here $v(f)$ is the residual spectrum obtained using Eq.(7.48). Since $v(f)$ is defined only for $f \geq 0$, we create the residual spectrum for $f < 0$, such that

$$v(-f) = v(f)^*,$$

where the superscript * indicates the complex conjugation. Since the Hermitian symmetry holds for $v(f)$, $u_R(t)$ is guaranteed to be a real-valued time course. The envelope of $u_R(t)$ is computed using

$$\mathcal{A}[u_R(t)] = A_R(t) e^{i\theta_R(t)}, \quad (7.104)$$

where $A_R(t)$ is the envelope of $u_R(t)$ called the residual envelope. Once $A_R(t)$ is computed, the residual envelope correlation Θ_R is computed using

$$\Theta_R = \frac{\sum_{j=1}^K A_R(t_j) A_S(t_j)}{\sqrt{\left[\sum_{j=1}^K A_R(t_j)^2\right] \left[\sum_{j=1}^K A_S(t_j)^2\right]}}. \quad (7.105)$$

Since in the residual time course the seed signal is regressed out, the residual envelope correlation is free from the spurious correlation caused by the algorithm leakage. Note that a method similar to the one mentioned here was reported in [16].

7.7.3 Envelope Coherence

The coherence can be computed between the envelope time courses $A_T(t)$ and $A_S(t)$. Let us define the spectra obtained from the envelope $A_S(t)$ and $A_T(t)$ as $\Sigma_S(f)$ and $\Sigma_T(f)$, respectively. The complex envelope coherence is computed using

$$\phi_E(f) = \frac{\langle \Sigma_T(f) \Sigma_S^*(f) \rangle}{\sqrt{\langle |\Sigma_T(f)|^2 \rangle \langle |\Sigma_S(f)|^2 \rangle}}. \quad (7.106)$$

The imaginary and corrected imaginary coherences are obtained as $\Im(\phi_E(f))$ and $\Im(\phi_E(f))/\sqrt{(1 - \Re(\phi_E(f))^2)}$.

7.8 Statistical Thresholding of Coherence Images

In practical applications, we need to assess the statistical significance of the obtained coherence images. The surrogate data method [17, 18] has been used for this assessment. In this method, the surrogate voxel spectra are created by multiplying random phases with the original voxel spectra. The surrogate spectra for the seed and target voxels, $\tilde{\sigma}_S(f)$ and $\tilde{\sigma}_T(f)$, are expressed as

$$\tilde{\sigma}_S = \widehat{\sigma}_S e^{i2\pi\delta_S} \quad \text{and} \quad \tilde{\sigma}_T = \widehat{\sigma}_T e^{i2\pi\delta_T}, \quad (7.107)$$

where δ_S and δ_T are uniform random numbers between 0 and 1, and the explicit notation of (f) is again omitted for simplicity. Note that the surrogate spectra have the same original power spectra but the phase relationship is destroyed by multiplying the random phases to the original spectra.

Any coherence-based metric, such as the magnitude/imaginary coherence or corrected imaginary coherence, can be computed using the surrogate spectra, $\tilde{\sigma}_S$ and $\tilde{\sigma}_T$. The metric computed using the surrogate spectra is denoted ω . As an example, the imaginary coherence is computed using the surrogate spectra. It is expressed as³

$$\omega = \frac{|\Im((\tilde{\sigma}_T \tilde{\sigma}_S^*))|}{\sqrt{\langle |\tilde{\sigma}_T|^2 \rangle \langle |\tilde{\sigma}_S|^2 \rangle}} = \frac{|\Im((\widehat{\sigma}_T \widehat{\sigma}_S^* e^{i2\pi\Delta\delta}))|}{\sqrt{\langle |\widehat{\sigma}_T|^2 \rangle \langle |\widehat{\sigma}_S|^2 \rangle}}, \quad (7.108)$$

where $\Delta\delta = \delta_T - \delta_S$. The generation of ω is repeated B times, and a total of B values of ω , denoted $\omega^1, \omega^2, \dots, \omega^B$, are obtained. These $\omega^1, \dots, \omega^B$ can form an empirical null distribution at each voxel.

We could derive a voxel-by-voxel statistical threshold using this empirical null distribution. However, the statistical threshold derived in this manner does not take the multiple comparisons into account and it generally leads to a situation in which many false-positive voxels arise, i.e., many voxels that contain no brain interaction are found to be interacting. To avoid this problem, the statistical significance is determined using a procedure that takes multiple comparisons into account. For this purpose, we can use the maximal statistics⁴ [19, 20].

To utilize maximum statistics, the values ω^ℓ ($\ell = 1, \dots, B$) are first standardized and converted into pseudo- t values, such that

³ The absolute value of the imaginary coherence is usually computed, because its sign has no meaning in expressing the connectivity.

⁴ We can apply another method such as the false discovery rate to this multiple comparison problem.

$$T^\ell = \frac{\omega^\ell - \langle \omega \rangle}{\sigma_\omega}, \quad (7.109)$$

where $\langle \omega \rangle$ and σ_ω^2 are the average and the variance of ω^ℓ ($\ell = 1, \dots, B$). Since these $\langle \omega \rangle$ and σ_ω^2 are obtained at each voxel, the values at the j th voxel are denoted $\langle \omega(j) \rangle$ and $\sigma_\omega^2(j)$. The maximum value of T^ℓ obtained at the j th voxel is denoted $T^{\max}(j)$. Defining a total number of voxels as N_V , we have $T^{\max}(1), \dots, T^{\max}(N_V)$ to form a null distribution. We then sort these values in an increasing order:

$$T^{\max}(\tilde{1}) \leq T^{\max}(\tilde{2}) \leq \dots \leq T^{\max}(\tilde{N}_V),$$

where $T^{\max}(\tilde{k})$ is the k th minimum value.

We set the level of the statistical significance to α , and choose $T^{\max}(\tilde{p})$ where $\tilde{p} = \lceil \alpha N_V \rceil$ and $\lceil \alpha N_V \rceil$ indicates the maximum integer not greater than αN_V . The threshold value for the j th voxel, $\omega^{\text{th}}(j)$, is finally derived as

$$\omega^{\text{th}}(j) = T^{\max}(\tilde{p})\sigma_\omega(j) + \langle \omega(j) \rangle. \quad (7.110)$$

At the j th voxel, we evaluate the statistical significance of the imaginary coherence value by comparing it with $\omega^{\text{th}}(j)$. When the metric value is greater than $\omega^{\text{th}}(j)$, it is considered to be statistically significant; if not, it is considered to be statistically insignificant.

7.9 Mean Imaginary Coherence (MIC) Mapping

Guggisberg et al. [3] have proposed to compute a metric called the mean imaginary coherence. Defining the coherence computed between the j th and k th voxels as $\hat{\phi}_{j,k}(f)$, the mean imaginary coherence for the j th voxel, $\mathcal{M}_j(f)$, is obtained using

$$\mathcal{M}_j(f) = \tanh \left[\frac{1}{N_V} \sum_{k=1}^{N_V} \tanh^{-1} (|\Im(\hat{\phi}_{j,k}(f))|) \right]. \quad (7.111)$$

On the right-hand side, the absolute value of the imaginary coherence $|\Im(\hat{\phi}_{j,k}(f))|$ is averaged across all voxel connections. In Eq. (7.111),

$$z = \tanh^{-1}(r) = \frac{1}{2} \log \left(\frac{1+r}{1-r} \right) \quad \text{and} \quad r = \tanh(z) = \frac{e^{2z} - 1}{e^{2z} + 1}$$

are the inverse hyperbolic and hyperbolic functions, respectively. The idea of using these functions is to average the voxel coherence values in the Fisher's Z-transform domain. We may use the corrected imaginary coherence instead of using the imaginary coherence in Eq. (7.111).

The mean imaginary coherence $\mathcal{M}_j(f)$ is hypothesized to express the communication capability (healthiness) of the brain tissue at the j th voxel location. This hypothesis has been tested by applying the mean imaginary coherence (MIC) mapping to various types of clinical data, and results positively confirming the validity of the hypothesis have been obtained. Such investigations include the MIC mapping for stroke-recovery patients [21], and patients with traumatic brain lesions [22]. In these applications, mapping of $\mathcal{M}_j(f)$ can predict the degree of a patient's recovery from stroke or brain injury. The MIC mapping has also been applied to MEG data from patients with schizophrenia [23] and Alzheimer's disease [24], and the values of $\mathcal{M}_j(f)$ are found to be significantly correlated with patient symptom scores at particular brain areas.

7.10 Numerical Examples

Computer simulation was performed to illustrate results of our arguments in this chapter. A sensor alignment of the 275-sensor array from the OmegaTM (VMS Medtech, Coquitlam, Canada) neuromagnetometer was used. The coordinate system and source-sensor configuration used in the computer simulation are depicted in Fig. 5.1 in Chap. 5. A vertical plane ($x = 0$ cm) was assumed at the middle of the whole-head sensor array, and three sources were assumed to exist on this plane.

Multiple-trial measurements were simulated, in which a total of 120 trials were generated with each trial consisting of 600 time points. The data were assumed to be collected with 2 ms sampling. We first conducted numerical experiments with

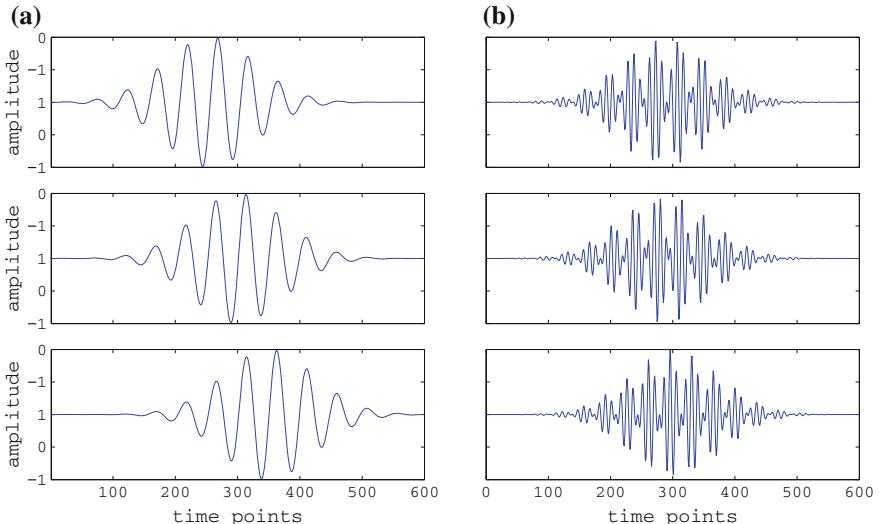


Fig. 7.2 Examples of the time courses **a** for the alpha-band experiments and **b** for the gamma-band experiments. The *top*, *middle*, and *bottom* panels, respectively, show the time courses for the first, second, and third sources

alpha-band activity. The time courses of the three sources for the first trial are shown in Fig. 7.2a. The time courses had trial-to-trial time jitters. The time jitters for the three time courses were generated using Gaussian random numbers with the same standard deviation of 20 time points. These time courses were projected onto the sensor space using the sensor lead field to obtain the signal magnetic recordings. The simulated sensor recordings were computed by adding spontaneous MEG data to the signal magnetic recordings. The signal-to-interference ratio was set to 0.5.

The voxel time courses were reconstructed using the vector-type narrow-band adaptive beamformer described in Chap. 3. The data-covariance matrix was tuned to the alpha frequency band. Reconstructed source power images on the plane $x = 0$ cm are shown in Fig. 7.3a. The three sources are clearly resolved. Since the spherical homogeneous conductor [25] was used for forward modeling, the vector adaptive beamformer reconstructed two time courses corresponding to the two tangential directions. We estimated a single representative time course at each voxel using the method described in Sect. 7.2, and computed voxel spectra. Coherence images were then computed with respect to the seed voxel, which was set at the second

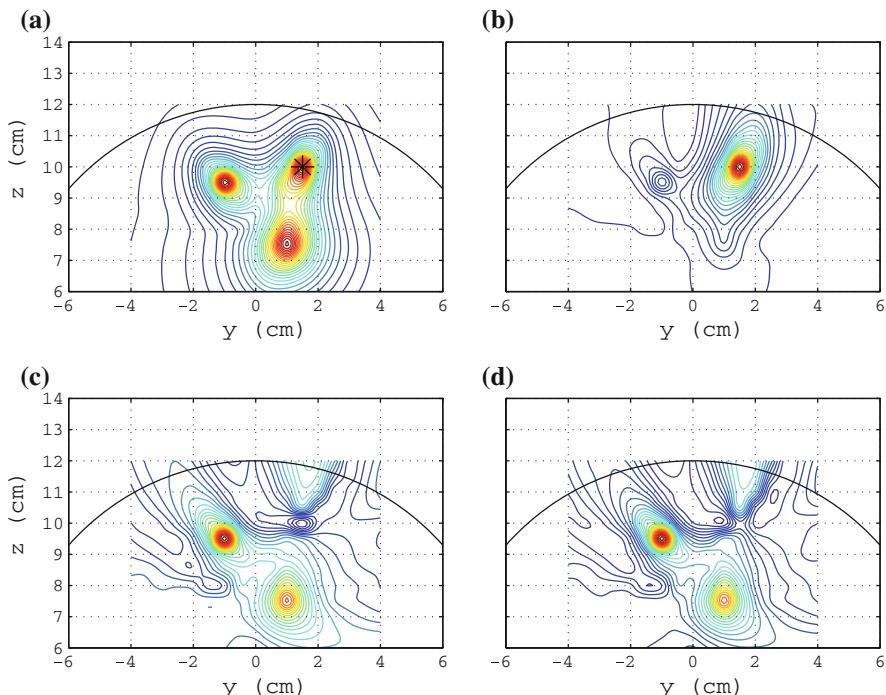


Fig. 7.3 Results of imaging the alpha-band voxel coherence on the plane $x = 0$ cm. **a** Results of source reconstruction. The asterisk shows the location of the seed voxel. **b** Magnitude coherence image. **c** Imaginary coherence image. **d** Corrected imaginary coherence image

source location. The magnitude coherence image, imaginary coherence image, and corrected imaginary coherence image are respectively shown in Fig. 7.3b–d.

In the magnitude coherence image (Fig. 7.3b), the seed blur dominates and obscures the other sources. On the contrary, in the imaginary coherence image (Fig. 7.3c), and in the corrected imaginary coherence image (7.3d), the intensity of the seed blur is much reduced and the two sources that interact with the second source can clearly be observed. Also, the imaginary and the corrected imaginary coherence images are very similar, because the magnitude coherence is as small as 0.26–0.35 in this computer simulation.

We next implemented the method described in Sect. 7.6.4 in which a single representative voxel time course is not computed, but instead, the canonical coherence is directly computed by using multiple voxel time courses. The image of canonical imaginary coherence ψ_I in Eq. (7.83) is shown in Fig. 7.4a. The image of mutual-information-based canonical imaginary coherence, $\tilde{\psi}_I$ in Eq. (7.85) is shown in Fig. 7.4b. The canonical magnitude coherence $|\psi|$ in Eq. (7.66) is shown in Fig. 7.4c. An image of mutual-information-based canonical magnitude coherence

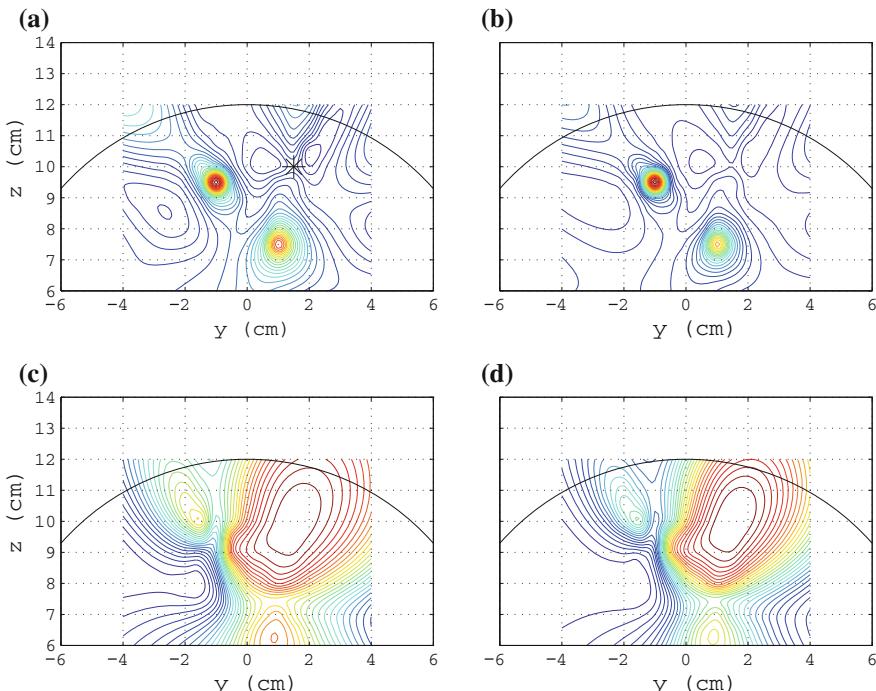


Fig. 7.4 Results of imaging canonical coherence on the plane $x = 0\text{cm}$. **a** Image of canonical imaginary coherence ψ_I in Eq. (7.83). The asterisk shows the location of the seed voxel. **b** Image of mutual-information-based canonical imaginary coherence, $\tilde{\psi}_I$ in Eq. (7.85). **c** Image of canonical magnitude coherence $|\psi|$ in Eq. (7.66). **d** Image of mutual-information-based canonical magnitude coherence $|\tilde{\psi}|$ in Eq. (7.68)

$|\tilde{\psi}|$ in Eq.(7.68) is shown in Fig. 7.4d. In these results, imaginary-coherence-based images ((a) and (b)), clearly detect the two interacting sources but the seed blur dominates in the magnitude-coherence-based images ((c) and (d)).

We carried out the numerical experiment simulating gamma-band signals with theta-band envelopes. The time courses of the three sources for the first trial are shown in Fig. 7.2b. The time courses had trial-to-trial time jitters generated using Gaussian random numbers with the same standard deviation of 20 time points.

The voxel time courses were estimated using the narrow-band adaptive beamformer with a data-covariance tuned to the gamma frequency band, and the coherence images were computed. The seed was set at the second source location. The coherence images are shown in Fig. 7.5a. Here, the magnitude, imaginary, and corrected imaginary coherence images are, respectively, shown in the top, middle, and bottom panels. In the magnitude coherence image, the seed blur dominates and only the seed source is detected. However, since the time jitter created a large phase jitter for the gamma-band signals, neither the imaginary nor the corrected imaginary coherence image contain meaningful information on the source connectivity.

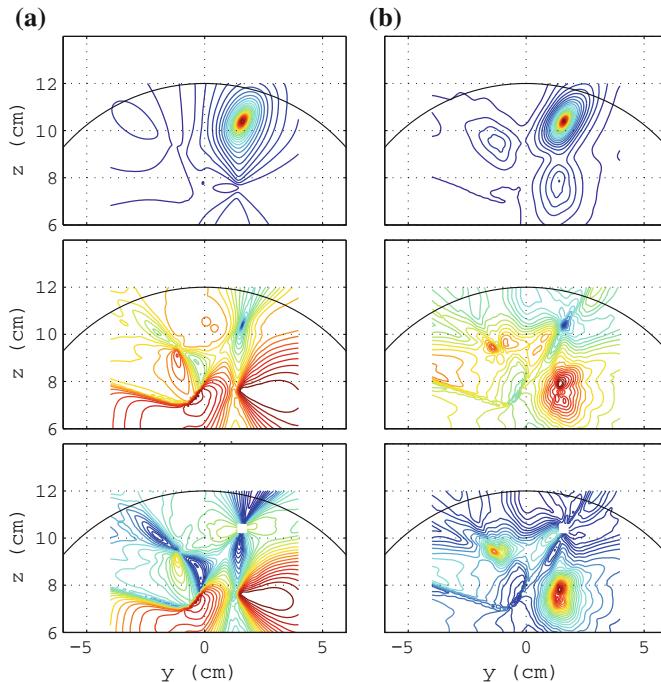


Fig. 7.5 Results of imaging the gamma-band source coherence on the plane $x = 0\text{ cm}$. The seed was set at the second source location. The coherence images were computed using the gamma-band voxel time courses. The *top, middle, and bottom panels* in **a**, respectively, show the magnitude, imaginary, and corrected imaginary coherence images. The *top, middle, and bottom panels* in **b**, respectively, show the magnitude, imaginary, and corrected imaginary envelope coherence images

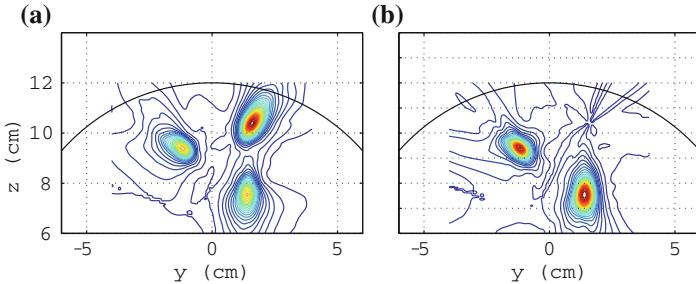


Fig. 7.6 Results of imaging envelope correlation on the plane $x = 0\text{ cm}$. The seed was set at the second source location. **a** Image of envelope correlation. **b** Image of residual envelope correlation

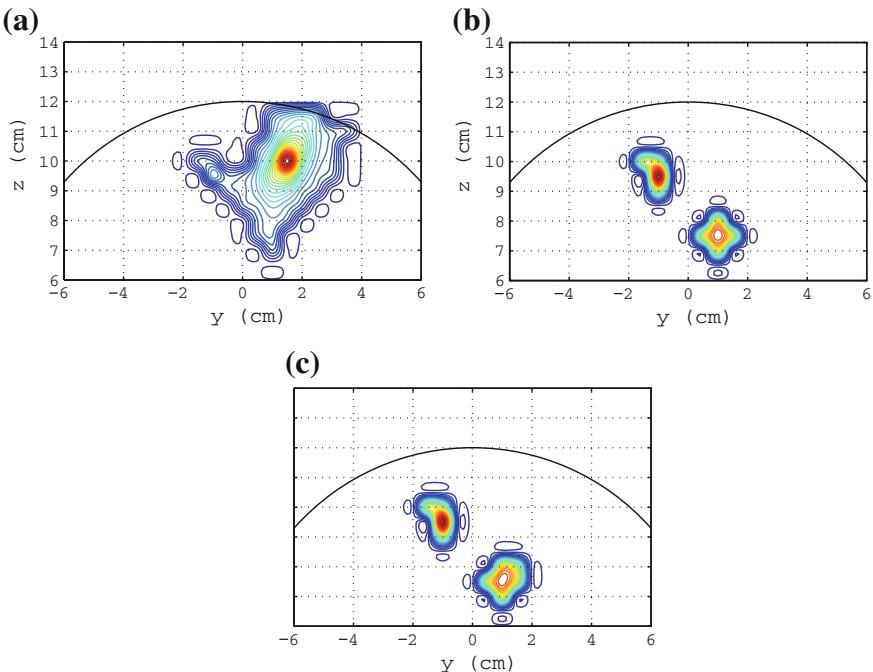


Fig. 7.7 Results of statistical thresholding of the alpha-band source coherence images. **a** Magnitude coherence image. **b** Imaginary coherence image. **c** Corrected imaginary coherence image. Statistical thresholding using the surrogate data method described in Sect. 7.8 was applied with the number of surrogate data set to 200 and the level of the statistical significance set to 0.99

We then computed the Hilbert envelopes of the reconstructed voxel time courses, and computed the envelope-to-envelope coherence. The results of the envelope-to-envelope coherence are shown in Fig. 7.5b in which the top, middle, and bottom panels show the magnitude, imaginary, and corrected imaginary coherence images, respectively. The magnitude coherence image detects only the seed source. In the imaginary coherence image, the interacting sources are not very clearly detected, but in the corrected imaginary coherence image, the first and the third sources are detected with reasonable clarity. In these results, the difference between the imaginary coherence and corrected imaginary coherence images is significantly large because the magnitude envelope coherence between the source time courses is as large as 0.7.

We computed the image of the envelope-to-envelope correlation, and the results are shown in Fig. 7.6. In this figure, the image of the envelope correlation is shown in (a) and the image of the residual envelope correlation is in (b). In both images, the first and the third sources that were interacting with the second (seed) source can be observed. While the original envelope correlation image contains seed blur, the residual envelope correlation image is free from such spurious activity.

We finally show the results of statistical thresholding method described in Sect. 7.8. The method was applied to the alpha-band coherence imaging results in Fig. 7.3. The thresholded images are shown in Fig. 7.7. The spurious baseline activity existing in the unthresholded images in Fig. 7.3 are removed, and it is much easier to interpret the results in the thresholded images in Fig. 7.7.

References

1. J.-M. Schoffelen, J. Gross, Source connectivity analysis with MEG and EEG. *Hum. Brain Mapp.* **30**, 1857–1865 (2009)
2. J. Gross, J. Kujara, M. Hämäläinen, L. Timmermann, A. Schnitzler, R. Salmelin, Dynamic imaging of coherent sources: studying neural interactions in the human brain. *Proc. Natl. Acad. Sci. U.S.A.* **98**, 694–699 (2001)
3. A.G. Guggisberg, S.M. Honma, A.M. Findlay, S.S. Dalal, H.E. Kirsch, M.S. Berger, S.S. Nagarajan, Mapping functional connectivity in patients with brain lesions. *Ann. Neurol.* **63**, 193–203 (2007)
4. P. Belardinelli, L. Ciancetta, M. Staudt, V. Pizzella, A. Londeg, N.B.G.L. Romani, C. Braun, Cerebro-muscular and cerebro-cerebral coherence in patients with pre- and perinatally acquired unilateral brain lesions. *NeuroImage* **37**, 1301–1314 (2007)
5. W.H.R. Miltner, C. Braun, M. Arnold, H. Witte, E. Taub, Coherence of gamma-band EEG activity as a basis for associative learning. *Nature* **397**, 434–436 (1999)
6. K. Sekihara, S.S. Nagarajan, *Adaptive Spatial Filters for Electromagnetic Brain Imaging* (Springer, Berlin, 2008)
7. K. Sekihara, J.P. Owen, S. Trisno, S.S. Nagarajan, Removal of spurious coherence in MEG source-space coherence analysis. *IEEE Trans. Biomed. Eng.* **58**, 3121–3129 (2011)
8. G. Nolte, O.B.L. Wheaton, Z. Mari, S. Vorbach, M. Hallett, Identifying true brain interaction from EEG data using the imaginary part of coherency. *Clin. Neurophysiol.* **115**, 2292–2307 (2004)
9. R.D. Pascual-Marqui, Instantaneous and lagged measurements of linear and nonlinear dependence between groups of multivariate time series: frequency decomposition (2007). arXiv preprint [arXiv:0711.1455](https://arxiv.org/abs/0711.1455)

10. E.J. Hannan, *Multiple Time Series*, vol. 38 (Wiley, New York, 2009)
11. D.R. Brillinger, *Time Series: Data Analysis and Theory*, vol. 36 (Siam, Philadelphia, 2001)
12. A. Ewald, L. Marzetti, F. Zappasodi, F.C. Meinecke, G. Nolte, Estimating true brain connectivity from EEG/MEG data invariant to linear and static transformations in sensor space. *NeuroImage* **60**(1), 476–488 (2012)
13. F. Shahbazi Avarvand, A. Ewald, G. Nolte, Localizing true brain interactions from EEG and MEG data with subspace methods and modified beamformers. *Comput. Math. Methods Med.* **2012**, 402341 (2012)
14. A. Bruns, R. Eckhorn, Task-related coupling from high- to low-frequency signals among visual cortical areas in human subdural recordings. *Int. J. Psychophysiol.* **51**, 97–116 (2004)
15. K. Sekihara, S.S. Nagarajan, Residual coherence and residual envelope correlation in MEG/EEG source-space connectivity analysis, in *Conference of Proceedings of the IEEE Engineering in Medicine and Biology Society*, pp. 4417–7 (2013)
16. J.F. Hipp, D.J. Hawellek, M. Corbetta, M. Siegel, A.K. Engel, Large-scale cortical correlation structure of spontaneous oscillatory activity. *Nat. Neurosci.* **15**(6), 884–890 (2012)
17. J. Theiler, S. Eubank, A. Longtin, B. Galdrikian, J.D. Farmer, Testing for nonlinearity in time series: the method of surrogate data. *Phys. D* **58**, 77–94 (1992)
18. L. Faes, G.D. Pinna, A. Porta, R. Maestri, G. Nollo, Surrogate data analysis for assessing the significance of the coherence function. *IEEE Trans. Biomed. Eng.* **51**, 1156–1166 (2004)
19. Y. Benjamini, Y. Hochberg, Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. B* **57**, 289–300 (1995)
20. T.E. Nichols, S. Hayasaka, Controlling the familywise error rate in functional neuroimaging: A_j neural eng. comparative review. *Stat. Methods Med. Res.* **12**, 419–446 (2003)
21. K.P. Westlake, L.B. Hinkley, M. Bucci, A.G. Guggisberg, A.M. Findlay, R.G. Henry, S.S. Nagarajan, N. Byl, Resting state alpha-band functional connectivity and recovery after stroke. *Exp. Neurol.* **237**(1), 160–169 (2012)
22. P.E. Tarapore, A.M. Findlay, S.C. LaHue, H. Lee, S.M. Honma, D. Mizuiri, T.L. Luks, G.T. Manley, S.S. Nagarajan, P. Mukherjee, Resting state magnetoencephalography functional connectivity in traumatic brain injury: clinical article. *J. Neurosurg.* **118**(6), 1306–1316 (2013)
23. L.B. Hinkley, J.P. Owen, M. Fisher, A.M. Findlay, S. Vinogradov, S.S. Nagarajan, Cognitive impairments in schizophrenia as assessed through activation and connectivity measures of magnetoencephalography (MEG) data. *Front. Hum. Neurosci.* **3**, 73 (2009)
24. K.G. Ranasinghe, L.B. Hinkley, A.J. Beagle, D. Mizuiri, A.F. Dowling, S.M. Honma, M.M. Finucane, C. Scherling, B.L. Miller, S.S. Nagarajan et al., Regional functional connectivity predicts distinct cognitive impairments in Alzheimers disease spectrum. *NeuroImage: Clin.* **5**, 385–395 (2014)
25. J. Sarvas, Basic mathematical and electromagnetic concepts of the biomagnetic inverse problem. *Phys. Med. Biol.* **32**, 11–22 (1987)

Chapter 8

Estimation of Causal Networks: Source-Space Causality Analysis

8.1 Introduction

This chapter reviews the methodology for estimating causal relationships among cortical activities in MEG/EEG source space analysis. The source space causality analysis computes some types of causality measures using the estimated time series of the source activities of interest. Commonly-used measures are the Granger causality related measures, which are based on the MVAR modeling of the voxel time series. We first describe these Granger causality based measures, including the original definition in the temporal domain [1, 2], Geweke's extension to the spectral domain [3, 4], and related measures such as the directed transfer function (DTF) [5] and the partial directed coherence (PDC) [6]. The transfer entropy [7], which is a representative non-parametric measure, is also discussed, with a proof of its equivalence to the Granger causality under the Gaussianity assumption. The last section of this chapter describes methods of estimating the MVAR coefficients. Here we present a conventional least-square-based algorithm and a method based on the sparse Bayesian inference for this MVAR estimation problem. The sparse Bayesian algorithm outperforms the conventional method when the sensor data contains a large amount of interferences.

8.2 Multivariate Vector Autoregressive (MVAR) Process

8.2.1 MVAR Modeling of Time Series

The Granger-causality measures make use of multivariate vector auto regressive (MVAR) process of the source time series. We first explain the modeling of time series using the MVAR process. A general approach to the source-space causality analysis first chooses a relatively small number of voxels corresponding to the source activities of interest, and the causal relationships among time series of the selected voxels are estimated.

Denoting the number of selected voxels q , we express the q -channel voxel time series using the vector $\mathbf{y}(t)$: $\mathbf{y}(t) = [y_1(t), \dots, y_q(t)]^T$ where $y_j(t)$ is the time series of the j th selected voxel at time t . Here the time t is expressed using a unit-less value. We impose multivariate vector autoregressive (MVAR) modeling on the time series $\mathbf{y}(t)$, such that

$$\mathbf{y}(t) = \sum_{p=1}^P \mathbf{A}(p) \mathbf{y}(t-p) + \mathbf{e}(t). \quad (8.1)$$

Here, $\mathbf{A}(p)$ is the AR coefficient matrix, P is the model order, and $\mathbf{e}(t)$ is the residual vector. The MVAR process is expressed in the frequency domain. By computing the Fourier transform of Eq. (8.1), we get

$$\mathbf{y}(f) = \sum_{p=1}^P \mathbf{A}(p) e^{-2\pi i p f} \mathbf{y}(f) + \mathbf{e}(f), \quad (8.2)$$

where the Fourier transforms of $\mathbf{y}(t)$ and $\mathbf{e}(t)$ are expressed in $\mathbf{y}(f)$ and $\mathbf{e}(f)$. We here use the relationship,

$$\int \mathbf{y}(t-p) \exp(-2\pi i f t) dt = e^{-2\pi i p f} \mathbf{y}(f). \quad (8.3)$$

Equation (8.2) is also expressed as

$$\left[\mathbf{I} - \sum_{p=1}^P \mathbf{A}(p) e^{-2\pi i p f} \right] \mathbf{y}(f) = \mathbf{e}(f). \quad (8.4)$$

Defining a $q \times q$ matrix $\bar{\mathbf{A}}(f)$ such that

$$\bar{\mathbf{A}}(f) = \mathbf{I} - \sum_{p=1}^P \mathbf{A}(p) e^{-2\pi i p f}, \quad (8.5)$$

we can obtain

$$\bar{\mathbf{A}}(f) \mathbf{y}(f) = \mathbf{e}(f). \quad (8.6)$$

Also, defining $\mathbf{H}(f) = \bar{\mathbf{A}}(f)^{-1}$, the relationship

$$\mathbf{y}(f) = \mathbf{H}(f) \mathbf{e}(f) \quad (8.7)$$

can be obtained. According to the equation above, $\mathbf{e}(f)$ and $\mathbf{y}(f)$ can be interpreted as the input and the output of a linear-system whose transfer function is $\mathbf{H}(f)$.

8.2.2 Coherence and Partial Coherence of the MVAR Process

Using Eq. (8.7), we can derive the relationship

$$\mathbf{S}(f) = \mathbf{H}(f)\boldsymbol{\Sigma}\mathbf{H}^H(f), \quad (8.8)$$

where the superscript H indicates the Hermitian transpose.¹ In the equation above, $\mathbf{S}(f)$ is the cross spectrum matrix, which is given by $\mathbf{S}(f) = \langle \mathbf{y}(f)\mathbf{y}^H(f) \rangle$ where $\langle \cdot \rangle$ indicates the ensemble average. Also, $\boldsymbol{\Sigma}$ is the covariance matrix of the residual, which is equal to $\boldsymbol{\Sigma} = \langle \mathbf{e}(f)\mathbf{e}^H(f) \rangle$. Using these definitions, the coherence between the time series of the j th and the k th channels is expressed as

$$\phi_{j,k}(f) = \frac{S_{j,k}}{\sqrt{S_{j,j} S_{k,k}}} = \frac{\mathbf{h}_j^H \boldsymbol{\Sigma} \mathbf{h}_k}{\sqrt{[\mathbf{h}_j^H \boldsymbol{\Sigma} \mathbf{h}_j][\mathbf{h}_k^H \boldsymbol{\Sigma} \mathbf{h}_k]}}, \quad (8.9)$$

where \mathbf{h}_j is the j th column of \mathbf{H}^H such that $\mathbf{H}^H = [\mathbf{h}_1, \dots, \mathbf{h}_q]$.

The partial coherence between the j th and the k th channels, $\kappa_{j,k}$, is expressed as [8, 9]

$$\kappa_{j,k}(f) = \frac{\mathbf{M}_{j,k}}{\sqrt{\mathbf{M}_{j,j} \mathbf{M}_{k,k}}}, \quad (8.10)$$

where $\mathbf{M}_{j,k}$ is the minor of the matrix $\mathbf{S}(f)$; the minor is a determinant value of a matrix formed by the j th row and the k th column removed from $\mathbf{S}(f)$. The partial coherence $\kappa_{j,k}(f)$ is a measure of the coherence between the time series of the j th and the k th channels where the influence of other channels is removed. In other words, it expresses the direct interaction between these two channels.

Let us derive a convenient formula for computing the partial coherence. We use the relationship called Cramer's rule [10]:

$$[\mathbf{S}^{-1}(f)]_{j,k} = \frac{(-1)^{j+k} \mathbf{M}_{k,j}}{|\mathbf{S}(f)|},$$

where $[\mathbf{S}^{-1}(f)]_{j,k}$ indicates the (j, k) component of the matrix $\mathbf{S}^{-1}(f)$. Considering the fact that $\mathbf{S}(f)$ is a positive semidefinite Hermitian matrix, the equation above can be changed to

$$\mathbf{M}_{j,k} = (-1)^{j+k} [\mathbf{S}^{-1}(f)]_{j,k} |\mathbf{S}(f)|. \quad (8.11)$$

¹ The Hermitian transpose is the matrix transpose with the complex conjugation.

We then rewrite Eq.(8.10) into

$$\begin{aligned}\kappa_{j,k}(f) &= \frac{(-1)^{j+k} [\mathbf{S}^{-1}(f)]_{j,k} |\mathbf{S}(f)|}{\sqrt{(-1)^{2j} [\mathbf{S}^{-1}(f)]_{j,j} |\mathbf{S}(f)| (-1)^{2k} [\mathbf{S}^{-1}(f)]_{k,k} |\mathbf{S}(f)|}} \\ &= \frac{[\mathbf{S}^{-1}(f)]_{j,k}}{\sqrt{[\mathbf{S}^{-1}(f)]_{j,j} [\mathbf{S}^{-1}(f)]_{k,k}}}.\end{aligned}\quad (8.12)$$

On the other hand, using Eq.(8.8), we have the relationship

$$\begin{aligned}\mathbf{S}^{-1}(f) &= [\mathbf{H}(f) \boldsymbol{\Sigma} \mathbf{H}^H(f)]^{-1} \\ &= [\mathbf{H}^{-1}(f)]^H \boldsymbol{\Sigma}^{-1} \mathbf{H}^{-1}(f) = \bar{\mathbf{A}}^H(f) \boldsymbol{\Sigma}^{-1} \bar{\mathbf{A}}(f).\end{aligned}\quad (8.13)$$

The k th column vector of $\bar{\mathbf{A}}(f)$ is denoted $\bar{\mathbf{a}}_k$, i.e., $\bar{\mathbf{A}}(f) = [\bar{\mathbf{a}}_1, \dots, \bar{\mathbf{a}}_q]$. We can then obtain

$$[\mathbf{S}^{-1}(f)]_{j,k} = \bar{\mathbf{a}}_j^H \boldsymbol{\Sigma}^{-1} \bar{\mathbf{a}}_k,$$

and thus derive

$$\kappa_{j,k}(f) = \frac{\bar{\mathbf{a}}_j^H \boldsymbol{\Sigma}^{-1} \bar{\mathbf{a}}_k}{\sqrt{[\bar{\mathbf{a}}_j^H \boldsymbol{\Sigma}^{-1} \bar{\mathbf{a}}_j][\bar{\mathbf{a}}_k^H \boldsymbol{\Sigma}^{-1} \bar{\mathbf{a}}_k]}}.\quad (8.14)$$

The equation above is used for deriving partial directed coherence in Sect. 8.5.3.

8.3 Time-Domain Granger Causality

8.3.1 Granger Causality for a Bivariate Process

Let us first consider bivariate cases in which only a pair of two time series $y_1(t)$ and $y_2(t)$ are considered. In the first case, $y_1(t)$ and $y_2(t)$ are determined by using only their own past values, i.e., the following relationships hold:

$$y_1(t) = \sum_{p=1}^{\infty} A_{1,1}(p) y_1(t-p) + e_1(t),\quad (8.15)$$

$$y_2(t) = \sum_{p=1}^{\infty} A_{2,2}(p) y_2(t-p) + e_2(t).\quad (8.16)$$

We then consider the second case in which the following relationships hold:

$$y_1(t) = \sum_{p=1}^{\infty} A_{1,1}(p)y_1(t-p) + \sum_{p=1}^{\infty} A_{1,2}(p)y_2(t-p) + \epsilon_1(t), \quad (8.17)$$

$$y_2(t) = \sum_{p=1}^{\infty} A_{2,1}(p)y_1(t-p) + \sum_{p=1}^{\infty} A_{2,2}(p)y_2(t-p) + \epsilon_2(t). \quad (8.18)$$

The difference between the first and the second cases above is that in the second case, time series $y_1(t)$ and $y_2(t)$ are determined by the past values of both the time series of $y_1(t)$ and $y_2(t)$, while in the first case, the time series $y_1(t)$ is determined only by the past values of $y_1(t)$ and the time series $y_2(t)$ is determined only by the past values of $y_2(t)$.

We define the variances of the residual terms appearing above, such that

$$V(\epsilon_1(t)) = \Sigma_{\epsilon}^1, \quad (8.19)$$

$$V(\epsilon_2(t)) = \Sigma_{\epsilon}^2, \quad (8.20)$$

$$V(\epsilon_1(t)) = \Sigma_{\epsilon}^1, \quad (8.21)$$

$$V(\epsilon_2(t)) = \Sigma_{\epsilon}^2, \quad (8.22)$$

where $V(\cdot)$ indicates the variance. Then, the Granger causality, $\mathcal{G}_{2 \rightarrow 1}$, is defined as

$$\mathcal{G}_{2 \rightarrow 1} = \log \left[\frac{\Sigma_{\epsilon}^1}{\Sigma_{\epsilon}^2} \right]. \quad (8.23)$$

The meaning of $\mathcal{G}_{2 \rightarrow 1}$ is that, if the residual variance is reduced by using both the past values of $y_1(t)$ and $y_2(t)$ when estimating the current value of $y_1(t)$, the past of $y_2(t)$ affects the current value of $y_1(t)$. In such a case, we can conclude that there is an information flow from the time series $y_2(t)$ to the time series $y_1(t)$. The above $\mathcal{G}_{2 \rightarrow 1}$ can quantify the amount of this information flow. Similarly, we can define $\mathcal{G}_{1 \rightarrow 2}$ such that

$$\mathcal{G}_{1 \rightarrow 2} = \log \left[\frac{\Sigma_{\epsilon}^2}{\Sigma_{\epsilon}^1} \right]. \quad (8.24)$$

This $\mathcal{G}_{1 \rightarrow 2}$ expresses the information flow from the time series $y_1(t)$ to the time series $y_2(t)$.

8.3.2 Multivariate Granger Causality

The idea of Granger causality is extended to a general multivariate case. Let us define q -dimensional time series as $\mathbf{y}(t)$ and r -dimensional time series as $\mathbf{x}(t)$. We consider

a case in which $\mathbf{x}(t)$ and $\mathbf{y}(t)$ obey the MVAR process such that

$$\mathbf{x}(t) = \sum_{p=1}^P \mathbf{A}_x(p) \mathbf{x}(t-p) + \mathbf{e}_x(t), \quad (8.25)$$

and

$$\mathbf{y}(t) = \sum_{p=1}^P \mathbf{B}_y(p) \mathbf{y}(t-p) + \mathbf{e}_y(t). \quad (8.26)$$

We define covariance matrices of the residuals in this case, such that

$$\boldsymbol{\Sigma}_e^x = \langle \mathbf{e}_x(t) \mathbf{e}_x^T(t) \rangle, \quad (8.27)$$

$$\boldsymbol{\Sigma}_e^y = \langle \mathbf{e}_y(t) \mathbf{e}_y^T(t) \rangle, \quad (8.28)$$

where $\langle \cdot \rangle$ indicates the ensemble average.

We next assume that $\mathbf{x}(t)$ and $\mathbf{y}(t)$ obey the following MVAR process

$$\mathbf{x}(t) = \sum_{p=1}^P \mathbf{A}_x(p) \mathbf{x}(t-p) + \sum_{p=1}^P \mathbf{B}_y(p) \mathbf{y}(t-p) + \boldsymbol{\epsilon}_x(t), \quad (8.29)$$

and

$$\mathbf{y}(t) = \sum_{p=1}^P \mathbf{A}_y(p) \mathbf{y}(t-p) + \sum_{p=1}^P \mathbf{B}_x(p) \mathbf{x}(t-p) + \boldsymbol{\epsilon}_y(t). \quad (8.30)$$

We can define covariance matrices of the residuals, such that

$$\boldsymbol{\Sigma}_{\epsilon}^x = \langle \boldsymbol{\epsilon}_x(t) \boldsymbol{\epsilon}_x^T(t) \rangle, \quad (8.31)$$

$$\boldsymbol{\Sigma}_{\epsilon}^y = \langle \boldsymbol{\epsilon}_y(t) \boldsymbol{\epsilon}_y^T(t) \rangle. \quad (8.32)$$

Using the same idea for Eqs.(8.23) and (8.24), the multivariate Granger causality $\mathcal{G}_{x \rightarrow y}$ and $\mathcal{G}_{y \rightarrow x}$ are given by

$$\mathcal{G}_{x \rightarrow y} = \log \frac{|\boldsymbol{\Sigma}_e^y|}{|\boldsymbol{\Sigma}_{\epsilon}^y|}, \quad (8.33)$$

$$\mathcal{G}_{y \rightarrow x} = \log \frac{|\boldsymbol{\Sigma}_e^x|}{|\boldsymbol{\Sigma}_{\epsilon}^x|}, \quad (8.34)$$

where $|\cdot|$ indicates the matrix determinant.

8.3.3 Total Interdependence

Let us define the $(q+r)$ -dimensional augmented time series $\mathbf{z}(t) = [\mathbf{x}^T(t), \mathbf{y}^T(t)]^T$, which obeys

$$\mathbf{z}(t) = \sum_{p=1}^P \mathbf{A}_z(p) \mathbf{z}(t-p) + \boldsymbol{\epsilon}_z(t). \quad (8.35)$$

The covariance matrix of the residual vector $\boldsymbol{\epsilon}_z(t)$ is defined as $\boldsymbol{\Sigma}_{\epsilon}^z$. When the time series $\mathbf{x}(t)$ and $\mathbf{y}(t)$ are independent, $\boldsymbol{\epsilon}_z(t)$ is expressed as $\boldsymbol{\epsilon}_z(t) = [\boldsymbol{\epsilon}_x(t), \boldsymbol{\epsilon}_y(t)]^T$, and $\boldsymbol{\Sigma}_{\epsilon}^z$ is expressed as

$$\boldsymbol{\Sigma}_{\epsilon}^z = \langle \boldsymbol{\epsilon}_z(t) \boldsymbol{\epsilon}_z^T(t) \rangle = \begin{bmatrix} \boldsymbol{\Sigma}_{\epsilon}^x & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{\epsilon}^y \end{bmatrix}. \quad (8.36)$$

When $\mathbf{x}(t)$ and $\mathbf{y}(t)$ are not independent, it is expressed as

$$\boldsymbol{\Sigma}_{\epsilon}^z = \begin{bmatrix} \boldsymbol{\Sigma}_{\epsilon}^x & \boldsymbol{\Sigma}_{\epsilon}^{xy} \\ \boldsymbol{\Sigma}_{\epsilon}^{yx} & \boldsymbol{\Sigma}_{\epsilon}^y \end{bmatrix}, \quad (8.37)$$

where

$$\boldsymbol{\Sigma}_{\epsilon}^{xy} = \langle \boldsymbol{\epsilon}_x(t) \boldsymbol{\epsilon}_y^T(t) \rangle, \quad (8.38)$$

$$\boldsymbol{\Sigma}_{\epsilon}^{yx} = \langle \boldsymbol{\epsilon}_y(t) \boldsymbol{\epsilon}_x^T(t) \rangle = (\boldsymbol{\Sigma}_{\epsilon}^{xy})^T. \quad (8.39)$$

The total interdependence between \mathbf{x} and \mathbf{y} , $\mathcal{I}_{\{\mathbf{x}, \mathbf{y}\}}$, is defined such that

$$\mathcal{I}_{\{\mathbf{x}, \mathbf{y}\}} = \log \frac{\left| \begin{bmatrix} \boldsymbol{\Sigma}_{\epsilon}^x & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{\epsilon}^y \end{bmatrix} \right|}{\left| \boldsymbol{\Sigma}_{\epsilon}^z \right|} = \log \frac{\left| \boldsymbol{\Sigma}_{\epsilon}^x \right| \left| \boldsymbol{\Sigma}_{\epsilon}^y \right|}{\left| \boldsymbol{\Sigma}_{\epsilon}^z \right|}. \quad (8.40)$$

This total interdependence expresses the deviation of the value of $\left| \boldsymbol{\Sigma}_{\epsilon}^z \right|$ from its value obtained when $\mathbf{x}(t)$ and $\mathbf{y}(t)$ are independent. Using Eqs. (8.33) and (8.34), it is easy to show the relationship

$$\mathcal{I}_{\{\mathbf{x}, \mathbf{y}\}} - \mathcal{G}_{\mathbf{x} \rightarrow \mathbf{y}} - \mathcal{G}_{\mathbf{y} \rightarrow \mathbf{x}} = \log \frac{\left| \boldsymbol{\Sigma}_{\epsilon}^x \right| \left| \boldsymbol{\Sigma}_{\epsilon}^y \right|}{\left| \boldsymbol{\Sigma}_{\epsilon}^z \right|}. \quad (8.41)$$

According to Geweke [3], the right-hand side of the equation above is defined as $\mathcal{I}_{\{\mathbf{x}, \mathbf{y}\}}$:

$$\mathcal{I}_{\{\mathbf{x}, \mathbf{y}\}} = \log \frac{\left| \boldsymbol{\Sigma}_{\epsilon}^x \right| \left| \boldsymbol{\Sigma}_{\epsilon}^y \right|}{\left| \boldsymbol{\Sigma}_{\epsilon}^z \right|}, \quad (8.42)$$

and this $\mathcal{I}_{\{x,y\}}$ is called the instantaneous dependence. The following relationship holds among the causality measures mentioned above,

$$\mathcal{I}_{\{x,y\}} = \mathcal{G}_{x \rightarrow y} + \mathcal{G}_{y \rightarrow x} + \mathcal{I}_{\{x \cdot y\}}.$$

That is, the total interdependence between the two time series $x(t)$ and $y(t)$ is expressed as the summation of the Granger causality from $x(t)$ to $y(t)$, the Granger causality from $y(t)$ to $x(t)$, and the instantaneous dependence.

The rationale of this instantaneous dependence can be explained in the following manner. Using the determinant identity in Eq.(C.94), Eq.(8.42) is rewritten as

$$\mathcal{I}_{\{x \cdot y\}} = \log |\boldsymbol{\Sigma}_\epsilon^x| - \log |\Delta|, \quad (8.43)$$

where

$$\Delta = \boldsymbol{\Sigma}_\epsilon^x - \boldsymbol{\Sigma}_\epsilon^{xy} (\boldsymbol{\Sigma}_\epsilon^y)^{-1} \boldsymbol{\Sigma}_\epsilon^{yx}. \quad (8.44)$$

Let us consider the linear regression in which the residual signal $\epsilon_x(t)$ is regressed by the other residual signal $\epsilon_y(t)$. In Eq.(8.44), Δ expresses the covariance matrix of the residual of this linear regression, according to the arguments in Sect.C.3.3.

Let us suppose a case where $x(t)$ and $y(t)$ contain a common instantaneous interaction $\nu(t)$, which does not exist in the past values of $x(t)$ and $y(t)$. In such a case, both the residual signals $\epsilon_x(t)$ and $\epsilon_y(t)$ contain the common component $\nu(t)$, and because of this, $\nu(t)$ is regressed out when $\epsilon_x(t)$ is regressed by $\epsilon_y(t)$. Therefore, $\log |\Delta|$ does not contain the influence of this common instantaneous component, while $\log |\boldsymbol{\Sigma}_\epsilon^x|$ does. Accordingly, $\mathcal{I}_{\{x \cdot y\}}$ represents the influence of the instantaneous component alone.

8.4 Spectral Granger Causality: Geweke Measures

8.4.1 Basic Relationships in the Frequency Domain

Since brain activities have spectral dependence, we naturally like to perform causality analysis in the spectral domain. The extension of Granger causality analysis into the spectral domain has been investigated by Geweke [3, 4]. The arguments in this section are according to Ding [11], and we restrict our arguments to a bivariate process²:

$$\begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \sum_{p=1}^{\infty} A(p) \begin{bmatrix} x(t-p) \\ y(t-p) \end{bmatrix} + \begin{bmatrix} e_x(t) \\ e_y(t) \end{bmatrix}, \quad (8.45)$$

² The extension of the arguments here to a general multivariate case remains unknown.

where $\mathbf{A}(p)$ is a (2×2) coefficient matrix. We can derive Eq. (8.7), which is explicitly written as

$$\begin{bmatrix} x(f) \\ y(f) \end{bmatrix} = \begin{bmatrix} H_{xx} & H_{xy} \\ H_{yx} & H_{yy} \end{bmatrix} \begin{bmatrix} e_x(f) \\ e_y(f) \end{bmatrix}. \quad (8.46)$$

In the expressions above, the explicit notation of the frequency f is omitted from the components of $\mathbf{H}(f)$. Denoting the residual covariance matrix such that

$$\begin{bmatrix} \langle |e_x(f)|^2 \rangle & \langle e_x(f)e_y(f)^* \rangle \\ \langle e_y(f)e_x(f)^* \rangle & \langle |e_y(f)|^2 \rangle \end{bmatrix} = \begin{bmatrix} \Sigma_x & \gamma \\ \gamma^* & \Sigma_y \end{bmatrix}, \quad (8.47)$$

where the superscript * indicates the complex conjugation. Equation (8.8) is explicitly written as

$$\begin{bmatrix} \langle |x(f)|^2 \rangle & \langle x(f)y(f)^* \rangle \\ \langle y(f)x(f)^* \rangle & \langle |y(f)|^2 \rangle \end{bmatrix} = \mathbf{H} \begin{bmatrix} \Sigma_x & \gamma \\ \gamma^* & \Sigma_y \end{bmatrix} \mathbf{H}^H. \quad (8.48)$$

8.4.2 Total Interdependence and Coherence

We first define the total interdependence between x and y , which corresponds to $\mathcal{I}_{\{x,y\}}$ in Eq. (8.40). Using the arguments in Sect. 8.3.3, the total interdependence between x and y in the spectral domain is defined as

$$\begin{aligned} f_{\{x,y\}} &= \log \frac{\begin{vmatrix} \langle |x(f)|^2 \rangle & 0 \\ 0 & \langle |y(f)|^2 \rangle \end{vmatrix}}{|S(f)|} \\ &= \log \frac{\langle |x(f)|^2 \rangle \langle |y(f)|^2 \rangle}{\langle |x(f)|^2 \rangle \langle |y(f)|^2 \rangle - |\langle x(f)y(f)^* \rangle|^2}. \end{aligned} \quad (8.49)$$

Comparing the equation above to Eq. (7.40) in Chap. 7, it can immediately be seen that under the Gaussianity assumption, the spectral domain total interdependence is equal to the mutual information discussed in Sect. 7.5.2. Therefore, using the same arguments, the total interdependence $f_{\{x,y\}}$ is related to the magnitude coherence $|\phi(f)|$ through

$$f_{\{x,y\}}(f) = -\log(1 - |\phi(f)|^2), \quad (8.50)$$

where

$$|\phi(f)|^2 = \frac{|\langle x(f)y(f)^* \rangle|^2}{\langle |x(f)|^2 \rangle \langle |y(f)|^2 \rangle}. \quad (8.51)$$

8.4.3 Deriving Causal Relationships in the Frequency Domain

In order to derive causal relationships between two channels, we should distinguish between the amount of the signal power coming from its own past and the amount of the power coming from the past of the other channel. On the basis of Eq. (8.46), we are able to do this separation. That is, using Eq. (8.46), we have

$$x(f) = H_{xx}e_x(f) + H_{xy}e_y(f). \quad (8.52)$$

On the right-hand side of the equation above, the first term is the intrinsic term representing the influence of the past of $x(t)$, and the second term represents the influence of the past of $y(t)$. We then decompose the power of the signal $x(f)$, such that

$$\langle |x(f)|^2 \rangle = H_{xx}\Sigma_x H_{xx}^* + H_{xy}\Sigma_y H_{xy}^* + 2\Re \left[\gamma H_{xx}H_{xy}^* \right]. \quad (8.53)$$

On the right-hand side above, the first term represents the intrinsic term and the second term expresses the influence from the past of $y(f)$, i.e., the causal term. However, it is unclear how much amount of the information represented by the third term is attributed to the intrinsic or the causal terms. Therefore, to attain the intrinsic/causal factorization, we first transform Eq. (8.46) into a domain where the innovation terms are uncorrelated, and, in this domain, we factorize the total signal power into the intrinsic and causal terms.

Let us factorize $\langle |x(f)|^2 \rangle$ in this manner. The transformation is performed in this case using $\boldsymbol{\Pi}$ defined such that

$$\boldsymbol{\Pi} = \begin{bmatrix} 1 & 0 \\ -\gamma^*/\Sigma_x & 1 \end{bmatrix}. \quad (8.54)$$

Using Eq. (8.46) and

$$\boldsymbol{\Pi}^{-1} = \begin{bmatrix} 1 & 0 \\ \gamma^*/\Sigma_x & 1 \end{bmatrix},$$

we have

$$\begin{bmatrix} x(f) \\ y(f) \end{bmatrix} = \begin{bmatrix} H_{xx} & H_{xy} \\ H_{yx} & H_{yy} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \gamma^*/\Sigma_x & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -\gamma^*/\Sigma_x & 1 \end{bmatrix} \begin{bmatrix} e_x(f) \\ e_y(f) \end{bmatrix}. \quad (8.55)$$

We then have

$$\begin{bmatrix} x(f) \\ y(f) \end{bmatrix} = \begin{bmatrix} H_{xx} + \frac{\gamma^*}{\Sigma_x} H_{xy} & H_{xy} \\ H_{yx} + \frac{\gamma^*}{\Sigma_x} H_{yy} & H_{yy} \end{bmatrix} \begin{bmatrix} e_x(f) \\ -\frac{\gamma^*}{\Sigma_x} e_x(f) + e_y(f) \end{bmatrix}. \quad (8.56)$$

The cross correlation of the two innovation components, $e_x(f)$ and $-\frac{\gamma^*}{\Sigma_x}e_x(f) + e_y(f)$ is zero, because

$$\begin{aligned} & \left\langle \left[-\frac{\gamma^*}{\Sigma_x}e_x(f) + e_y(f) \right] e_x(f)^* \right\rangle \\ &= -\frac{\gamma^*}{\Sigma_x} \langle |e_x(f)|^2 \rangle + \langle e_y(f)e_x(f)^* \rangle = -\frac{\gamma^*}{\Sigma_x} \Sigma_x + \gamma^* = 0 \end{aligned} \quad (8.57)$$

We can thus derive

$$\begin{aligned} & \begin{bmatrix} \langle |x(f)|^2 \rangle & \langle x(f)y(f)^* \rangle \\ \langle y(f)x(f)^* \rangle & \langle |y(f)|^2 \rangle \end{bmatrix} \\ &= \begin{bmatrix} H_{xx} + \frac{\gamma^*}{\Sigma_x}H_{xy} & H_{xy} \\ H_{yx} + \frac{\gamma^*}{\Sigma_x}H_{yy} & H_{yy} \end{bmatrix} \begin{bmatrix} \Sigma_x & 0 \\ 0 & \Sigma_y - \frac{|\gamma|^2}{\Sigma_x} \end{bmatrix} \begin{bmatrix} H_{xx}^* + \frac{\gamma}{\Sigma_x}H_{xy}^* & H_{yx}^* + \frac{\gamma}{\Sigma_x}H_{yy}^* \\ H_{xy}^* & H_{yy}^* \end{bmatrix}. \end{aligned} \quad (8.58)$$

Therefore, the signal power $\langle |x(f)|^2 \rangle$ is factorized to

$$\langle |x(f)|^2 \rangle = \left(H_{xx} + \frac{\gamma^*}{\Sigma_x}H_{xy} \right) \Sigma_x \left(H_{xx}^* + \frac{\gamma}{\Sigma_x}H_{xy}^* \right) + H_{xy} \left(\Sigma_y - \frac{|\gamma|^2}{\Sigma_x} \right) H_{xy}^*. \quad (8.59)$$

On the right-hand side of this equation, the first term is interpreted as the intrinsic term, which represents the influence of the past of $x(t)$ on the power spectrum $\langle |x(f)|^2 \rangle$. The second term is interpreted as the causal influence term, which represents the influence of the past of $y(t)$ on $\langle |x(f)|^2 \rangle$. Thus, the spectral Granger causality, $f_{y \rightarrow x}(f)$, is defined as

$$\begin{aligned} f_{y \rightarrow x}(f) &= \log \frac{\langle |x(f)|^2 \rangle}{(H_{xx} + \frac{\gamma^*}{\Sigma_x}H_{xy}) \Sigma_x (H_{xx}^* + \frac{\gamma}{\Sigma_x}H_{xy}^*)} \\ &= -\log \frac{\langle |x(f)|^2 \rangle - H_{xy}(\Sigma_y - \frac{|\gamma|^2}{\Sigma_x})H_{xy}^*}{\langle |x(f)|^2 \rangle} \\ &= -\log \left[1 - \frac{(\Sigma_y - \frac{|\gamma|^2}{\Sigma_x})|H_{xy}(f)|^2}{\langle |x(f)|^2 \rangle} \right]. \end{aligned} \quad (8.60)$$

The spectral Granger causality, $f_{x \rightarrow y}(f)$, can be derived in exactly the same manner, using

$$\boldsymbol{\Pi} = \begin{bmatrix} 1 & -\gamma^*/\Sigma_y \\ 0 & 1 \end{bmatrix},$$

and the results are given by

$$f_{x \rightarrow y}(f) = -\log \left[1 - \frac{(\Sigma_x - \frac{|\gamma|^2}{\Sigma_y})|H_{yx}(f)|^2}{\langle |y(f)|^2 \rangle} \right]. \quad (8.61)$$

8.5 Other MVAR-Modeling-Based Measures

8.5.1 Directed Transfer Function (DTF)

The directed transfer function (DTF) [5] is derived from a definition of causality that is somewhat different from the Granger causality. For a bivariate AR process, we have the relationships

$$\begin{aligned} y_1(t) &= \sum_{p=1}^P A_{1,1}(p)y_1(t-p) + \sum_{p=1}^P A_{1,2}(p)y_2(t-p) + \epsilon_1(t) \\ y_2(t) &= \sum_{p=1}^P A_{2,1}(p)y_1(t-p) + \sum_{p=1}^P A_{2,2}(p)y_2(t-p) + \epsilon_2(t). \end{aligned}$$

It can be intuitively clear that, for example, the influence of the past values of $y_2(t)$ on the current value of $y_1(t)$ can be assessed by the values of the AR coefficients $A_{1,2}(p)$, (where $p = 1, \dots, P$). Namely, the causal influence of $y_2(t)$ on $y_1(t)$ can be assessed by using

$$\sqrt{\sum_{p=1}^P A_{1,2}^2(p)}. \quad (8.62)$$

The spectral domain quantities that play a role similar to the above quantity are

$$|\tilde{A}_{1,2}(f)| \text{ and } |H_{1,2}(f)|, \quad (8.63)$$

where the matrix $\tilde{A}(f)$ is defined in Eq.(8.5) and the transfer matrix $H(f)$ is defined as $H(f) = \tilde{A}(f)^{-1}$. The directed transfer function (DTF) makes use of $H(f)$ to express the causal relationship.

We define the (unnormalized) directed transfer function (DTF) using the elements of the transfer function $H(f)$. Namely, the unnormalized directed transfer function that represents the causal influence of $y_2(t)$ on $y_1(t)$ is defined as $H_{1,2}(f)$. For a bivariate case, The relationship,

$$H_{1,2}(f) = \frac{\bar{A}_{1,2}(f)}{|\bar{A}(f)|}$$

indicates that the directed transfer function is equal to $\bar{A}_{1,2}(f)$ aside from the scaling constant $|\bar{A}(f)|$.

One advantage of this definition is that the extension to a general multivariate case is straightforward. In a general multivariate case, the unnormalized transfer function is equal to

$$|H_{1,2}(f)| = \frac{|\mathbf{M}_{2,1}(f)|}{|\bar{A}(f)|}, \quad (8.64)$$

where $\mathbf{M}_{2,1}(f)$ is a minor of the $(2, 1)$ element of the matrix $\bar{A}(f)$. In a trivariate case, using Eq. (8.64), the unnormalized transfer function is derived as

$$|H_{1,2}(f)| = \frac{|\bar{A}_{1,2}\bar{A}_{3,3} - \bar{A}_{3,1}\bar{A}_{2,3}|}{|\bar{A}(f)|}, \quad (8.65)$$

where the explicit notation (f) is omitted in the numerator. On the right hand side, $\bar{A}_{1,2}$ represents the direct influence from channel #1 to channel #2. The equation above shows that even when $\bar{A}_{1,2}$ is equal to zero, $|H_{1,2}|$ cannot equal zero because the term $\bar{A}_{3,1}\bar{A}_{2,3}$ can also not be zero. This term represents the indirect influence of channel #1 on channel #2 via channel #3. That is, DTF contains the indirect causal influence, as well as the direct causal influence.

The normalized DTF can be defined by normalizing the causal influence on the j th channel from all other channels. That is, when total q channels exist, the normalized DTF is defined as

$$\mu_{k \rightarrow j} = \frac{H_{j,k}(f)}{\sqrt{[\mathbf{h}_j^H \mathbf{h}_j]}} = \frac{H_{j,k}(f)}{\sqrt{\sum_{m=1}^q |H_{j,m}(f)|^2}}. \quad (8.66)$$

8.5.2 Relationship Between DTF and Coherence

The coherence can be expressed as a result of factorization of the normalized DTF. To show this, we define the directed coherence $\gamma_{j,k}$ such that

$$\gamma_{j,k} = \frac{H_{j,k}(f)}{\sqrt{\mathbf{h}_j^H \boldsymbol{\Sigma} \mathbf{h}_j}}. \quad (8.67)$$

We then define the column vector

$$\boldsymbol{\gamma}_j(f) = [\gamma_{j,1}^*, \dots, \gamma_{j,q}^*]^T, \quad (8.68)$$

and using this column vector, coherence is expressed as³

$$\phi_{j,k}(f) = \frac{S_{j,k}}{\sqrt{S_{j,j} S_{k,k}}} = \frac{\mathbf{h}_j^H \boldsymbol{\Sigma} \mathbf{h}_k}{\sqrt{[\mathbf{h}_j^H \boldsymbol{\Sigma} \mathbf{h}_j][\mathbf{h}_k^H \boldsymbol{\Sigma} \mathbf{h}_k]}} = \gamma_j^H(f) \boldsymbol{\Sigma} \gamma_k(f). \quad (8.69)$$

This directed coherence $\gamma_{j,k}$ is known to represent the directional influence from the k th to the j th channels. It contains the residual covariance matrix $\boldsymbol{\Sigma}$ whose off-diagonal terms may represent the instantaneous interaction between two channels, according to the arguments in Sect. 8.3.3. Setting $\boldsymbol{\Sigma}$ equal to \mathbf{I} , we have

$$\gamma_{j,k} = \frac{H_{j,k}(f)}{\sqrt{|\mathbf{h}_j^H \mathbf{h}_j|}}, \quad (8.70)$$

which is exactly equal to the normalized DTF in Eq.(8.66). Namely, the DTF is a measure for directional influences obtained by removing instantaneous components from the directed coherence. Under the assumption that $\boldsymbol{\Sigma} = \mathbf{I}$, coherence $\phi_{j,k}(f)$ can be decomposed of the sum of the product of DTFs, such that

$$\phi_{j,k}(f) = \frac{\sum_{m=1}^q H_{j,m}(f) H_{k,m}^*(f)}{\sqrt{[\mathbf{h}_j^H \mathbf{h}_j][\mathbf{h}_k^H \mathbf{h}_k]}} = \sum_{m=1}^q \mu_{m \rightarrow j} \mu_{m \rightarrow k}^*. \quad (8.71)$$

8.5.3 Partial Directed Coherence (PDC)

In Sect. 8.5.2, DTF is expressed using a factorization of coherence. We apply similar factorization to partial coherence to obtain the partial directed coherence (PDC) [6, 12]. The starting point is Eq. (8.14), which is re-written as,

$$\kappa_{j,k}(f) = \frac{\bar{\mathbf{a}}_j^H \boldsymbol{\Sigma}^{-1} \bar{\mathbf{a}}_k}{\sqrt{[\bar{\mathbf{a}}_j^H \boldsymbol{\Sigma}^{-1} \bar{\mathbf{a}}_j][\bar{\mathbf{a}}_k^H \boldsymbol{\Sigma}^{-1} \bar{\mathbf{a}}_k]}}. \quad (8.72)$$

Here, recall that $\bar{\mathbf{a}}_j$ is the j th column of the matrix $\bar{\mathbf{A}}(f)$ and $\bar{\mathbf{a}}_j = [\bar{A}_{1,j}, \dots, \bar{A}_{q,j}]$. We define

$$\bar{\pi}_{j,k}(f) = \frac{\bar{A}_{j,k}(f)}{\sqrt{\bar{\mathbf{a}}_k^H \boldsymbol{\Sigma}^{-1} \bar{\mathbf{a}}_k}}, \quad (8.73)$$

³ Note that since \mathbf{h}_j is the j th column vector of \mathbf{H}^H , \mathbf{h}_j is equal to $\mathbf{h}_j = [H_{j,1}^*, \dots, H_{j,q}^*]^T$.

and a column vector $\bar{\pi}_k(f)$ as

$$\bar{\pi}_k(f) = [\bar{\pi}_{1,k}, \dots, \bar{\pi}_{q,k}]^T. \quad (8.74)$$

We can then express the partial coherence as

$$\kappa_{j,k}(f) = \bar{\pi}_j^H(f) \boldsymbol{\Sigma}^{-1} \bar{\pi}_k(f). \quad (8.75)$$

The equation above indicates that the partial coherence $\kappa_{j,k}(f)$ is factorized, and the form of the factorization is very similar to the factorization of the coherence in Eq. (8.69).

The partial directed coherence (PDC) is defined using $\bar{\pi}_{j,k}$ in Eq. (8.73) with replacing $\boldsymbol{\Sigma}$ with \mathbf{I} for removing the instantaneous interaction. That is, the PDC from the k th to the j th channels, $\pi_{k \rightarrow j}$, is defined as

$$\pi_{k \rightarrow j} = \frac{\bar{A}_{j,k}(f)}{\sqrt{\bar{a}_k^H \bar{a}_k}}. \quad (8.76)$$

Assuming $\boldsymbol{\Sigma}^{-1} = \mathbf{I}$, the partial coherence $\kappa_{j,k}(f)$ is factorized using PDC, such that

$$\kappa_{j,k}(f) = \frac{\sum_{m=1}^q \bar{A}_{m,j}^*(f) \bar{A}_{m,k}(f)}{\sqrt{[\bar{a}_j^H \bar{a}_j][\bar{a}_k^H \bar{a}_k]}} = \sum_{m=1}^q \pi_{j \rightarrow m}^* \pi_{k \rightarrow m}. \quad (8.77)$$

As discussed in the previous sections, the non-diagonal elements of the MVAR coefficient matrix $\bar{A}(f)$ should contain information on causal interaction. The PDC directly makes use of this information. The DTF also uses this information but it does so in an indirect manner through the transfer matrix $\mathbf{H}(f)$, which is the inverse of $\bar{A}(f)$. Although PDC is derived in a somewhat heuristic manner, it represents causality through $\bar{A}_{j,k}(f)$. The major difference between PDC and DTF is that PDC only represents the direct causal influence and is not affected by the indirect influence.

8.6 Transfer Entropy

8.6.1 Definition

Granger causality and its related measures such as DTF and PDC relies on the MVAR modeling of voxel time series. In that sense, these measures are model-based. In this section, we introduce transfer entropy [7, 13], which does not use the MVAR modeling of voxel time series. Explanations on fundamentals of entropy and mutual information are found in Sect. C.3.2 in the Appendix.

We denote two random vector time series $\mathbf{x}(t)$ and $\mathbf{y}(t)$. Let us define vectors $\tilde{\mathbf{x}}(t)$ and $\tilde{\mathbf{y}}(t)$ as those made by concatenating their past values such that

$$\tilde{\mathbf{x}}(t) = \begin{bmatrix} \mathbf{x}(t-1) \\ \mathbf{x}(t-2) \\ \vdots \\ \mathbf{x}(t-P) \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{y}}(t) = \begin{bmatrix} \mathbf{y}(t-1) \\ \mathbf{y}(t-2) \\ \vdots \\ \mathbf{y}(t-P) \end{bmatrix}. \quad (8.78)$$

We then define the conditional entropy of $\mathbf{y}(t)$, given its past values $\tilde{\mathbf{y}}(t)$, such that

$$\mathcal{H}(\mathbf{y}|\tilde{\mathbf{y}}) = - \iint p(\mathbf{y}, \tilde{\mathbf{y}}) \log p(\mathbf{y}|\tilde{\mathbf{y}}) d\mathbf{y} d\tilde{\mathbf{y}}. \quad (8.79)$$

Similarly, we define the conditional entropy of $\mathbf{y}(t)$, given the past values $\tilde{\mathbf{x}}(t)$ and $\tilde{\mathbf{y}}(t)$, such that

$$\mathcal{H}(\mathbf{y}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = - \iint p(\tilde{\mathbf{x}}, \mathbf{y}, \tilde{\mathbf{y}}) \log p(\mathbf{y}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) d\mathbf{y} d\tilde{\mathbf{x}} d\tilde{\mathbf{y}}. \quad (8.80)$$

The transfer entropy $\mathcal{H}_{\mathbf{x} \rightarrow \mathbf{y}}$ is defined as

$$\begin{aligned} \mathcal{H}_{\mathbf{x} \rightarrow \mathbf{y}} &= \mathcal{H}(\mathbf{y}|\tilde{\mathbf{y}}) - \mathcal{H}(\mathbf{y}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \\ &= \iint p(\tilde{\mathbf{x}}, \mathbf{y}, \tilde{\mathbf{y}}) \log \frac{\log p(\mathbf{y}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})}{\log p(\mathbf{y}|\tilde{\mathbf{y}})} d\mathbf{y} d\tilde{\mathbf{x}} d\tilde{\mathbf{y}}. \end{aligned} \quad (8.81)$$

In the equations above, the explicit time notation (t) is omitted for simplicity. In Eq.(8.81), $\mathcal{H}(\mathbf{y}|\tilde{\mathbf{y}})$ represents the uncertainty on the current value of \mathbf{y} , when we know $\tilde{\mathbf{y}}$, which is the past values of \mathbf{y} . Also, $\mathcal{H}(\mathbf{y}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ represents the uncertainty on the current value of \mathbf{y} , when we know both $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{y}}$, which are the past values of \mathbf{x} and \mathbf{y} . Therefore, the transfer entropy is equal to the reduction of uncertainty of the current value of \mathbf{y} as a result of knowing the past values of \mathbf{x} .

8.6.2 Transfer Entropy Under Gaussianity Assumption

Assuming that the random vectors \mathbf{x} and \mathbf{y} follow a Gaussian distribution, using Eq.(C.61), the conditional entropy $\mathcal{H}(\mathbf{y}|\tilde{\mathbf{y}})$ is expressed as [14]

$$\mathcal{H}(\mathbf{y}|\tilde{\mathbf{y}}) = \frac{1}{2} \log \left| \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{y\tilde{y}} \boldsymbol{\Sigma}_{\tilde{y}\tilde{y}}^{-1} \boldsymbol{\Sigma}_{y\tilde{y}}^T \right|, \quad (8.82)$$

where⁴

$$\boldsymbol{\Sigma}_{yy} = \langle \mathbf{y} \mathbf{y}^T \rangle, \quad (8.83)$$

$$\boldsymbol{\Sigma}_{y\tilde{y}} = \langle \mathbf{y} \tilde{\mathbf{y}}^T \rangle, \quad (8.84)$$

$$\boldsymbol{\Sigma}_{\tilde{y}\tilde{y}} = \langle \tilde{\mathbf{y}} \tilde{\mathbf{y}}^T \rangle. \quad (8.85)$$

Similarly, the conditional entropy $\mathcal{H}(\mathbf{y}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ is expressed as

$$\mathcal{H}(\mathbf{y}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = \mathcal{H}(\mathbf{y}|\tilde{\mathbf{z}}) = \frac{1}{2} \log \left| \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{y\tilde{z}} \boldsymbol{\Sigma}_{\tilde{z}\tilde{z}}^{-1} \boldsymbol{\Sigma}_{y\tilde{z}}^T \right|, \quad (8.86)$$

where $\tilde{\mathbf{z}} = [\tilde{\mathbf{y}}^T, \tilde{\mathbf{x}}^T]^T$ and

$$\boldsymbol{\Sigma}_{y\tilde{z}} = \langle \mathbf{y} \tilde{\mathbf{z}}^T \rangle, \quad (8.87)$$

$$\boldsymbol{\Sigma}_{\tilde{z}\tilde{z}} = \langle \tilde{\mathbf{z}} \tilde{\mathbf{z}}^T \rangle. \quad (8.88)$$

Thus, the transfer entropy $\mathcal{H}_{x \rightarrow y}$ is given by

$$\begin{aligned} \mathcal{H}_{x \rightarrow y} &= \mathcal{H}(\mathbf{y}|\tilde{\mathbf{y}}) - \mathcal{H}(\mathbf{y}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \\ &= \mathcal{H}(\mathbf{y}|\tilde{\mathbf{y}}) - \mathcal{H}(\mathbf{y}|\tilde{\mathbf{z}}) = \frac{1}{2} \log \frac{\left| \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{y\tilde{y}} \boldsymbol{\Sigma}_{\tilde{y}\tilde{y}}^{-1} \boldsymbol{\Sigma}_{y\tilde{y}}^T \right|}{\left| \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{y\tilde{z}} \boldsymbol{\Sigma}_{\tilde{z}\tilde{z}}^{-1} \boldsymbol{\Sigma}_{y\tilde{z}}^T \right|}. \end{aligned} \quad (8.89)$$

8.6.3 Equivalence Between Transfer Entropy and Granger Causality

We will show that the transfer entropy and Granger causality are equivalent under the Gaussianity assumption. The arguments here follow those in [14]. As discussed in Sect. 8.3.2, we consider the two forms of the regression to define Granger causality. In the first regression, $\mathbf{y}(t)$ is regressed using only its past values, such that

$$\mathbf{y}(t) = \sum_{p=1}^P \mathbf{A}(p) \mathbf{y}(t-p) + \mathbf{e} = \mathbf{A}\tilde{\mathbf{y}}(t) + \mathbf{e}, \quad (8.90)$$

where $\mathbf{A} = [\mathbf{A}(1), \dots, \mathbf{A}(P)]$, $\tilde{\mathbf{y}}(t)$ is defined in Eq. (8.78), and \mathbf{e} is a residual vector. In the second regression, $\mathbf{y}(t)$ is regressed using not only its past values but also the past values of $\mathbf{x}(t)$, such that

⁴ Note that in Sect. C.3.3 the expectation operator $E[\cdot]$ is used, instead of the averaging operator $\langle \cdot \rangle$. They have the same meaning in the arguments here.

$$\begin{aligned} \mathbf{y}(t) &= \sum_{p=1}^P \mathbf{A}(p)\mathbf{y}(t-p) + \sum_{p=1}^P \mathbf{B}(p)\mathbf{x}(t-p) + \boldsymbol{\epsilon} \\ &= \mathbf{A}\tilde{\mathbf{y}}(t) + \mathbf{B}\tilde{\mathbf{x}}(t) + \boldsymbol{\epsilon} = \mathbf{C}\tilde{\mathbf{z}} + \boldsymbol{\epsilon}. \end{aligned} \quad (8.91)$$

where $\mathbf{B} = [\mathbf{B}(1), \dots, \mathbf{B}(P)]$, $\mathbf{C} = [\mathbf{A}, \mathbf{B}]$, and $\boldsymbol{\epsilon}$ is a residual vector.

The Granger causality from the time series \mathbf{x} to \mathbf{y} , $\mathcal{G}_{x \rightarrow y}$ is defined as

$$\mathcal{G}_{x \rightarrow y} = \log \frac{|\boldsymbol{\Sigma}_e|}{|\boldsymbol{\Sigma}_\epsilon|}, \quad (8.92)$$

where $\boldsymbol{\Sigma}_e$ is the covariance matrix of the residual $\boldsymbol{\epsilon}$ in Eq.(8.90), and $\boldsymbol{\Sigma}_\epsilon$ is the covariance matrix of the residual $\boldsymbol{\epsilon}$ in Eq.(8.91). Using Eq.(C.57), we have

$$|\boldsymbol{\Sigma}_e| = |\boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{y\tilde{y}}\boldsymbol{\Sigma}_{\tilde{y}\tilde{y}}^{-1}\boldsymbol{\Sigma}_{y\tilde{y}}^T|, \quad (8.93)$$

and

$$|\boldsymbol{\Sigma}_\epsilon| = |\boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{y\tilde{z}}\boldsymbol{\Sigma}_{\tilde{z}\tilde{z}}^{-1}\boldsymbol{\Sigma}_{y\tilde{z}}^T|. \quad (8.94)$$

Substituting Eqs. (8.93) and (8.94) into (8.92), we get

$$\mathcal{G}_{x \rightarrow y} = \log \frac{|\boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{y\tilde{y}}\boldsymbol{\Sigma}_{\tilde{y}\tilde{y}}^{-1}\boldsymbol{\Sigma}_{y\tilde{y}}^T|}{|\boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{y\tilde{z}}\boldsymbol{\Sigma}_{\tilde{z}\tilde{z}}^{-1}\boldsymbol{\Sigma}_{y\tilde{z}}^T|}. \quad (8.95)$$

The right-hand side of the equation above is exactly equal to that of Eq.(8.89) except for the multiplicative constant $1/2$, indicating that these two measures are equivalent.

8.6.4 Computation of Transfer Entropy

In Sect.C.3.2 in the Appendix, we show that, assuming the Gaussian processes for the $r \times 1$ real random vector \mathbf{x} and the $q \times 1$ real random vector \mathbf{y} , the mutual information between these two vectors is expressed as

$$\mathcal{I}(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \log \frac{1}{\left| \mathbf{I} - \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx}^T \right|} = \frac{1}{2} \sum_{j=1}^d \log \frac{1}{1 - \lambda_j}, \quad (8.96)$$

where we assume $d = \min\{q, r\}$, and λ_j is the j th eigenvalue of the matrix,

$$\boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx}^T.$$

We can derive a similar formula to compute the transfer entropy. To do so, we use the similarity between the transfer entropy and mutual information. Using Eq. (C.45), the conditional entropy $\mathcal{H}(y|\tilde{x}, \tilde{y})$ is rewritten as

$$\mathcal{H}(y|\tilde{x}, \tilde{y}) = \mathcal{H}(y, \tilde{x}|\tilde{y}) - \mathcal{H}(\tilde{x}|\tilde{y}). \quad (8.97)$$

Substituting this equation into Eq. (8.81), we get

$$\mathcal{H}_{x \rightarrow y} = \mathcal{H}(\tilde{x}|\tilde{y}) + \mathcal{H}(y|\tilde{y}) - \mathcal{H}(y, \tilde{x}|\tilde{y}). \quad (8.98)$$

Comparing the equation above with Eq. (C.46), it can be seen that the transfer entropy is equal to the mutual information between $y(t)$ and $\tilde{x}(t)$, when $\tilde{y}(t)$ is given.

We define $\Sigma_{u,v|w}$ as

$$\Sigma_{u,v|w} = \Sigma_{uv} - \Sigma_{uw} \Sigma_{ww}^{-1} \Sigma_{vw}^T. \quad (8.99)$$

Then, using Eq. (C.61), we can express $\mathcal{H}(y|\tilde{y})$, and $\mathcal{H}(\tilde{x}|\tilde{y})$, such that

$$\mathcal{H}(y|\tilde{y}) = \frac{1}{2} \log |\Sigma_{y,y|\tilde{y}}|, \quad (8.100)$$

$$\mathcal{H}(\tilde{x}|\tilde{y}) = \frac{1}{2} \log |\Sigma_{\tilde{x},\tilde{x}|\tilde{y}}|. \quad (8.101)$$

On the basis of Eq. (C.59), we also derive

$$\mathcal{H}(y, \tilde{x}|\tilde{y}) = \frac{1}{2} \log \left| \begin{array}{c} \Sigma_{y,y|\tilde{y}} \Sigma_{\tilde{x},y|\tilde{y}}^T \\ \Sigma_{\tilde{x},y|\tilde{y}} \Sigma_{\tilde{x},\tilde{x}|\tilde{y}} \end{array} \right|. \quad (8.102)$$

Thus, substituting the equations above into Eq. (8.98), we obtain

$$\mathcal{H}_{x \rightarrow y} = \frac{1}{2} \log \frac{1}{\left| I - \Sigma_{y,y|\tilde{y}}^{-1} \Sigma_{y,\tilde{x}|\tilde{y}} \Sigma_{\tilde{x},\tilde{x}|\tilde{y}}^{-1} \Sigma_{y,\tilde{x}|\tilde{y}}^T \right|}. \quad (8.103)$$

Accordingly, defining the eigenvalues of a matrix,

$$\Sigma_{y,y|\tilde{y}}^{-1} \Sigma_{y,\tilde{x}|\tilde{y}} \Sigma_{\tilde{x},\tilde{x}|\tilde{y}}^{-1} \Sigma_{y,\tilde{x}|\tilde{y}}^T, \quad (8.104)$$

as χ_j ($j = 1, \dots, d$), the transfer entropy is given by

$$\mathcal{H}_{x \rightarrow y} = \frac{1}{2} \sum_{j=1}^d \log \frac{1}{1 - \chi_j}. \quad (8.105)$$

8.7 Estimation of MVAR Coefficients

8.7.1 Least-Squares Algorithm

The Granger-causality and related measures rely on the modeling of the multivariate vector auto-regressive (MVAR) process of the source time series. Use of these measures requires estimating the MVAR coefficient matrices. This section deals with the estimation of MVAR coefficients. A $q \times 1$ random vector $\mathbf{y}(t)$ is modeled using the MVAR process, such that

$$\mathbf{y}(t) = \sum_{p=1}^P A(p)\mathbf{y}(t-p) + \mathbf{e}(t). \quad (8.106)$$

We can estimate the MVAR coefficients $A_{i,j}(p)$ where $i, j = 1, \dots, q$ and $p = 1, \dots, P$ based on the least-squares principle. To derive the least-squares equation, let us explicitly write the MVAR process for the ℓ th component $y_\ell(t)$ as

$$\begin{aligned} y_\ell(t) &= \sum_{j=1}^q A_{\ell,j}(1)y_j(t-1) + \sum_{j=1}^q A_{\ell,j}(2)y_j(t-2) \\ &\quad + \cdots + \sum_{j=1}^q A_{\ell,j}(P)y_j(t-P) + e_\ell(t). \end{aligned} \quad (8.107)$$

We assume that the source time series are obtained at $t = 1, \dots, K$ where $K \gg q \times P$. Then, since Eq. (8.107) holds for $t = (P+1), \dots, K$, a total of $K - P$ linear equations are obtained by setting $t = (P+1), \dots, K$ in Eq. (8.107).

These equations are formulated in a matrix form,

$$\mathbf{y}_\ell = \mathbf{G}\mathbf{x}_\ell + \mathbf{e}_\ell. \quad (8.108)$$

Here, the $(K - P) \times 1$ column vector \mathbf{y}_ℓ is defined as

$$\mathbf{y}_\ell = [y_\ell(P+1), y_\ell(P+2), \dots, y_\ell(K)]^T. \quad (8.109)$$

In Eq. (8.108), \mathbf{G} is a $(K - P) \times Pq$ matrix expressed as

$$\mathbf{G} = \begin{bmatrix} y_1(P) & \cdots & y_q(P) & \cdots & y_1(1) & \cdots & y_q(1) \\ y_1(P+1) & \cdots & y_q(P+1) & \cdots & y_1(2) & \cdots & y_q(2) \\ \vdots & & & & & & \\ y_1(K-1) & \cdots & y_q(K-1) & \cdots & y_1(K-P) & \cdots & y_q(K-P) \end{bmatrix}. \quad (8.110)$$

The column vector \mathbf{x}_ℓ is expressed as

$$\mathbf{x}_\ell = [A_{\ell,1}(1), \dots, A_{\ell,q}(1), \dots, A_{\ell,1}(P), \dots, A_{\ell,q}(P)]^T. \quad (8.111)$$

The residual vector \mathbf{e}_ℓ is given by

$$\mathbf{e}_\ell = [e_\ell(1), \dots, e_\ell(P+1), \dots, e_\ell(K)]^T. \quad (8.112)$$

Equation (8.108) is called the Yule-Walker equation. The least-squares estimate of \mathbf{x}_ℓ , $\hat{\mathbf{x}}_\ell$, is then obtained using,

$$\hat{\mathbf{x}}_\ell = (\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T \mathbf{y}_\ell. \quad (8.113)$$

8.7.2 Sparse Bayesian (Champagne) Algorithm

Since causality analysis is generally performed using the source time series estimated from non-averaged data, the estimated time series inevitably contains a large influence of noise, which may cause errors in the MVAR-coefficient estimation. One approach to reduce such errors is to impose a sparsity constraint when estimating the MVAR coefficients. The key assumption here is that true brain interaction causes a small number of MVAR coefficients to have non-zero values, and most of the MVAR coefficients remain zero. If this is true, the sparsity constraint should be able to prevent MVAR coefficients that must be equal to zero from having erroneous non-zero values.

We can apply a simpler version of the Champagne algorithm described in Chap. 4 to this MVAR coefficient estimation. In the Champagne algorithm, the prior probability distribution of \mathbf{x}_ℓ , $p(\mathbf{x}_\ell)$, is assumed to be Gaussian:

$$p(\mathbf{x}_\ell) = \mathcal{N}(\mathbf{x}_\ell | \mathbf{0}, \boldsymbol{\Phi}^{-1}), \quad (8.114)$$

where $\boldsymbol{\Phi}$ is a diagonal precision matrix. The probability distribution of \mathbf{y}_ℓ given \mathbf{x}_ℓ , $p(\mathbf{y}_\ell | \mathbf{x}_\ell)$, is also assumed to be Gaussian:

$$p(\mathbf{y}_\ell | \mathbf{x}_\ell) = \mathcal{N}(\mathbf{y}_\ell | \mathbf{G}\mathbf{x}_\ell, \boldsymbol{\Lambda}^{-1}), \quad (8.115)$$

where $\boldsymbol{\Lambda}$ is a diagonal noise precision matrix. Then, the posterior distribution of \mathbf{x}_ℓ , $p(\mathbf{x}_\ell | \mathbf{y}_\ell)$, is shown to be Gaussian, and it is expressed as

$$p(\mathbf{x}_\ell | \mathbf{y}_\ell) = \mathcal{N}(\mathbf{x}_\ell | \bar{\mathbf{x}}_\ell, \boldsymbol{\Gamma}^{-1}). \quad (8.116)$$

Unlike the Champagne algorithm in Chap. 4, the update equations for $\boldsymbol{\Phi}$ and $\boldsymbol{\Lambda}$ can be derived using the expectation-maximization (EM) algorithm. This is because,

since $(K - P) \gg Pq$ holds, the estimation problem is not an ill-posed problem. The EM algorithm for the Gaussian model is described in Sect. B.5 in the Appendix. According to Eqs. (B.24) and (B.25), the E-step update equations are given by

$$\boldsymbol{\Gamma} = \mathbf{G}^T \boldsymbol{\Lambda} \mathbf{G} + \boldsymbol{\Phi} \quad (8.117)$$

$$\bar{\mathbf{x}}_\ell = \boldsymbol{\Gamma}^{-1} \mathbf{G}^T \boldsymbol{\Lambda} \mathbf{y}_\ell \quad (8.118)$$

The parameters, $\boldsymbol{\Phi}$ and $\boldsymbol{\Lambda}$, are estimated in the M step of the EM algorithm. The update equation for $\boldsymbol{\Phi}$ is obtained from Eq. (B.42) with setting $K = 1$ in that equation, resulting in

$$\boldsymbol{\Phi}^{-1} = \text{diag} \left[\bar{\mathbf{x}}_\ell \bar{\mathbf{x}}_\ell^T + \boldsymbol{\Gamma}^{-1} \right], \quad (8.119)$$

where $\text{diag}[\cdot]$ indicates a diagonal matrix whose diagonal elements are equal to those of a matrix in the parentheses. According to Eq. (B.45), the noise precision matrix $\boldsymbol{\Lambda}$ is given by

$$\boldsymbol{\Lambda}^{-1} = \text{diag} \left[\|\mathbf{y}_\ell - \mathbf{G}\bar{\mathbf{x}}_\ell\|^2 + \mathbf{G}\boldsymbol{\Gamma}^{-1}\mathbf{G}^T \right]. \quad (8.120)$$

The estimate of the MVAR coefficients is obtained as $\bar{\mathbf{x}}_\ell$, after the EM iteration is terminated. This algorithm is similar to, but considerably simpler than the one proposed in [15].

8.8 Numerical Examples

8.8.1 Experiments Using Bivariate Causal Time Series

Numerical experiments were performed to illustrate the properties of the causality measures described in this chapter. The source-space causality analysis was applied in which source time series are first estimated from simulated MEG recordings, and the MVAR coefficients are then computed using the time series at selected voxels. Here, we use a sensor alignment of the 275 whole-head MEG sensor array from OmegaTM (VMS Medtech, Coquitlam, Canada) neuromagnetometer. Three sources are assumed to exist on the vertical single plane of $x = 0$ cm, as in the numerical experiments in the previous chapters. The source-sensor configuration and the coordinate system are depicted in Fig. 8.1. As shown in this figure, we assume three sources and the time series of these three sources are denoted $s_1(t)$, $s_2(t)$, and $s_3(t)$. The first experiments assume a causal relationship between bivariate time series, and an information flow exists from $s_1(t)$ to $s_2(t)$. The time series $s_1(t)$ and $s_2(t)$ are generated by using the MVAR process [11]

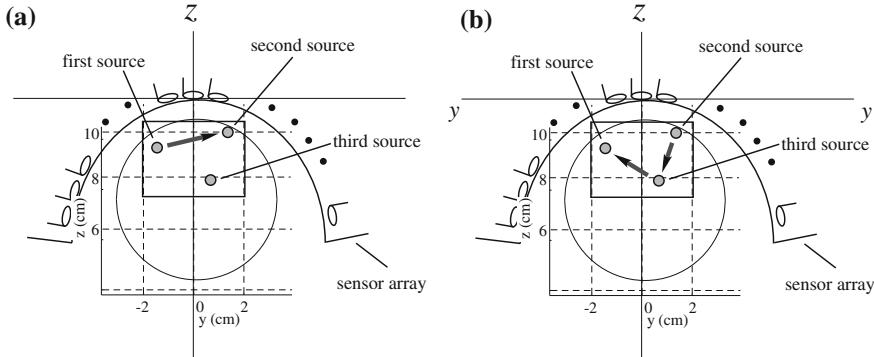


Fig. 8.1 The coordinate system and source-sensor configuration used in the numerical experiments. The plane at $x = 0\text{ cm}$ is shown. The *small circles* show the locations of the three sources, and the *bold arrows* schematically show their causal relationships assumed in the experiments. **a** The first experiment with bivariate causal time series. **b** The second experiments with trivariate causal time series

$$\begin{bmatrix} s_1(t) \\ s_2(t) \end{bmatrix} = \begin{bmatrix} 0.9 & 0 \\ 0.16 & 0.8 \end{bmatrix} \begin{bmatrix} s_1(t-1) \\ s_2(t-1) \end{bmatrix} + \begin{bmatrix} -0.5 & 0 \\ -0.2 & -0.5 \end{bmatrix} \begin{bmatrix} s_1(t-2) \\ s_2(t-2) \end{bmatrix} + e(t). \quad (8.121)$$

The time series of the third source, $s_3(t)$, was generated using Gaussian random numbers. Thus, the third source activity was independent from either the first or the second source activities. The causal relationship assumed in this experiment is depicted in Fig. 8.1a.

Then, simulated sensor recordings were computed by projecting the time series of the three sources onto the sensor space by using the sensor lead field. A small amount of simulated sensor noise was added. The Champagne source reconstruction algorithm was applied to the simulated sensor recordings. Here, three-dimensional reconstruction was performed on a region defined as $-4 \leq x \leq 4$, $-4 \leq y \leq 4$, and $6 \leq z \leq 12\text{ cm}$ with a voxel interval equal to 0.5 cm . Reconstructed source time series $\hat{s}_1(t)$, $\hat{s}_2(t)$, and $\hat{s}_3(t)$ were obtained as the time series at voxels nearest to the assumed source locations.

Once $\hat{s}_1(t)$, $\hat{s}_2(t)$ were obtained, the MVAR coefficients between these time series were estimated by using the least-squares method in Sect. 8.7.1. Using the estimated MVAR coefficients, we computed the spectral Geweke causality described in Sect. 8.4, as well as coherence. The results are shown in Fig. 8.2a.

In these results, coherence (Eq. (8.51)) can detect an interaction between the first and second source activities. The spectral Geweke causality in Eqs. (8.60) and (8.61) detects the unidirectional information flow from the first source activity to the second source activity. We also computed the partial directed coherence (PDC) and the directed transfer function (DTF), with results shown in Fig. 8.2b. The PDC and DTF

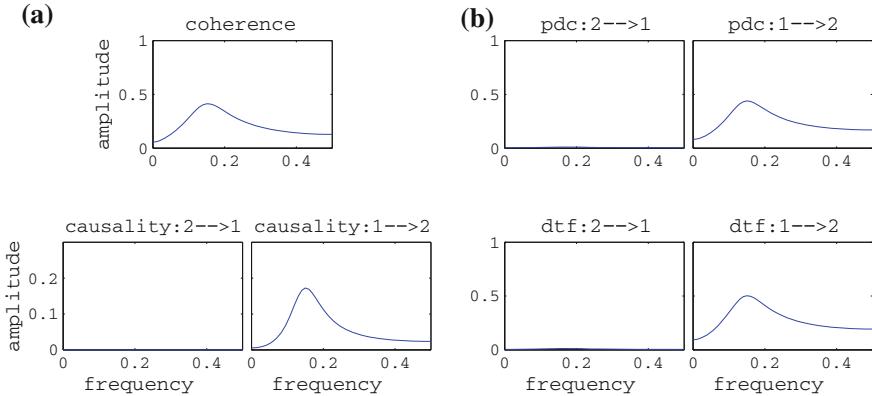


Fig. 8.2 Results of computing bivariate measures between $\hat{s}_1(t)$ and $\hat{s}_2(t)$ in the first experiment. **a** Results of computing spectral causality measures and the coherence. Results for the coherence are shown in the *top panel*. Results for the spectral Geweke causality are shown on the *bottom-left*, and *bottom-right* panels. **b** Results of computing the partial directed coherence (PDC) (*top panels*) and the directed transfer function (DTF) (*bottom panels*). The abscissas of plots indicate the normalized frequency. (The frequency is normalized by the sampling frequency in which the value 0.5 indicates the Nyquist frequency.) The ordinates indicate the relative amplitudes of corresponding measures

can also detect the information flow from the first to the second sources. In this two-channel experiment, the PDC and DTF provide identical results, because, aside from the scaling, both measures are equal to $|\bar{A}_{1,2}(f)|$ in a bivariate case.

The causality measures were computed between $\hat{s}_1(t)$ and $\hat{s}_3(t)$, and these results are shown in Fig. 8.3. Since there was no interaction between $s_1(t)$ and $s_3(t)$, all the causal measures used in the experiment, as well as the coherence, are close to zero in these results.

8.8.2 Experiments Using Trivariate Causal Time Series

Next, we conducted experiments using trivariate causal time series. In these experiments, the time series of the three sources are generated using the MVAR process introduced in [11],

$$\begin{bmatrix} s_1(t) \\ s_2(t) \\ s_3(t) \end{bmatrix} = \begin{bmatrix} 0.8 & 0 & 0.4 \\ 0 & 0.9 & 0 \\ 0 & 0.5 & 0.5 \end{bmatrix} \begin{bmatrix} s_1(t-1) \\ s_2(t-1) \\ s_3(t-1) \end{bmatrix} + \begin{bmatrix} -0.5 & 0 & 0 \\ 0 & -0.8 & 0 \\ 0 & 0 & -0.2 \end{bmatrix} \begin{bmatrix} s_1(t-2) \\ s_2(t-2) \\ s_3(t-2) \end{bmatrix} + \boldsymbol{\epsilon}(t). \quad (8.122)$$

This MVAR process represents the causal relationship depicted in Fig. 8.1b, in which the second source has a directional causal influence on the third source and the third

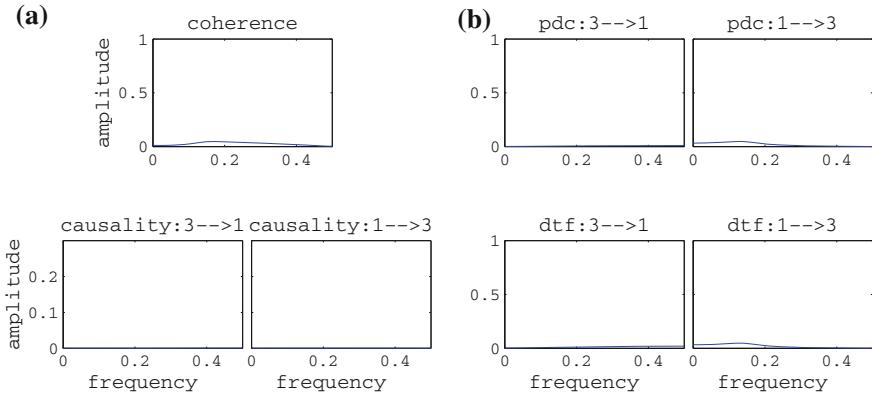


Fig. 8.3 Results of computing bivariate measures between $\hat{s}_1(t)$ and $\hat{s}_3(t)$ in the first experiments. **a** Results of computing spectral causality measures and the coherence. Results for the coherence are shown in the *top panel*. Results for the spectral Geweke causality are shown in the *bottom panels*. **b** Results of computing the partial directed coherence (PDC) (*top panels*) and the directed transfer function (DTF) (*bottom panels*). Explanations on the ordinate and abscissa variables are found in the caption for Fig. 8.2

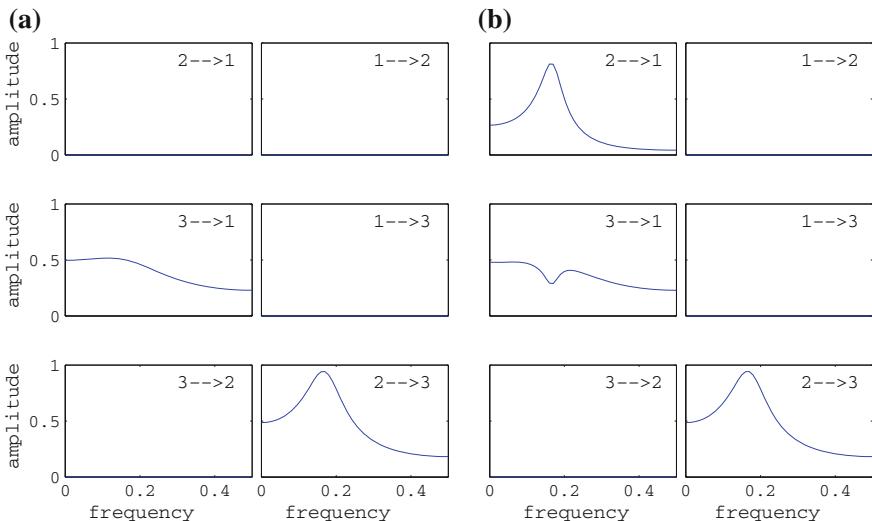


Fig. 8.4 **a** The plot of partial directed coherence (PDC) and **b** directed transfer function (DTF) computed using the model MVAR coefficients in Eq.(8.122). Explanations on the ordinate and abscissa variables are found in the caption for Fig. 8.2

source has a directional influence on the first source. The PDC and DTF computed using the model MVAR coefficients which appear in Eq. (8.122) are shown in Fig. 8.4. The results in Fig. 8.4 show the ground truth in the following experiments.

The signal part of simulated sensor recordings, $\mathbf{b}_s(t)$, was generated by projecting the time series of the three sources onto the sensor space using the sensor lead field. The final form of simulated MEG recordings $\mathbf{b}(t)$ was computed by adding spontaneous MEG to $\mathbf{b}_s(t)$, such that $\mathbf{b}(t) = \mathbf{b}_s(t) + \varrho \mathbf{b}_I(t)$, where $\mathbf{b}_I(t)$ is the spontaneous MEG measured using the same sensor array, and ϱ is a constant that controls the signal-to-interference ratio(SIR) of the generated sensor recordings. We first set ϱ in order for SIR to be equal to 8, and computed the simulated MEG recordings $\mathbf{b}(t)$. The Champagne source reconstruction algorithm was applied to $\mathbf{b}(t)$, and the estimated time series of the three sources, $\hat{s}_1(t)$, $\hat{s}_2(t)$, and $\hat{s}_3(t)$ were obtained. The MVAR coefficients were estimated using the least-squares method. The results of computing PDC and DTF are shown in Fig. 8.5.

In Fig. 8.5, results very close to the ground truth in Fig. 8.4 were obtained. The DTF detects the information flow from the second to the first sources, which is the indirect causal coupling via $s_3(t)$, but the PDC does not detect this indirect coupling. These results are consistent with the explanation in Sect. 8.5.1.

We next generated the simulated sensor recordings with setting SIR at 2. The results of computing the PDC and DTF are shown in Fig. 8.6. Here, the MVAR coefficients were estimated using the least-squares method. The figure shows that large spurious causal relationships exist, due to the low SIR of the generated data. We then applied the sparse Bayesian algorithm described in Sect. 8.7.2 for estimating the MVAR coefficients, and computed the PDC and DTF. The results are shown in

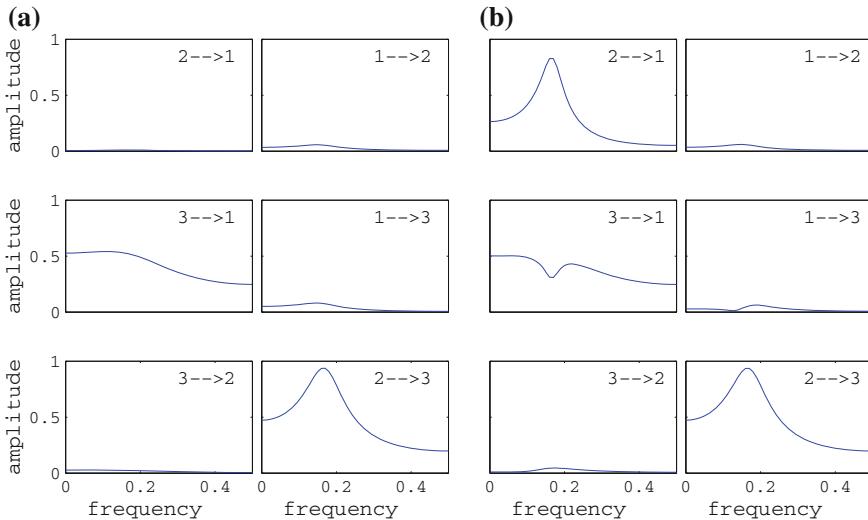


Fig. 8.5 **a** The plot of partial directed coherence (PDC) and **b** the directed transfer function (DTF) computed using MVAR coefficients estimated from the simulated MEG data with SIR equal to 8. The least-squares method in Sect. 8.7.1 was used for estimating MVAR coefficients. Explanations on the ordinate and abscissa variables are found in the caption for Fig. 8.2

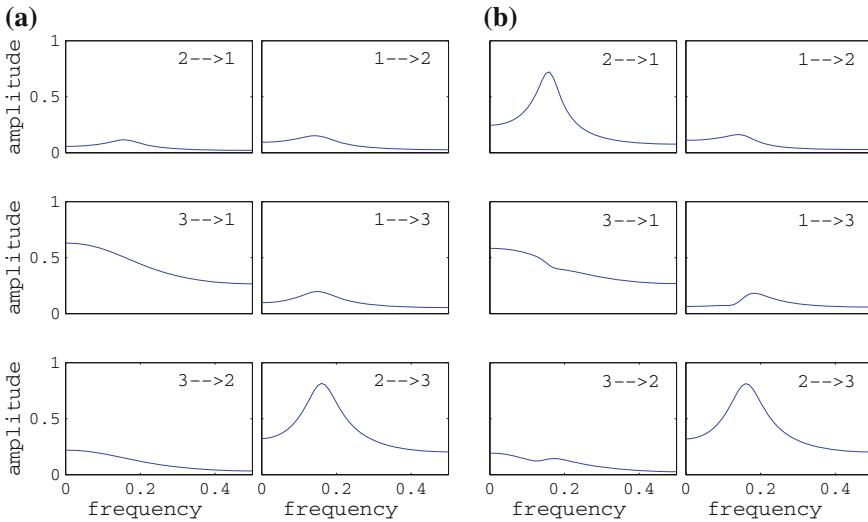


Fig. 8.6 **a** Plots of partial directed coherence (PDC) and **b** directed transfer function (DTF) computed using MVAR coefficients estimated from the simulated MEG data with SIR equal to 2. The least-squares method in Sect. 8.7.1 was used for estimating MVAR coefficients. Explanations on the ordinate and abscissa variables are found in the caption for Fig. 8.2

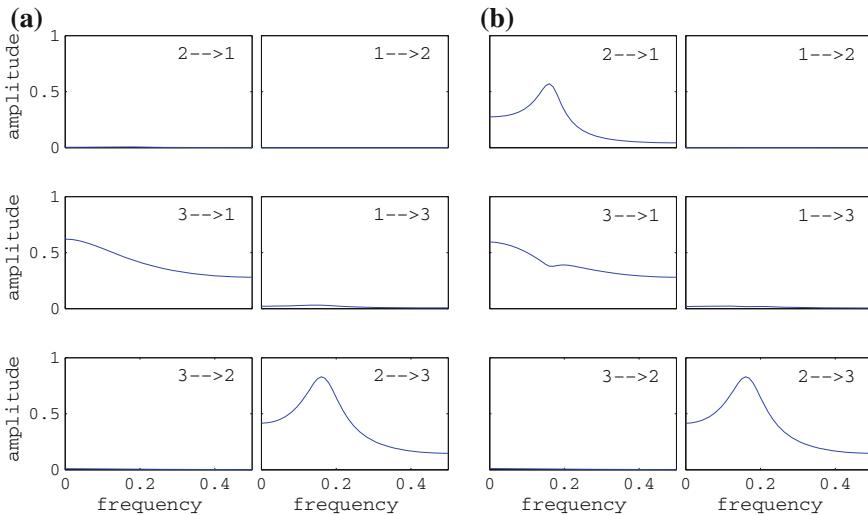


Fig. 8.7 **a** Plots of partial directed coherence (PDC) and **b** the directed transfer function (DTF) computed using MVAR coefficients estimated from the simulated MEG data with SIR equal to 2. The sparse Bayesian algorithm described in Sect. 8.7.2 was used for estimating MVAR coefficients. Explanations on the ordinate and abscissa variables are found in the caption for Fig. 8.2

Fig. 8.7. The results are very close to the ground truth in Fig. 8.4, demonstrating the effectiveness of the sparse Bayesian method in the MVAR estimation.

References

1. C.W. Granger, Investigating causal relations by econometric models and cross-spectral methods. *Econom.: J. Econom. Soc.* **37**(3), 424–438 (1969)
2. C. Granger, Investigating causal relations by econometric models and cross-spectral methods. *Econom. Soc. Monogr.* **33**, 31–47 (2001)
3. J. Geweke, Measures of conditional linear dependence and feedback between time series. *J. Am. Stat. Assoc.* **77**, 304–313 (1982)
4. J. Geweke, Inference and causality in economic time series models, in *Handbook of Econometrics*, vol. 2 (Elsevier, Amsterdam, 1984), pp. 1101–1144
5. M.J. Kamiński, K.J. Blinowska, A new method of the description of the information flow in the brain structure. *Biol. Cybern.* **65**, 203–210 (1991)
6. L.A. Baccalá, K. Sameshima, Partial directed coherence: a new concept in neural structure determination. *Biol. Cybern.* **84**, 463–474 (2001)
7. R. Vicente, M. Wibral, M. Lindner, G. Pipa, Transfer entropy—a model-free measure of effective connectivity for the neurosciences. *J. Comput. Neurosci.* **30**, 45–67 (2011)
8. K.J. Blinowska, M. Kamiński, Multivariate signal analysis by parametric models, in *Handbook of Time Series Analysis: Recent Theoretical Developments and Applications* (2006), p. 373
9. P. Franaszczuk, K. Blinowska, M. Kowalczyk, The application of parametric multichannel spectral estimates in the study of electrical brain activity. *Biol. Cybern.* **51**(4), 239–247 (1985)
10. C.D. Meyer, *Matrix Analysis and Applied Linear Algebra* (Society for Industrial and Applied Mathematics, Philadelphia, 2000)
11. M. Ding, Y. Chen, S.L. Bressler, Granger causality: basic theory and application to neuroscience, in *Handbook of Time Series Analysis*, ed. by B. Schelter et al. (Wiley-VCH, Weinheim, 2006), pp. 500–600
12. K. Sameshima, L.A. Baccalá, Using partial directed coherence to describe neuronal ensemble interactions. *J. Neurosci. Methods* **94**(1), 93–103 (1999)
13. T. Schreiber, Measuring information transfer. *Phys. Rev. Lett.* **85**, 461–464 (2000)
14. L. Barnett, A.B. Barrett, A.K. Seth, Granger causality and transfer entropy are equivalent for Gaussian variables. *Phys. Rev. Lett.* **103**, 238701 (2009)
15. W.D. Penny, S.J. Roberts, Bayesian multivariate autoregressive models with structured priors. *IEE Proc.* **149**, 33–41 (2002)

Chapter 9

Detection of Phase–Amplitude Coupling in MEG Source Space: An Empirical Study

9.1 Introduction

Neural oscillations across multiple frequency bands have consistently been observed in EEG and MEG recordings. The alpha rhythm (8–13 Hz), the best known brain oscillation, is observed throughout the brain. Other well-known oscillations include the delta (1–3 Hz), theta (4–7 Hz), beta (13–30 Hz), and gamma (>30 Hz) bands. Gamma-band oscillations are considered to be those most closely associated with firing of cortical neurons [1].

Recently, neural substrates responsible for the genesis of such brain rhythms and their functional role in brain information processing have been the subject of intense investigations. As a result, a number of hypotheses for explaining the role of neural oscillations have emerged, some including the concept that sensory input information is coded using multiple temporal scales [2–6]. One of these hypotheses claims that temporal coding through regulation of oscillators’ firing is essential for the effective exchange of information in the mammalian brain [7]. Another hypothesis claims that temporal coding organized by coupled alpha and gamma oscillations prioritizes visual processing [8].

A recent study [9] suggests that information gating via selective attention is closely related to temporal encoding. Another study suggests that temporal encoding combined with the encoding by the phase of firing is expected to carry more information than temporal encoding by spike rates or spike patterns [10]. Therefore, temporal encoding due to phase of firing is considered to be one of the most plausible explanations behind mechanisms of information processing in the brain.

On the other hand, if the brain uses temporal encoding based on the phase of firing, there should be an oscillator that regulates neural firing based on its phase. Therefore, the hypothesis for temporal encoding based on phase of firing suggests a possible coupling between oscillators with different frequencies. That is, the amplitude

The authors of this chapter are Eiichi Okumura, Ph.D. and Takashi Asakawa, Ph.D. who are with Neuroscience Project Office, Innovation Headquarters, Yokogawa Electric Corporation.

of the high-frequency oscillator increases at a certain phase of the low-frequency oscillator. Such relationship between oscillators with different frequencies is called cross-frequency coupling [11], which is the main topic of this chapter.

The aim of this chapter is to demonstrate that MEG source space analysis is able to detect such cross-frequency interactions. It is a widely accepted notion that low-frequency oscillations, such as alpha, are “encoders” of information processing, and high-frequency oscillations, such as gamma, represent particular neuronal activities. Therefore, analysis of cross-frequency interactions between different cortical regions would be extremely important for the investigation of information processing mechanisms in a human brain.

9.2 Types of Cross-Frequency Coupling

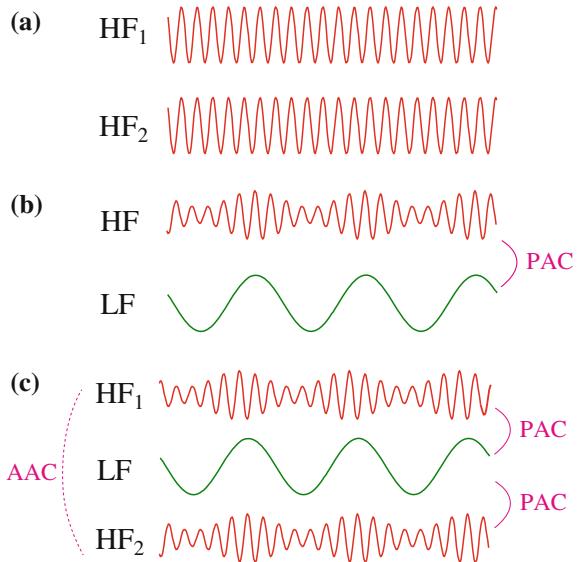
Connections between brain oscillators at the same frequency, termed as same-frequency coupling (SFC) in this chapter, have been investigated in various functional connectivity studies. In such studies, the coherence (or phase-locking value) between target oscillations is commonly used as a measure of connectivity.¹ In contrast to SFC, the cross-frequency coupling (CFC) is a phenomenon in which oscillations with different frequencies interact with each other. In general, CFC is classified into the following coupling types between the oscillators:

- Phase–amplitude coupling (PAC): The phase of a low-frequency (LF) oscillation drives the amplitude of a high-frequency (HF) oscillation with the highest amplitude of the HF oscillation occurring at a specific phase of the LF oscillation, which is referred to as the “preferred coupling phase”.
- Amplitude–amplitude coupling (AAC): The amplitude of an oscillation correlates the amplitude of other oscillations. Even if the oscillators are not directly coupled, their amplitude can be co-modulated.
- Phase–phase coupling (PPC): The phase of an LF oscillation is related to the phase of an HF oscillation. This is sometimes referred to as “ $n:m$ phase locking” where phase locking occurs between n cycles of the HF oscillation and m cycles of the LF oscillation.

Figure 9.1 depicts typical time courses of oscillations when a coupling exists. When a pair of oscillations is interacting with each other, the interaction causes a constant phase relationship between the two HF oscillations, resulting in same-frequency coupling. This situation is depicted in Fig. 9.1a. When phase–amplitude coupling exists between two oscillators, the phase of the LF oscillation regulates the amplitude of the HF oscillation. In this case, the HF oscillation increases in amplitude at a certain phase of the LF oscillation. Such a situation is depicted in Fig. 9.1b. One special case of the PAC is that, when a common LF oscillation simultaneously drives the amplitudes of two HF oscillators, amplitude–amplitude coupling (AAC) between

¹ Detection of SFC in the MEG source space is the scope of Chap. 7.

Fig. 9.1 **a** Typical time courses of two oscillators when same-frequency coupling (SFC) exists. **b** Typical time courses of two oscillators when phase-amplitude coupling (PAC) exists. **c** Typical time courses of two oscillators when PAC exists and a common lower frequency (LF) component simultaneously drives the amplitudes of two higher frequency (HF) oscillators. In this case, an amplitude-amplitude coupling (AAC) between the two HF oscillations exists even when they are not directly coupled



the two oscillators exists. Namely, AAC can exist, even when two HF oscillators are not directly coupled. This situation is depicted in Fig. 9.1c.

This chapter is based on the hypothesis that, in brain information processing, the phases of the low-frequency oscillations encode the amplitudes of the high-frequency oscillations that represent particular neuronal activities. If this hypothesis is true, analyzing the relationship between the phases of the LF oscillations and the amplitudes of the HF oscillations, i.e., analyzing the phase–amplitude coupling (PAC) between the LF and HF signals is extremely important in investigating brain information processing. Therefore, we focus on the PAC and related topics in this chapter.

9.3 Local PAC and Cross-Location PAC

When phase–amplitude coupling is observed between the HF and LF signals measured at the same location in a brain, this PAC is called the local PAC, and termed IPAC. On the other hand, when the PAC is observed between the HF and LF signals measured at different locations, this PAC is called the cross-location PAC, and here termed xPAC.

So far, most studies have investigated IPAC, but some recent studies have extended their scope to xPAC. For example, in a study using an electrocorticogram (ECoG) recording, the presence of xPACs during working memory tasks was reported [12]. The xPAC between cortical gamma oscillations and an LF oscillation from the

thalamus was reported in [13, 14]. We also show, in Sect. 9.5, that lPAC and xPAC can be detected from a healthy subject using hand-motor MEG data.

9.4 Quantification of Phase–Amplitude Coupling

9.4.1 Instantaneous Amplitude and Phase

In the PAC analysis, it is essential to calculate the instantaneous amplitude of the high-frequency signal and the instantaneous phase of the low-frequency signal. The procedure to obtain these quantities in the source space analysis is summarized as follows. Let us define the voxel time course as $x(t)$. The voxel time course is band-pass filtered into specific HF and LF bands of interest, resulting in the LF signal $x_L(t)$ and the HF signal $x_H(t)$.

The Hilbert transform is then applied to extract their instantaneous amplitudes and phases, such that

$$\mathcal{A}[x_L(t)] = A_L(t) \exp[i\theta_L(t)], \quad (9.1)$$

$$\mathcal{A}[x_H(t)] = A_H(t) \exp[i\theta_H(t)], \quad (9.2)$$

where $\mathcal{A}[\cdot]$ indicates an operator that converts a real-valued signal into its analytic signal. In the equations above, $A_L(t)$ and $A_H(t)$ are the instantaneous amplitudes of the LF and HF signals, and $\theta_L(t)$ and $\theta_H(t)$ are the instantaneous phases of the LF and HF signals. Once amplitudes and phases of the LF and HF signals are obtained, we proceed with the evaluation of the strength of the PAC using measures described in the following subsections.

9.4.2 Amplitude–Phase Diagram

To evaluate the strength of the PAC, it is necessary to derive an empirical distribution of the HF signal amplitude, $A_H(t)$, with respect to the phase of the LF signal, $\theta_L(t)$. To obtain such an empirical distribution, we observe $A_H(t)$ for a time window containing several cycles of the LF signal, and the amplitude values $A_H(t)$ are classified according to the phase of the LF signal. Let us divide the phase value into total q intervals such as $\Delta_1, \Delta_2, \dots, \Delta_q$, where $\Delta_j = [2\pi \frac{j-1}{q}, 2\pi \frac{j}{q}]$. There are multiple (and usually many) values of $A_H(t)$ obtained when $\theta_L(t)$ has a value within Δ_j . The mean of such values of $A_H(t)$ is computed, and it is denoted Ψ_j , i.e.,

$$\Psi_j = \langle A_H(t) \rangle_{\theta_L(t) \in \Delta_j}, \quad (9.3)$$

where $\langle \cdot \rangle_{\theta_L(t) \in \Delta_j}$ indicates taking the mean of amplitude values obtained when $\theta_L(t)$ belongs to Δ_j . The mean amplitude Ψ_j is computed for all the phase bins $\Delta_1, \dots, \Delta_q$, and the resultant set of Ψ_1, \dots, Ψ_q represents the distribution of the mean HF signal amplitude with respect to the phase of the LF signal.

The plot of these mean amplitudes Ψ_j with respect to the phase bins is called the amplitude–phase diagram, which directly expresses the dependence of the HF signal amplitude on the phase of the LF signal. Therefore, if no PAC occurs, $A_H(t)$ does not depend on $\theta_L(t)$. Thus, the value of Ψ_j is independent from the index j , resulting in a uniform amplitude–phase diagram. On the contrary, if a PAC occurs, the amplitude of the HF signal becomes stronger (or weaker) at specific phase values of the LF signal, and the amplitude–phase diagram deviates from the uniform distribution. We show several examples of the amplitude–phase diagram in Sect. 9.5.

9.4.3 Modulation Index (MI)

The modulation index [15] is a measure that quantizes the deviation of the amplitude–phase diagram from a uniform distribution. To compute the modulation index, we first normalize the mean amplitude values, Ψ_1, \dots, Ψ_q , by their total sum. That is, the normalized mean amplitude value for the j th phase bin, expressed as $p(j)$, is obtained as

$$p(j) = \frac{\Psi_j}{\sum_{i=1}^q \Psi_i}. \quad (9.4)$$

This $p(j)$ is the normalized version of the amplitude–phase diagram, which can be interpreted as the empirically derived probability distribution of the occurrence of $A_H(t)$ obtained when $\theta_L(t)$ belongs to Δ_j . If no PAC occurs, there is no specific relationship between $A_H(t)$ and $\theta_L(t)$, resulting in $p(j)$ having a uniform value. Thus, the difference of the empirical distribution $p(j)$ from the uniform distribution can be a measure of the strength of PAC. The modulation index, \mathcal{M}_I , employs the Kullback–Leibler distance to assess this difference, and is defined as

$$\mathcal{M}_I = \mathcal{K}[p(j) \| u(j)] = \sum_{j=1}^q p(j) \log \left[\frac{p(j)}{u(j)} \right], \quad (9.5)$$

where $u(j)$ is the uniform distribution, which is $u(j) = 1/q$ ($j = 1, \dots, q$).

The time variation of the modulation index expresses temporal dynamics of the PAC, which represents the changes in the functional states of a brain. To compute the modulation index time variation, the source time course is divided into multiple time windows, and a value of the modulation index is obtained from each time window. Assuming that a total of K windows are used, we obtain a series of modulation index values, $\mathcal{M}_I(t_1), \mathcal{M}_I(t_2), \dots, \mathcal{M}_I(t_K)$, which expresses the temporal change of the PAC strength.

9.4.4 Phase-Informed Time-Frequency Map

A “phase-informed” time-frequency map is the time-frequency map of the HF signal $x_H(t)$ displayed with the overlay of the sinusoidal time course $\cos(\theta_L(t))$. By displaying the time-frequency map of $x_H(t)$ in this manner, the relationship between the amplitude of the HF signal and the phase of the LF signal can be visualized in an intuitive manner. Particularly, the preference of the HF signal amplitude to the phase of the LF signal is clearly shown. This type of map was first introduced by Canolty et al. [11] and it is often called as the Canolty map. Examples of the phase-informed TF map obtained using hand-motor MEG data are shown in Sect. 9.5.

9.5 Source Space PAC Analysis: An Example Study Using Hand-Motor MEG Data

9.5.1 Experimental Design and Recordings

Here, we present results of an example study to demonstrate the noninvasive detection of PAC using the MEG source space analysis. We used a hand-motor task similar to the task used in [16]. A healthy, right-handed subject (30 years old male) participated in the MEG study.² The subject was asked to perform grasping movements (grasping a soft toy) with his left hand every 10 s. An auditory cue was provided to inform the grasp timing. The subject repeated grasping 64 times with his eyes open.

MEG data was acquired using a whole-head 160-channel MEG system (PQA160C, Yokogawa Electric Corporation, Japan) with a sample rate of 2,000 Hz. A band-pass filter with a bandwidth of 0.16–500 Hz was applied. Electromyogram was recorded from the subject’s left forearm to monitor subject’s hand movements. During the recordings, the subject lay in the supine position on a bed in a magnetically shielded room. The subject’s T1-weighted MRI was obtained, which consists of 166 slices of 1.2 mm thickness. A single-slice image has 512×512 voxels in a field of view of 261×261 mm.

9.5.2 Data Analysis

9.5.2.1 Preprocessing

An ICA-based interference removal procedure was applied to remove eyeblink and MCG artifacts. An offline digital notch filter with the frequency of 60, 120, and 180 Hz

² Informed consent was obtained from the subject before the experiment. Handedness was determined using the Edinburgh Handedness Inventory [17].

was applied to remove the power line interference. Continuous MEG recordings were converted into multiple trials. A single trial had a time period between -6 and $+6$ s allowing an overlap of about 2 s to its adjacent trials. Here, the time origin was defined as the time of starting the initial rise in the electromyogram waveform. We excluded trials with the magnetic field strength exceeding 2.5 pT in any sensor channels.

9.5.2.2 Source Estimation

The watershed algorithm in FreeSurfer [18] was applied to the subject's MRI to generate triangulations on the inner skull surface. The dense triangulation on the folded cortical surface was divided into a grid of 15,000 voxels. The source space includes both hemispheres. The average distance between voxels is 2 mm. The sensor lead field was computed using the boundary element model provided by OpenMEG [19]. The source estimation was performed using the depth-weighted minimum L2 norm method implemented in Brainstorm [20]. The anatomical labels were identified automatically by using the FreeSurfer software with the Desikan–Killiany atlas [21]. Since in our study, we focused only on the contralateral sensorimotor area, we defined the region of interest as one containing subareas such as precentral, postcentral, paracentral, caudal media frontal areas, and a part of the superior frontal area. The total number of voxels in the ROI resulted in 1,258. We computed a single time course at each voxel using the procedure described in Sect. 7.2 in Chap. 7. Since the estimated source orientation has a 180° ambiguity, we use either the estimated orientation or its opposite orientation that is closer to the direction of the cortical surface provided by Freesurfer.

9.5.2.3 Modulation Index Analysis

The LF and HF signals are extracted from voxel time courses by applying a finite impulse response band-pass filter. The frequency band of the LF signal ranges from 8.5 to 14.5 Hz, which was divided into 11 sub-bands with a bandwidth of 1 Hz; each sub-band half-overlapped with adjacent sub-bands. The frequency band of the HF signal ranges from 50 to 200 Hz; the band is divided into five half-overlapped sub-bands with the bandwidth of 50 Hz. We computed the modulation index described in Sect. 9.4.3 using all the combinations of the LF and HF sub-bands. The number of the phase bins, q , was set to 12.

When computing the time course of the modulation index, we divided a 12-second-long interval of trial data into 21 time windows in which the window width is 2 s with a 1.5 s overlap. The modulation index \mathcal{M}_I was computed using all 21 sets of the time-windowed data to obtain the time course, $\mathcal{M}_I(t_1), \dots, \mathcal{M}_I(t_{21})$, where t_1, \dots, t_{21} indicates the midpoints of the 21 time windows. We also computed the maximum of the modulation index, such that

$$\mathcal{M}_P = \max (\mathcal{M}_I(t_k)) \quad (9.6)$$

This \mathcal{M}_P is used as a measure of the strength of the PAC. The surrogate-data method, similar to the one described in Sect. 7.8, is applied to assess the statistical significance of \mathcal{M}_P .

9.5.3 Results of Local PAC Analysis

9.5.3.1 Time Course Estimation and Spatial Mapping of the Modulation Index

We first analyzed local PAC (IPAC), and computed the modulation index time courses for all the sub-bands of the LF signal at all voxel locations. The bandwidth of the HF signal was fixed to 75 ± 25 Hz. The resultant time courses were classified using K-means clustering analysis, which thereby found three types of clusters. The first kind of the cluster, called type I, consists of the time courses that peaked before the onset of the hand movements. The time course averaged across all the type I time courses is shown in Fig. 9.2a. We computed \mathcal{M}_P in Eq. (9.6) using the type I time courses, and assessed the statistical significance of \mathcal{M}_P at all voxel locations. Voxels containing statistically significant IPAC activity were found with the LF signal of 12.5 ± 0.5 Hz near the primary motor area. These voxels are shown in Fig. 9.2b.

The second kind of the modulation index time courses, called type II, peaked immediately before or during the hand movements. The mean time course of this type is shown in Fig. 9.2c. Voxels containing statistically significant \mathcal{M}_P (computed using the type II time courses) were found with the LF signal of 10.5 ± 0.5 Hz near the lateral part of the precentral area, as shown in Fig. 9.2d. The third kind of time courses, type III, peaked after the hand movements. The mean time course is shown in Fig. 9.2e. Voxels containing significant \mathcal{M}_P (computed using the type III time courses) were found with the LF signal of 12.5 ± 0.5 Hz near the postcentral or central sulcus areas. These voxels are shown in Fig. 9.2f.

In summary, we found significant IPAC activities near the contralateral precentral area before the execution. However, following the execution, IPAC activities were found near the postcentral area and central sulcus. These findings may indicate that a certain relationship exists between these three types of IPAC activities, and the preparation and rebound of the motor and sensory systems.

9.5.3.2 Amplitude–Phase Diagram and Phase-Informed TF Map

An amplitude–phase diagram was derived using a voxel time course in the area labeled as Region A in Fig. 9.3a. Note that Region A shows significant IPAC activity showing the type I time course. In this derivation, the LF signal was set to 12.5 ± 0.5 Hz, and the HF signal to 75 ± 25 Hz. The data in the time window from -4.5 to -2.5 s (when the type-I IPAC was nearly maximum) was used. The resultant amplitude–phase diagram is shown in Fig. 9.3b in which each bar indicates the mean

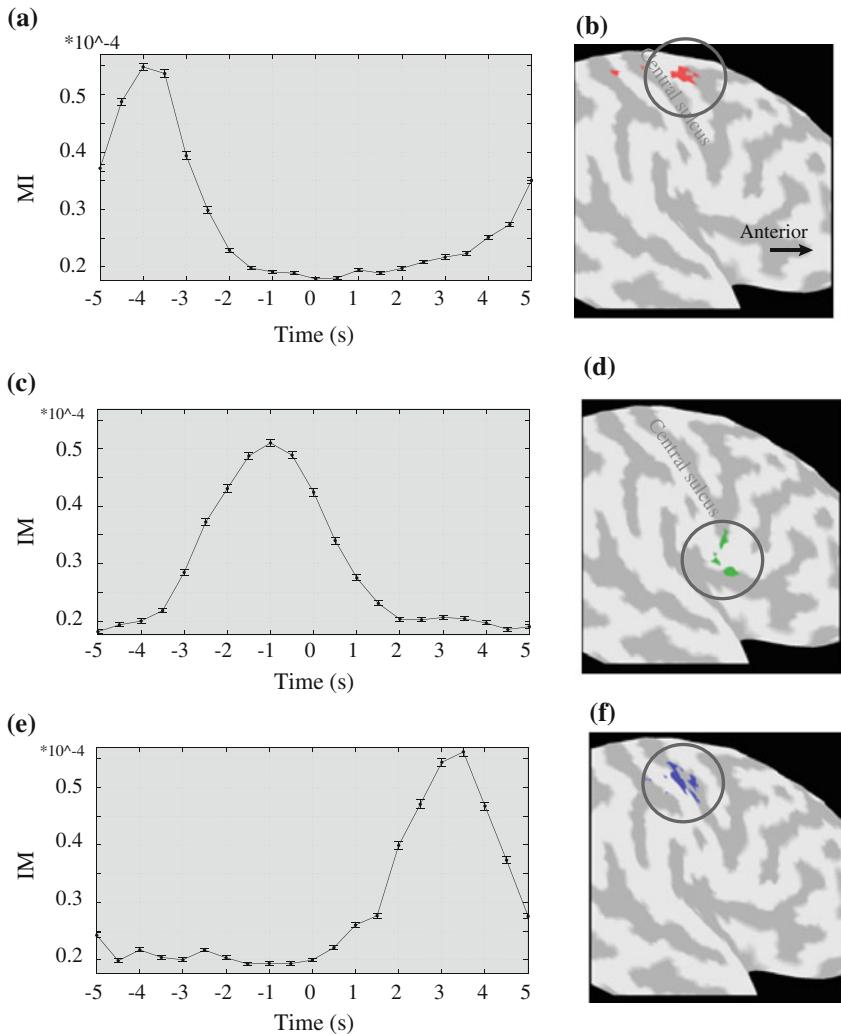


Fig. 9.2 **a** Modulation index time course of the type-I IPAC activity. **b** Voxels containing statistically significant IPAC activity having the type I time course. **c** Modulation index time course of the type II IPAC activity. **d** Voxels containing statistically significant IPAC activity having the type II time course. **e** Modulation index time course of the type III IPAC activity. **f** Voxels containing statistically significant IPAC activity having the type III time course. The significant activities in (b) and (f) were found when the LF signal of 12.5 ± 0.5 Hz and the HF signal of 75 ± 25 Hz were used. The significant activities in (d) were found when the LF signal of 10.5 ± 0.5 Hz and the HF signal of 75 ± 25 Hz were used. In (a), (c), and (e), the average time course is shown and the *error bar* indicates the range of \pm standard deviation. Here, the time origin is equal to the onset of the hand movement. In (b), (d), and (f), the *blank circles* are used to indicate the locations of the voxels with significant activity

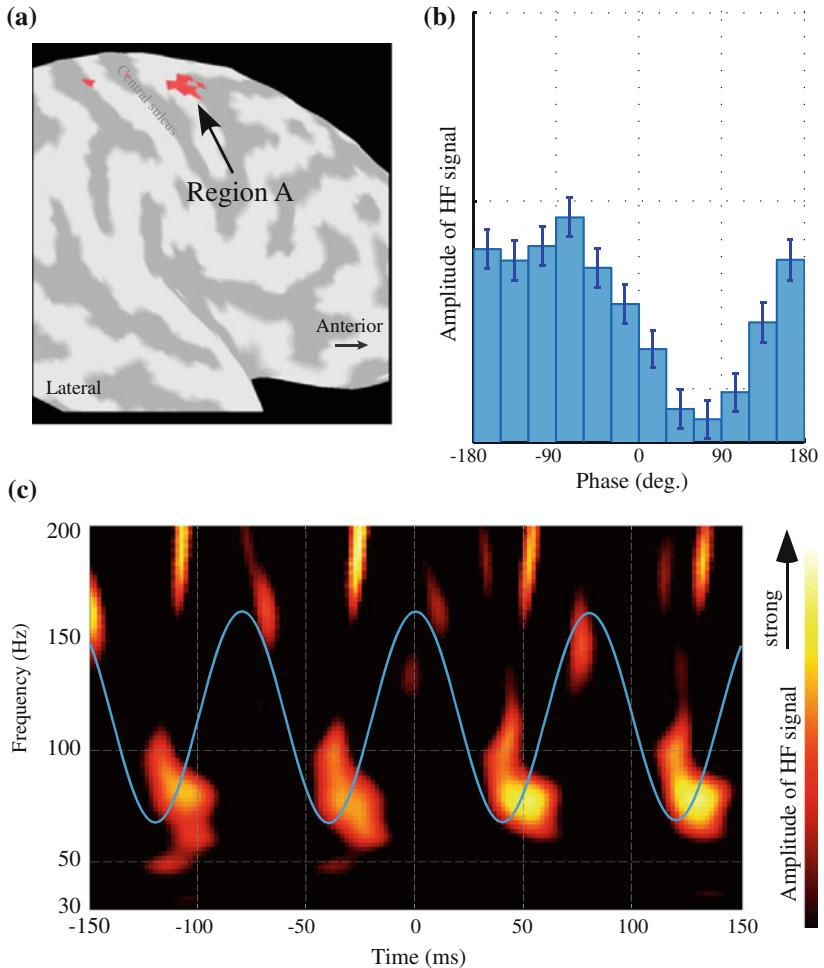


Fig. 9.3 Results of local PAC (IPAC) analysis. **a** Location of Region A where a significant type I IPAC activity was observed. The IPAC activity at Region A was analyzed. **b** Amplitude–phase diagram. Each bar indicates the mean amplitude value at that phase bin. The error bar indicates the \pm standard deviation. In this analysis, the LF signal was set to 12.5 ± 0.5 Hz, and the HF signal was set to 75 ± 25 Hz. The data in the time window from -4.5 to -2.5 s was analyzed. **c** The phase-informed TF map. The amplitude of the HF signal is color-coded and displayed according to the color bar. The sinusoidal plot shows $\cos(\theta_L(t))$ where $\theta_L(t)$ is the instantaneous phase of the LF signal. Only the portion between ± 150 ms from the LF peaks in the time window is displayed

amplitude value Ψ_j at that phase bin. It can be seen that the distribution of Ψ_j significantly deviates from the uniform distribution, indicating that a strong IPAC activity occurs. The preferred phase of the HF signal can be observed. That is, the amplitude of the HF signal increases significantly at a specific phase of the LF signal around -150° , and decreases around $+60^\circ$.

The relationship between the HF signal amplitude and the LF signal phase can be seen more clearly in the phase-informed TF map in Fig. 9.3c. Here, the same voxel data from Region A has been used. The phase-informed TF map intuitively shows the preference of the HF signal amplitude to a specific phase of the LF signal. The map clearly indicates that the amplitude of the HF signal from 60 to 100 Hz significantly increases near -150° of the phase of the LF signal.

9.5.3.3 LPAC Temporal Trend and Comparison to Event-Related Power Changes

Event-related power changes such as desynchronization (ERD) and synchronization (ERS) are representative aspects of brain activities. Comparison between PAC temporal trends and ERD/ERS should provide useful insights into brain information processing. Using the IPAC signal from Region A, we computed modulation index time courses with all 11 sub-bands of the LF signal. The HF signal was fixed to the frequency band of 75 ± 25 Hz. The results are shown in Fig. 9.4a. In this figure, the sizes of the circles represent the strength of the modulation index, and their colors represent the LF sub-bands.

The results in Fig. 9.4a show that LF signals at two sub-bands, 10 ± 0.5 and 12.5 ± 0.5 Hz, have significantly large modulation index values. The IPAC with LF signal of 12.5 ± 0.5 Hz reaches its maximum between -5 and -3 s but the IPAC with the LF signal of 10 ± 0.5 Hz reaches its maximum between -0.5 and $+0.5$ s. This observation suggests the presence of multiple brain rhythm scales underlying motor planning and execution.

A time-frequency representation of the broadband MEG signal from Region A is shown in Fig. 9.4b. In this figure, we can observe that the ERD whose center frequency is approximately equal to 12.5 Hz starts near -3 s. On the other hand, Fig. 9.4a shows that IPAC with LF of 12.5 ± 0.5 Hz starts decreasing around -3 s. Comparison between Fig. 9.4a, b shows that a strong 12.5 Hz ERS starts at $+3.5$ s at which the IPAC with the LF signal of 12.5 ± 0.5 Hz restarts. On the other hand, a stable IPAC activity with the LF band between 9 and 11 Hz (marked by the dotted circle in Fig. 9.4a) is observed at the hand-movement period when a significantly strong ERD is observed at the same frequency band near 10 Hz. Such coincidence between ERD/ERS activities and IPAC time courses suggests that IPAC and ERD/ERS activities are interrelated.

9.5.4 Results of Analyzing Cross-Location PACs

Finally, the cross-location PAC (xPAC) was analyzed using the same hand-motor MEG data. In this analysis, the LF signal is measured in voxels of the precentral area; these voxels are labeled Voxel B in Fig. 9.5a. The HF signal is measured at the voxel near the primary motor area; the voxel is labeled Voxel A in Fig. 9.5a. We first

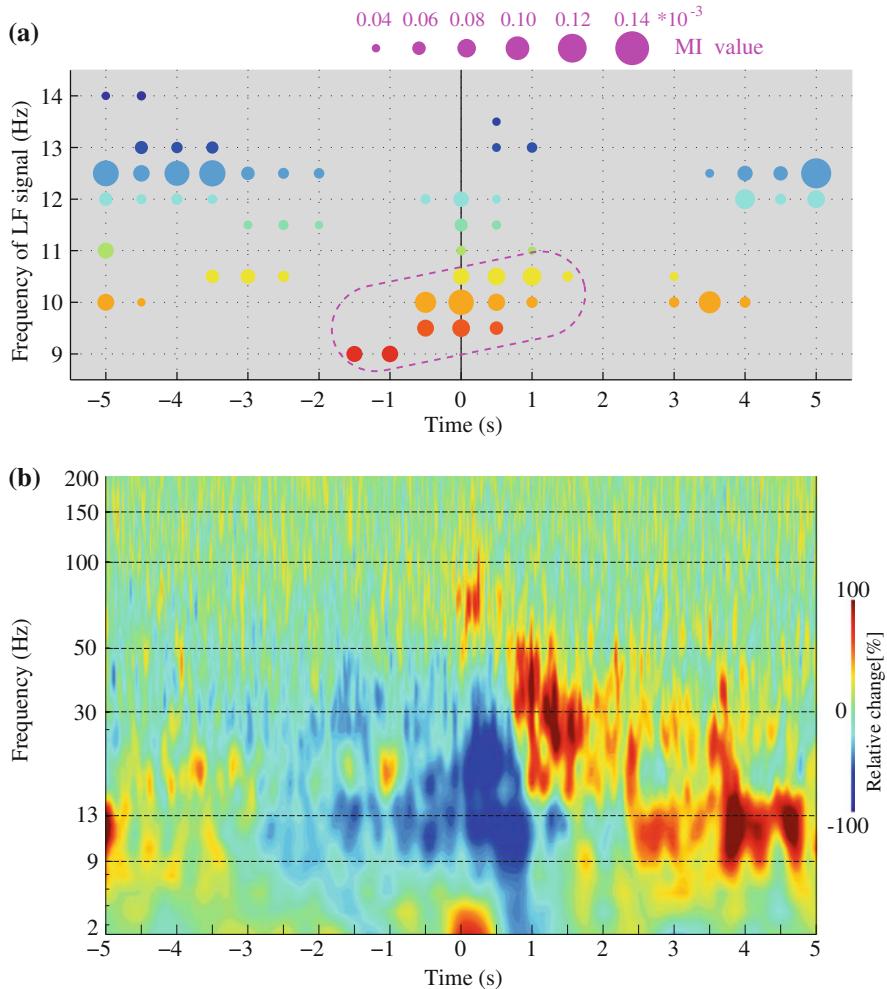


Fig. 9.4 **a** Modulation index (MI) time courses for the 11 sub-bands of the LF signal. The IPAC activity at Region A was analyzed, and the HF signal is fixed to 75 ± 25 Hz. Sizes of the circles represent the values of the modulation index, and their colors represent the LF sub-bands. A dotted circle indicates an IPAC activity with the LF band near 10 Hz observed when a significantly strong ERD is observed near 10 Hz. **b** The time-frequency representation of the (broadband) voxel time courses at Region A. The TF representation shows relative power change from the baseline, which is color-coded and displayed according the color bar. The time origin is defined as the onset of hand grasping

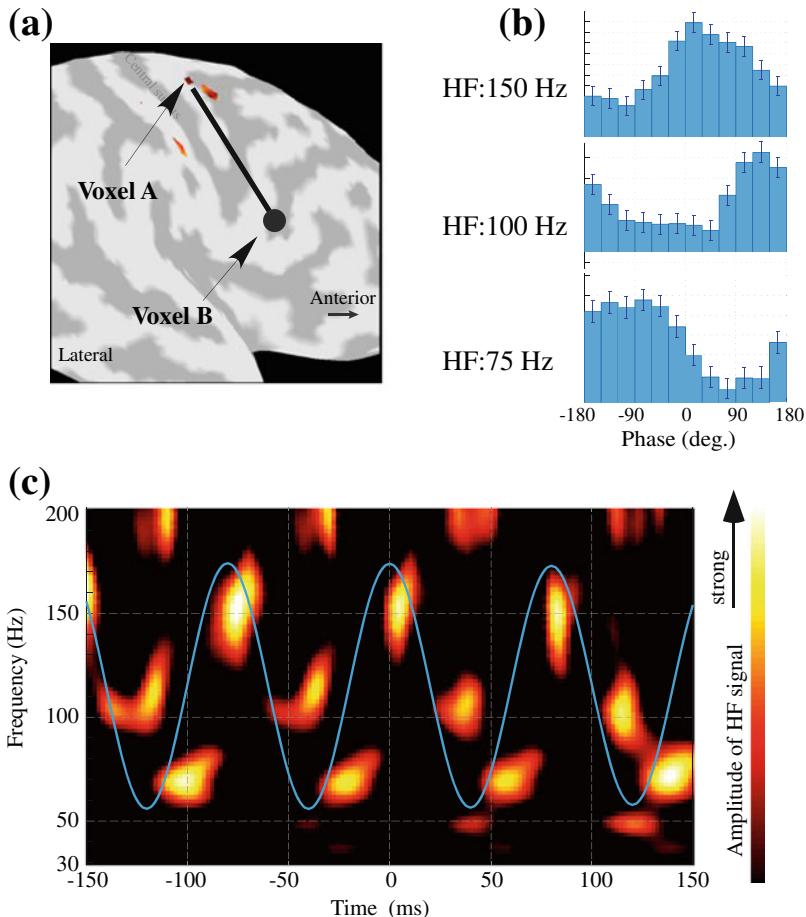


Fig. 9.5 Results of cross-location PAC (xPAC) analysis. **a** Voxel locations for measuring the LF and HF signals. The LF signal is measured at Voxel B and the HF signal is measured at Voxel A. **b** The amplitude–phase diagram for three HF sub-bands: 75 ± 25 Hz (top panel), 100 ± 25 Hz (middle panel), and 150 ± 25 Hz (bottom panel). Each bar indicates the mean amplitude value at that phase bin. The error bar indicates the \pm standard deviation. The data in the time window from -4.5 to -2.5 s was analyzed and the frequency of the LF signal was fixed at 12.5 ± 0.5 Hz. **c** The phase-informed TF map. The amplitude of the HF signal is color-coded and displayed according to the color bar. The sinusoidal plot shows $\cos(\theta_L(t))$ where $\theta_L(t)$ is the instantaneous phase of the LF signal. Only the portion between ± 150 ms from the LF peaks in the time window is displayed

computed the amplitude–phase diagrams using three sub-bands of the HF signal: 75 ± 25 , 100 ± 25 , and 150 ± 25 Hz. The data in the time window of -4 ± 1 s was used, and the frequency band of the LF signal was fixed at 12.5 ± 0.5 Hz. The results are shown in Fig. 9.5b.

In this figure, significant increases in HF signal amplitudes were observed at specific phases of the LF signal. For the HF band of 75 ± 25 Hz, the amplitude of the HF signal increases significantly near -105° of the LF phase. For the HF band of the 100 ± 25 Hz, the amplitude of the HF signal increases significantly near 135° of the LF phase. For the HF band of the 150 ± 25 Hz, the amplitude of the HF signal reaches its maximum near 15° of the LF phase.

Figure 9.5c shows the phase-informed TF map in which three HF components are observed near 75, 100, and 150 Hz. These HF components have different preferred phases for the LF signal, indicating that the three HF components have different timing dependencies on the same LF signal. Such LF phase dependence of the HF signal amplitude across locations suggests the possibility of dynamic multiplexing neuronal codes that exploit timing differences across regions due to synaptic or propagation delays.

9.6 Summary

We have shown that the phase–amplitude coupling (PAC), both local PAC and cross-location PAC, can successfully be detected in MEG source space analysis. PAC is considered to reflect temporal coding of a brain. The measurement of PAC could therefore be a promising tool for monitoring information processing in the human brain. Specifically, PAC analysis can provide coupling-related information, including the frequency of brain rhythms, and the specific phases of LF components, which are preferred by the HF components. These preferred phases are considered to reflect the timing of information exchange. We believe that exploring xPAC is one of the most promising approaches to reveal the mechanisms of brain information processing. Our studies suggest that MEG source space analysis could be a powerful tool in xPAC studies, because of MEG’s wide coverage of the brain and its noninvasiveness.

References

1. E. Privman, R. Malach, Y. Yeshurun, Modeling the electrical field created by mass neural activity. *Neural Netw.* **40**, 44–51 (2013)
2. A.K. Engel, P. Roelfsema, P. Fries, M. Brecht, W. Singer, Role of the temporal domain for response selection and perceptual binding. *Cereb. Cortex* **7**(6), 571–582 (1997)
3. A.K. Engel, P. Fries, P. König, M. Brecht, W. Singer, Temporal binding, binocular rivalry, and consciousness. *Conscious. Cogn.* **8**(2), 128–151 (1999)
4. P. Lakatos, G. Karmos, A.D. Mehta, I. Ulbert, C.E. Schroeder, Entrainment of neuronal oscillations as a mechanism of attentional selection. *Science* **320**(5872), 110–113 (2008)
5. S. Panzeri, N. Brunel, N.K. Logothetis, C. Kayser, Sensory neural codes using multiplexed temporal scales. *Trends Neurosci.* **33**(3), 111–120 (2010)
6. T. Akam, D.M. Kullmann, Oscillatory multiplexing of population codes for selective communication in the mammalian brain. *Nat. Rev. Neurosci.* **15**(2), 111–122 (2014)
7. G. Buzsaki, *Rhythms of the Brain* (Oxford University Press, New York, 2006)

8. O. Jensen, B. Gips, T.O. Bergmann, M. Bonnefond, Temporal coding organized by coupled alpha and gamma oscillations prioritize visual processing. *Trends Neurosci.* **37**(7), 357–369 (2014)
9. E.M. Zion Golumbic, N. Ding, S. Bickel, P. Lakatos, C.A. Schevon, G.M. McKhann, R.R. Goodman, R. Emerson, A.D. Mehta, J.Z. Simon et al., Mechanisms underlying selective neuronal tracking of attended speech at a cocktail party. *Neuron* **77**(5), 980–991 (2013)
10. C. Kayser, M.A. Montemurro, N.K. Logothetis, S. Panzeri, Spike-phase coding boosts and stabilizes information carried by spatial and temporal spike patterns. *Neuron* **61**(4), 597–608 (2009)
11. R.T. Canolty, E. Edwards, S.S. Dalal, M. Soltani, S.S. Nagarajan, H.E. Kirsch, M.S. Berger, N.M. Barbaro, R.T. Knight, High gamma power is phase-locked to theta oscillations in human neocortex. *Science* **313**(5793), 1626–1628 (2006)
12. R. van der Meij, M. Kahana, E. Maris, Phase-amplitude coupling in human electrocorticography is spatially distributed and phase diverse. *J. Neurosci.* **32**(1), 111–123 (2012)
13. T. Staudigl, T. Zaehle, J. Voges, S. Hanslmayr, C. Esslinger, H. Hinrichs, F.C. Schmitt, H.-J. Heinze, A. Richardson-Klavehn, Memory signals from the thalamus: early thalamocortical phase synchronization entrains gamma oscillations during long-term memory retrieval. *Neuropsychologia* **50**(14), 3519–3527 (2012)
14. T.H. Fitzgerald, A. Valentin, R. Selway, M.P. Richardson, Cross-frequency coupling within and between the human thalamus and neocortex. *Front. Hum. Neurosci.* **7**, 84 (2013)
15. A.B. Tort, R. Komorowski, H. Eichenbaum, N. Kopell, Measuring phase-amplitude coupling between neuronal oscillations of different frequencies. *J. Neurophysiol.* **104**(2), 1195–1210 (2010)
16. T. Yanagisawa, O. Yamashita, M. Hirata, H. Kishima, Y. Saitoh, T. Goto, T. Yoshimine, Y. Kamitani, Regulation of motor representation by phase-amplitude coupling in the sensorimotor cortex. *J. Neurosci.* **32**(44), 15467–15475 (2012)
17. R.C. Oldfield, The assessment and analysis of handedness: the Edinburgh inventory. *Neuropsychologia* **9**(1), 97–113 (1971)
18. B. Fischl, Freesurfer. *NeuroImage* **62**(2), 774–781 (2012)
19. A. Gramfort, T. Papadopoulou, E. Olivi, M. Clerc et al., OpenMEEG: opensource software for quasistatic bioelectromagnetics. *Biomed. Eng. Online* **9**(1), 45 (2010)
20. F. Tadel, S. Baillet, J.C. Mosher, D. Pantazis, R.M. Leahy, Brainstorm: a user-friendly application for MEG/EEG analysis. *Comput. Intell. Neurosci.* **2011**, 8 (2011)
21. R.S. Desikan, F. Ségonne, B. Fischl, B.T. Quinn, B.C. Dickerson, D. Blacker, R.L. Buckner, A.M. Dale, R.P. Maguire, B.T. Hyman et al., An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *NeuroImage* **31**(3), 968–980 (2006)

Appendix A

Bioelectromagnetic Forward Modeling

A.1 Neuronal Basis of MEG and EEG Signals

Neurons in the brain function electrically, as well as chemically, and have associated electric and magnetic fields which can be detected outside the head. We first discuss the anatomical and physiological properties of neurons in the mammalian brain, that contribute to the electric and magnetic fields detected by MEG and EEG. Neuronal currents are considered at a level sufficient for understanding the primary source of these fields.

A.1.1 *Neurons and Synapses*

Neurons are cells that are highly specialized for signal processing and conduction via electrochemical and electrical processes. The morphological structure of a neuron includes a cell body, called the soma, and elaborate branching structures called dendrites and axons that enable communication with sensory receptors, distant neurons, etc. Inputs to a neuron are collected in a continuous fashion by the dendrites, and represented as a spatially and temporally continuous variation of the transmembrane voltage. Multiple inputs are summed in the dendritic tree, and the net input is often represented as transmembrane voltage at the soma. When the somatic voltage reaches some threshold, a discrete voltage pulse is generated, called an action potential, which propagates down the axon as a discrete output. The end of the axon has elaborate branching to enable communication with target neurons.

Neuronal communication occurs for the most part via chemical transmission through synapses, although a small minority of neurons also communicate electrically through gap junctions. When an action potential reaches the presynaptic terminal of a synapse with another neuron, it causes the release of neurotransmitter into the synaptic cleft between two cells. The released neurotransmitter binds with some probability to sites of the postsynaptic terminal, stimulating a flow of current

across the membrane of the postsynaptic neuron. Postsynaptic currents may flow inward or outward, depending upon the type of neurotransmitter involved (glutamate or GABA), and the type of ion that flows in response. Excitatory synapses have the effect of increasing the membrane potential of the postsynaptic cell, i.e., making it more positive, while inhibitory synaptic connections decrease it. Neurons are exclusively inhibitory or excitatory.

A.1.2 *Cortical Anatomy*

The mammalian cerebral cortex is the outer mantle of cells surrounding the central structures, e.g., brainstem and thalamus. It is unique to mammals, and is believed to be necessary for most higher-level brain functions. Topologically the cerebral cortex is comprised of two spherical shells, corresponding to the two hemispheres and connected by the corpus callosum. Cortical thickness varies between 2 and 3 mm in the human, and is folded around the subcortical structures so as to appear wrinkled. Its average surface area is about 3000 cm^2 . It is estimated that there are roughly 10^{11} neurons in the human brain, and 10^{10} of these in the cortex.

The diversity of structures of neurons is extraordinary. Approximately 85 % are pyramidal cells whose dendritic trees have a distinctive, elongated geometry that makes possible the generation of extracellular fields at large distances. The remaining 15 % may be broadly classified as stellate cells, whose dendritic trees are approximately spherical, and make little or no contribution to distant field. Both cell types are interconnected and together form a single dynamical network, but it is believed that the fields at large distances are dominated by pyramidal cells because of their size and number.

Synaptic connections in the cortex are dense. Each cortical neuron receives 10^4 – 10^5 synaptic connections, with most inputs coming from distinct neurons. Pyramidal cells make excitatory connections to both cell types. They make intracortical connections over lengths ranging 0.5–3 mm, and cortico-cortical connections over lengths ranging 1–20 cm. Stellate cells make inhibitory connections to both cell types. They make intracortical connections over lengths ranging only 0.02–0.03 mm, much shorter than pyramidal cells. Thus connections in the cerebral cortex are said to exhibit long-range excitation and short-range inhibition.

An interesting distinction exists between intracortical and cortico-cortical connections. Intracortical connections are made locally between neighboring neurons, such that the probability of a connection between two neurons falls off smoothly as a function of distance. In contrast, corticocortical connections are made via subcortical white-matter fibers, and behave non-locally in the sense that connection to a distant neuron does not imply connections with intermediate neurons. Because of the diversity of scales of these synaptic connections, and the nonlocal nature of the cortico-cortical connections, the cerebral cortex exhibits rich spatio-temporal dynamics spanning a wide range of length and time scales.

A.1.3 Neuronal Currents

Electric currents in biological tissue are primarily due to ions, e.g., K+, Na+, Cl-, Ca2+, etc. These ions flow in response to the local electric field, according to Ohm's law, but also in response to their local concentration gradient, according to Fick's law [1]. In the resting state of the membrane, the concentration gradients and electric field are due to ion channel pumps, which use energy acquired from ATP to move ions across the membrane against their diffusion gradient.

To a good approximation, the concentration of each ion inside and outside the membrane may be assumed constant in time. The transmembrane voltage, however, changes radically in time, the strongest example being the action potential. Thus for the purposes of discussing neural activity, we take the transmembrane potential V_m as the primary dynamic state variable to be considered. By convention, V_m is defined as the potential inside relative to that outside, i.e., $V_m = V_i - V_o$. The current flowing across the membrane maybe be viewed as a function of V_m , and therefore as the basis of the extracellular fields detected by MEG and EEG.

In the resting state of the neuron, the balance between electric and diffusive flows determine the resting membrane potential of the neuron [2]. To consider this balance, we define two related quantities: the charge current density, \mathbf{J} [C/m² s], and the ionic flux \mathbf{j} [mol/m² s]. For a single ionic species, these are trivially related by: $\mathbf{J} = zF\mathbf{j}$, where $z = q/e$ is the signed integer number of charges carried by an ion, and F is the Faraday's constant (96,500 C/mol). The flux \mathbf{j} has two contributions, arising from the local concentration gradient ∇C and the local electric field \mathbf{E} , where C indicates the ionic concentration for different ions defined as the number of ions (in mol) per unit volume. Fick's law states that ions diffuse down their concentration gradient according to the linear relation $\mathbf{j} = -D\nabla C$, where D is the diffusion coefficient.

Furthermore, an ion within an applied electric field \mathbf{E} accelerates initially and reaches a terminal velocity v given by $v = \mu(z/|z|)\mathbf{E}$ where the factor $(z/|z|)$ accounts for the fact that negative ions travel in the opposite direction of the electric field. The quantity μ is called the mobility and the diffusion coefficient, $D = \frac{RT}{|z|F}\mu$, where $R = 8.314 \text{ J/(mol K)}$ is the ideal gas constant. Counting flow of particles through a parallelepiped, the ionic flux is $\mu C \frac{z}{|z|} \mathbf{E}$. This equation is related to the more common way of expressing the flow of charges in an electric field via Ohm's law, which in a volume conductor is written as $\mathbf{J} = \sigma\mathbf{E}$. Therefore, the total ionic flux is the linear sum of its diffusive and electrical contributions:

$$\mathbf{j} = \mathbf{j}_S + \mathbf{j}_E = -\frac{\mu}{|z|} \left[\frac{RT}{F} \nabla C + zC \nabla \Phi \right], \quad (\text{A.1})$$

where the electric potential Φ is defined as $\mathbf{E} = -\nabla\Phi$. This equation can be used to derive the Nernst equation, and the Goldman-Hodgkin-Katz equation [3], for the resting potential of a neuron, and also form the basis for modeling synaptic activity.

A.1.4 Neuronal Populations

Neuronal activity gives rise to extracellular electric and magnetic fields, which are detected in MEG and EEG. The fields generated by a single neuron are much too small to be detected at the scalp, but the fields generated by synchronously active neurons, with advantageous geometric alignment, can be detected. Dendrites are more able than axons to generate fields which are detectable at large distances. Stellate cells have approximately spherical dendritic trees, so the resulting extracellular fields tend to add with all possible orientations, and effectively cancel at distance.

Pyramidal cells have similar dendritic trees, but the tree branches are connected to the cell body (soma) by a long trunk, called the apical dendrite. It is a fortuitous anatomical feature of the cortex that pyramidal cells have their apical dendrites aligned systematically along the local normal to the cortical surface. In this way, the fields of synchronously active pyramidal neurons superimpose geometrically to be measurable at the scalp. Consider an approximately 1 cm^3 region of cortex, containing an order of 10^7 aligned pyramidal cells. If only 1 % of these neurons were synchronously active, then the relative contribution of synchronous to asynchronous neurons would be $10^5 / \sqrt{10^7} \sim 30$. Thus, scalp MEG and EEG are largely dominated by synchronous neural activity.

A.2 Electromagnetic Fields in Conductive Media

A.2.1 Maxwell's Equations

Outside the scalp, we measure the net fields produced by synchronously active neurons. For calculations at this scale, the brain and other head tissues can be considered bulk materials, characterized by properties such as electric conductivity σ . The challenge is to compute magnetic fields and scalp potentials as a function of neuronal source currents in a head-shaped conductive medium. We begin here with the fundamentals of electric and magnetic fields in matter, and later specialize to biological tissue.

The physics of electric and magnetic fields are summarized by Maxwell's equations. In matter the macroscopic fields obey

$$\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon}, \quad (\text{A.2})$$

$$\nabla \cdot \mathbf{B} = 0, \quad (\text{A.3})$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}, \quad (\text{A.4})$$

$$\nabla \times \mathbf{B} = \mu \mathbf{J} + \mu \epsilon \frac{\partial \mathbf{E}}{\partial t}. \quad (\text{A.5})$$

where \mathbf{E} is the electric field, \mathbf{B} the magnetic field, ϵ the dielectric permittivity, and μ the magnetic permeability. The Maxwell's equations above have source terms given by the charge density ρ and the current density \mathbf{J} . Additional contributions arise from the time derivatives of the fields.

According to Eq. (A.1), the electric current density \mathbf{J} is expressed as

$$\mathbf{J} = \mathbf{J}_S + \mathbf{J}_E = \mathbf{J}_S + \sigma \mathbf{E}, \quad (\text{A.6})$$

where \mathbf{J}_E is the “ohmic” current which flows in response to the local electric field, and \mathbf{J}_S caused by the ionic flux \mathbf{j}_S is the “source” current (or often called the “impressed” current), which flows in response to transmembrane concentration gradients whose sum corresponds to the net currents summed over all ions. In Eq. (A.6), σ is the bulk conductivity of the material, which is the main parameter governing the spread of ionic currents through the volume. Variations in conductivity affect both EEG and MEG spatial structure, although MEG is less sensitive to this parameter than EEG.

Charge conservation law is derived from Maxwell's equations in the following manner. Taking the divergence of Eq. (A.5) and using Eq. (A.2), we obtain

$$\nabla \cdot \mathbf{J} + \epsilon \frac{\partial \nabla \cdot \mathbf{E}}{\partial t} = \nabla \cdot \mathbf{J} + \frac{\partial \rho}{\partial t} = 0, \quad (\text{A.7})$$

where the relationship $\nabla \cdot (\nabla \times \mathbf{B}) = 0$ is used. Integrating over a closed volume V bounded by a surface S and using the Gauss theorem yield:

$$\oint_S \mathbf{J} \cdot \hat{\mathbf{n}} dS = -\frac{\partial}{\partial t} \int_V \rho dV, \quad (\text{A.8})$$

where $\hat{\mathbf{n}}$ is the outward unit normal to S . The integral on the left is the total current flowing outward across the S , and the integral on the right is the total charge in volume, so this relation states that the current flowing outward across S is equal to the rate of change of the charge in the volume.

A.2.1.1 Potential Formulation

It is convenient to re-express Maxwell's equations in terms of potential functions, which are related to the fields by simple derivatives. Because $\nabla \cdot \mathbf{B} = 0$, it is possible to write \mathbf{B} as a gradient of some other vector field \mathbf{A} , such that

$$\mathbf{B} = \nabla \times \mathbf{A}, \quad (\text{A.9})$$

where \mathbf{A} is called the magnetic vector potential. Substituting the equation above into Eq. (A.4), we obtain

$$\nabla \times \mathbf{E} + \frac{\partial}{\partial t} \nabla \times \mathbf{A} = \nabla \times \left(\mathbf{E} + \frac{\partial \mathbf{A}}{\partial t} \right) = 0. \quad (\text{A.10})$$

Using a scalar potential Φ , we can write

$$\mathbf{E} + \frac{\partial \mathbf{A}}{\partial t} = -\nabla \Phi, \quad (\text{A.11})$$

because for any scalar field Φ , the identity $\nabla \times (\nabla \Phi) = 0$ holds. Thus, we have

$$\mathbf{E} = -\nabla \Phi - \frac{\partial \mathbf{A}}{\partial t}. \quad (\text{A.12})$$

A.2.1.2 Gauge Transformations

Since electric and magnetic fields are defined as derivatives of the potentials, the potentials are not uniquely specified by their respective electromagnetic fields. This flexibility in defining \mathbf{A} and Φ is referred to as the choice of gauge. Considering again the relationship $\nabla \times \nabla \xi = 0$ for any scalar field ξ , the vector potential \mathbf{A} is only determined up to the gradient of a scalar function, i.e., we can make the replacement $\mathbf{A} \rightarrow \mathbf{A} - \nabla \xi$ without changing the value of the magnetic field \mathbf{B} . This implies that,

$$\begin{aligned} \mathbf{E} &\rightarrow -\nabla \Phi - \frac{\partial}{\partial t}(\mathbf{A} - \nabla \xi) \\ &= -\nabla(\Phi - \frac{\partial \xi}{\partial t}) - \frac{\partial \mathbf{A}}{\partial t}. \end{aligned} \quad (\text{A.13})$$

In other words, to keep \mathbf{E} from being unchanged, we must also make the replacement of

$$\Phi \rightarrow \Phi + \frac{\partial \xi}{\partial t}. \quad (\text{A.14})$$

This freedom to choose ξ allows us to impose an additional relationship between \mathbf{A} and Φ . Common choices include the “Coulomb” gauge where $\nabla \cdot \mathbf{A} = 0$ or the “Lorentz” gauge, where $\nabla \cdot \mathbf{A} + \mu\epsilon(\partial\Phi/\partial t) = 0$. In the following, we use the convenient “Gulrajani” gauge [4] for decoupling the differential equations for \mathbf{A} and Φ , thereby reflecting the fundamental symmetry between electric potentials and magnetic fields. The Gulrajani gauge is expressed as

$$\nabla \cdot \mathbf{A} + \mu\epsilon \frac{\partial \Phi}{\partial t} + \mu\sigma \Phi = 0. \quad (\text{A.15})$$

A.2.2 Magnetic Field and Electric Potential in an Unbounded Conducting Medium

A.2.2.1 Potential Equations

We have shown that using the vector and scalar potentials, the Maxwell's equations in Eqs. (A.3) and (A.4) are rewritten, respectively, as

$$\mathbf{B} = \nabla \times \mathbf{A}, \quad (\text{A.16})$$

$$\mathbf{E} = -\nabla \Phi - \frac{\partial \mathbf{A}}{\partial t}. \quad (\text{A.17})$$

Let us express the other two Maxwell's equations using these potentials and the source current \mathbf{J}_S . Let us substitute the two equations above into Eq. (A.5), resulting in

$$\nabla \times (\nabla \times \mathbf{A}) = \mu \mathbf{J} - \mu \epsilon \nabla \frac{\partial \Phi}{\partial t} - \mu \epsilon \frac{\partial^2 \mathbf{A}}{\partial t^2}. \quad (\text{A.18})$$

The relationship in Eq. (A.6) is expressed using the potentials such that

$$\mathbf{J} = \mathbf{J}_S - \sigma \nabla \Phi - \sigma \frac{\partial \mathbf{A}}{\partial t}. \quad (\text{A.19})$$

Substituting Eq. (A.19) into (A.18), and using the identity $\nabla \times (\nabla \times \mathbf{A}) = \nabla(\nabla \cdot \mathbf{A}) - \nabla^2 \mathbf{A}$, we obtain

$$\nabla^2 \mathbf{A} - \nabla \cdot \left[\nabla \cdot \mathbf{A} + \mu \epsilon \frac{\partial \Phi}{\partial t} + \mu \sigma \Phi \right] - \mu \epsilon \frac{\partial^2 \mathbf{A}}{\partial t^2} - \mu \sigma \frac{\partial \mathbf{A}}{\partial t} = -\mu \mathbf{J}_S, \quad (\text{A.20})$$

and finally

$$\nabla^2 \mathbf{A} - \mu \epsilon \frac{\partial^2 \mathbf{A}}{\partial t^2} - \mu \sigma \frac{\partial \mathbf{A}}{\partial t} = -\mu \mathbf{J}_S, \quad (\text{A.21})$$

where the gauge relationship in Eq. (A.15) is used in Eq. (A.20). From Maxwell's equation in Eq. (A.2), a similar derivation leads to

$$\nabla^2 \Phi - \mu \epsilon \frac{\partial^2 \Phi}{\partial t^2} - \mu \sigma \frac{\partial \Phi}{\partial t} = -\frac{\rho}{\epsilon}. \quad (\text{A.22})$$

Finally, we use Eqs. (A.6), (A.7), and (A.11), to obtain the time derivative of the charge ρ , expressed as

$$\begin{aligned}
\frac{\partial \rho}{\partial t} &= -\nabla \cdot \mathbf{J} \\
&= -\nabla \cdot (\sigma \mathbf{E} + \mathbf{J}_S) \\
&= -\sigma \nabla \cdot \left(-\nabla \Phi - \frac{\partial \mathbf{A}}{\partial t} \right) - \nabla \cdot \mathbf{J}_S \\
&= \sigma \nabla^2 \Phi + \sigma \frac{\partial \nabla \cdot \mathbf{A}}{\partial t} - \nabla \cdot \mathbf{J}_S \\
&= \sigma \nabla^2 \Phi - \sigma \frac{\partial}{\partial t} \left(\mu \epsilon \frac{\partial \Phi}{\partial t} + \mu \sigma \Phi \right) - \nabla \cdot \mathbf{J}_S \\
&= \sigma \nabla^2 \Phi - \epsilon \mu \sigma \frac{\partial^2 \Phi}{\partial t^2} - \mu \sigma^2 \frac{\partial \Phi}{\partial t} - \nabla \cdot \mathbf{J}_S,
\end{aligned} \tag{A.23}$$

where we again use the gauge relationship in Eq. (A.15).

A.2.2.2 Derivation of Potentials

To facilitate the solution of the equations for potentials, we can exploit linearity and also assume harmonic time dependence with angular frequency ω . This assumption does not require that the signals have perfect sinusoidal dependence, and is merely a mathematical convenience that helps subsequent discussion for determining the quasi-static regime. Namely, we express the solutions of the partial differential equations in Eqs. (A.21), (A.22), and (A.23), such that

$$\tilde{\mathbf{A}}(t, \mathbf{r}) = \mathbf{A}(\mathbf{r}) e^{i\omega t}, \tag{A.24}$$

$$\tilde{\Phi}(t, \mathbf{r}) = \Phi(\mathbf{r}) e^{i\omega t}, \tag{A.25}$$

$$\tilde{\rho}(t, \mathbf{r}) = \rho(\mathbf{r}) e^{i\omega t}, \tag{A.26}$$

$$\tilde{\mathbf{J}}_S(t, \mathbf{r}) = \mathbf{J}_S(\mathbf{r}) e^{i\omega t}. \tag{A.27}$$

Substituting $\tilde{\mathbf{A}}(t, \mathbf{r})$ and $\tilde{\mathbf{J}}_S(t, \mathbf{r})$ into Eq. (A.21), we obtain the following partial differential equation of space alone,

$$\nabla^2 \mathbf{A}(\mathbf{r}) + k^2 \mathbf{A}(\mathbf{r}) = -\mu \mathbf{J}_S(\mathbf{r}), \tag{A.28}$$

where, $k^2 = -i\omega\mu(\sigma + i\omega\epsilon)$.

Substitution of $\tilde{\Phi}(t, \mathbf{r})$ and $\tilde{\rho}(t, \mathbf{r})$ into Eq. (A.22) also results in

$$\nabla^2 \Phi(\mathbf{r}) + k^2 \Phi(\mathbf{r}) = -\frac{\rho(\mathbf{r})}{\epsilon}. \tag{A.29}$$

Substitution of $\tilde{\Phi}(t, \mathbf{r})$, $\tilde{\rho}(t, \mathbf{r})$, and $\tilde{\mathbf{J}}_S(t, \mathbf{r})$ into Eq. (A.23) gives the relationship:

$$i\omega\rho(\mathbf{r}) = \sigma \nabla^2 \Phi(\mathbf{r}) - \epsilon \mu \sigma (i\omega)^2 \Phi(\mathbf{r}) - \mu \sigma^2 (i\omega) \Phi(\mathbf{r}) - \nabla \cdot \mathbf{J}_S(\mathbf{r}),$$

which is simplified to

$$\rho(\mathbf{r}) = \frac{1}{i\omega} \sigma \nabla^2 \Phi(\mathbf{r}) - \left(\epsilon \mu \sigma(i\omega) + \mu \sigma^2 \right) \Phi(\mathbf{r}) - \frac{1}{i\omega} \nabla \cdot \mathbf{J}_S(\mathbf{r}). \quad (\text{A.30})$$

By substituting Eq. (A.30) into (A.29), we obtain

$$\begin{aligned} \nabla^2 \Phi(\mathbf{r}) - i\omega \mu(\sigma + i\omega\epsilon) \Phi(\mathbf{r}) \\ = -\frac{1}{i\omega\epsilon} \sigma \nabla^2 \Phi(\mathbf{r}) + \left(\mu\sigma(i\omega) + \frac{\mu\sigma^2}{\epsilon} \right) \Phi(\mathbf{r}) + \frac{1}{i\omega\epsilon} \nabla \cdot \mathbf{J}_S(\mathbf{r}). \end{aligned}$$

The equation above results in

$$\nabla^2 \Phi(\mathbf{r}) + k^2 \Phi(\mathbf{r}) = \frac{\nabla \cdot \mathbf{J}_S(\mathbf{r})}{\sigma + i\omega\epsilon}. \quad (\text{A.31})$$

In the absence of boundary conditions, Eqs. (A.28) and (A.31) have the following well-known solutions.

$$\mathbf{A}(\mathbf{r}) = \frac{\mu}{4\pi} \int_V \frac{\mathbf{J}_S(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} e^{-ik|\mathbf{r} - \mathbf{r}'|} d^3 r', \quad (\text{A.32})$$

and

$$\Phi(\mathbf{r}) = -\frac{1}{4\pi(\sigma + i\omega\epsilon)} \int_V \frac{\nabla' \cdot \mathbf{J}_S(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} e^{-ik|\mathbf{r} - \mathbf{r}'|} d^3 r'. \quad (\text{A.33})$$

Note that the derivation of these solutions has not involved any assumptions about the frequency ω or the tissue parameters σ, μ, ϵ . The typical values of these parameters, together with assumptions about the frequency ω , allow simplifications of the above equations suitable for EEG and MEG. The phase shifts of the electric potentials due to capacitive effects are accounted by the denominator of Eq. (A.33):

$$(\sigma + i\omega\epsilon) = \sigma \left(1 + i \frac{\omega\epsilon}{\sigma} \right). \quad (\text{A.34})$$

Thus capacitive effects may be ignored if $\omega\epsilon/\sigma \ll 1$. Although the very highest frequencies generated by the neural activity corresponding to action potentials may reach 1000 Hz, MEG and EEG signals are expected to be dominated by synaptic and dendritic activity that involves slower time scales less than 100 Hz ($\omega < 200\pi$ rad/s).

We make use of measured biological parameters to evaluate the dimensionless quantity $\omega\epsilon/\sigma$. To a best approximation, biological tissues are not magnetizable, thus $\mu = 4\pi \times 10^{-7}$ H/m, the value of vacuum. The dielectric permittivity ϵ is tissue and frequency dependent such that for brain and scalp, $\omega\epsilon/\sigma \simeq 0.01$ at 100 Hz, and $\omega\epsilon/\sigma \simeq 0.03$ at 1000 Hz. Thus, capacitive effects may be ignored to within 1% errors.

We then apply the Taylor series expansion:

$$e^{-ik|\mathbf{r}-\mathbf{r}'|} \simeq 1 - ik|\mathbf{r}-\mathbf{r}'| + \dots \quad (\text{A.35})$$

For MEG and EEG applications, we assume that $|\mathbf{r}-\mathbf{r}'| < 0.1$ m, typical for head-radius and distance from cortical sources to the sensors. With $\omega\epsilon/\sigma \ll 1$, we have $|k| \simeq \sqrt{\omega\mu\epsilon}$. Head tissue conductivity is highly variable, but the nominal value is $\sigma \approx 0.1$ s/m. Using this value leads to $|k||\mathbf{r}-\mathbf{r}'| \approx 0.0013$. Thus, propagation delays may also be ignored to within a 1 % error. With these simplifications, we can obtain

$$\mathbf{A}(\mathbf{r}) = \frac{\mu_0}{4\pi} \int_V \frac{\mathbf{J}_S(\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|} d^3 r', \quad (\text{A.36})$$

and

$$\boldsymbol{\Phi}(\mathbf{r}) = -\frac{1}{4\pi\sigma} \int_V \frac{\nabla' \cdot \mathbf{J}_S(\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|} d^3 r'. \quad (\text{A.37})$$

These simplified solutions are termed the quasi-static solutions, because the dependence on ω has been eliminated. This rigorous derivations often ignored in MEG and EEG literature, are given here to make explicit assumptions about tissue parameters and source frequencies which underly the quasi-static formulations of EEG and MEG.

A.2.2.3 Computation of Magnetic Field and Electric Potential

From the magnetic vector potential, the magnetic field is computed using:

$$\mathbf{B}(\mathbf{r}) = \nabla \times \mathbf{A}(\mathbf{r}) = \frac{\mu_0}{4\pi} \int_V \nabla \times \frac{\mathbf{J}_S(\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|} d^3 r'. \quad (\text{A.38})$$

We use the identity:

$$\begin{aligned} \nabla \times \frac{\mathbf{J}_S(\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|} &= \nabla \frac{1}{|\mathbf{r}-\mathbf{r}'|} \times \mathbf{J}_S(\mathbf{r}') + \frac{1}{|\mathbf{r}-\mathbf{r}'|} \nabla \times \mathbf{J}_S(\mathbf{r}') \\ &= \frac{\mathbf{J}_S(\mathbf{r}') \times (\mathbf{r}-\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|^3}. \end{aligned} \quad (\text{A.39})$$

Note that since ∇ is applied to \mathbf{r} , the relationships $\nabla|\mathbf{r}-\mathbf{r}'|^{-1} = -(\mathbf{r}-\mathbf{r}')/|\mathbf{r}-\mathbf{r}'|^{-3}$ and $\nabla \times \mathbf{J}_S(\mathbf{r}') = 0$ hold. Therefore, the magnetic field in an unbounded homogeneous medium is expressed as

$$\mathbf{B}(\mathbf{r}) = \frac{\mu_0}{4\pi} \int_V \frac{\mathbf{J}_S(\mathbf{r}') \times (\mathbf{r}-\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|^3} d^3 r'. \quad (\text{A.40})$$

This equation shows that the magnetic field in an unbounded homogeneous medium is expressed by the famous Biot-Savart law with replacing the total current \mathbf{J} with the source (impressed) current \mathbf{J}_S . In other words, the ohmic current \mathbf{J}_E does not contribute to the magnetic field in an unbounded homogeneous medium.

The expression for the scalar potential in Eq. (A.37) can be simplified in a following manner. Let us use the identity

$$\int_V \nabla' \cdot \frac{\mathbf{J}_S(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d^3 r' = \int_V \frac{\nabla' \cdot \mathbf{J}_S(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d^3 r' + \int_V \mathbf{J}_S(\mathbf{r}') \cdot \nabla' \frac{1}{|\mathbf{r} - \mathbf{r}'|} d^3 r'. \quad (\text{A.41})$$

The Gauss theorem is expressed as

$$\int_V \nabla' \cdot \frac{\mathbf{J}_S(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d^3 r' = \oint_S \frac{\mathbf{J}_S(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} \cdot \hat{\mathbf{n}} dS',$$

where $\hat{\mathbf{n}}$ is again the outward unit normal to the surface S . If we assume that $\mathbf{J}_S(\mathbf{r}')$ becomes zero on S , the surface integral on the right-hand side of the equation above becomes zero, and we can then get the relationship:

$$\int_V \frac{\nabla' \cdot \mathbf{J}_S(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d^3 r' = - \int_V \mathbf{J}_S(\mathbf{r}') \cdot \nabla' \frac{1}{|\mathbf{r} - \mathbf{r}'|} d^3 r'. \quad (\text{A.42})$$

Therefore, using

$$\nabla' \frac{1}{|\mathbf{r} - \mathbf{r}'|} = \frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3}, \quad (\text{A.43})$$

Eq. (A.37) is rewritten as

$$\Phi(\mathbf{r}) = \frac{1}{4\pi\sigma} \int_V \mathbf{J}_S(\mathbf{r}') \cdot \frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3} d^3 r'. \quad (\text{A.44})$$

The above equation is the expression to compute the electric potential in an unbounded homogeneous medium.

A.2.2.4 Dipoles in an Unbounded Homogeneous Medium

The transmembrane current density \mathbf{J}_S arises due to concentration gradients. Source models for \mathbf{J}_S fall into two categories: rigorous and phenomenological. Although a rigorous source model accounts reasonably accurately for each of the microscopic currents, it is difficult to derive such source models. A phenomenological source model is one which produces the same external fields, but is artificial in the sense that it does not actually reflect the microscopic details of the problem.

A representative phenomenological source model is the dipole model. The dipole is the simplest source for both Φ and \mathbf{B} and can be written as:

$$\mathbf{J}_S = \mathbf{Q}(\mathbf{r}_0) = \mathbf{Q}\delta(\mathbf{r} - \mathbf{r}_0). \quad (\text{A.45})$$

Substituting this for the equations for electric potential in Eq. (A.44) and magnetic field in Eq. (A.40) in a homogenous infinite volume conductor, we get

$$\boldsymbol{\Phi}(\mathbf{r}) = \frac{1}{4\pi\sigma} \mathbf{Q} \cdot \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|^3}, \quad (\text{A.46})$$

and

$$\mathbf{B}(\mathbf{r}) = \frac{\mu_0}{4\pi} \mathbf{Q} \times \frac{\mathbf{r} - \mathbf{r}_0}{|\mathbf{r} - \mathbf{r}_0|^3}. \quad (\text{A.47})$$

The equations above are the expressions for computing the electric potential and magnetic field produced by an current dipole at \mathbf{r}_0 embedded in an unbounded homogeneous conductive medium.

A.2.3 Magnetic Field from a Bounded Conductor with Piecewise-Constant Conductivity

We next consider the magnetic field generated by an inhomogeneous conductor. We assume that the region V can be divided into subregions $V_j, j = 1, \dots, \varpi$, and the region V_j has conductivity σ_j . The surface of V_j is denoted S_j . We also assume that the conductivity, $\sigma(\mathbf{r})$, is zero outside V . We start the derivation from the Bio-Savart law:

$$\mathbf{B}(\mathbf{r}) = \frac{\mu_0}{4\pi} \int_V \mathbf{J}(\mathbf{r}') \times \frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3} d^3 r' = \frac{\mu_0}{4\pi} \int_V \mathbf{J}(\mathbf{r}') \times \mathbf{G}(\mathbf{r}, \mathbf{r}') d^3 r', \quad (\text{A.48})$$

where

$$\mathbf{G}(\mathbf{r}, \mathbf{r}') = \frac{\mathbf{r} - \mathbf{r}'}{|\mathbf{r} - \mathbf{r}'|^3}.$$

Substituting $\mathbf{J}(\mathbf{r}') = \mathbf{J}_S(\mathbf{r}') - \sigma(\mathbf{r}') \nabla \boldsymbol{\Phi}(\mathbf{r}')$ into Eq. (A.48), we can obtain

$$\begin{aligned} \mathbf{B}(\mathbf{r}) &= \frac{\mu_0}{4\pi} \int_V [\mathbf{J}_S(\mathbf{r}') - \sigma(\mathbf{r}') \nabla \boldsymbol{\Phi}(\mathbf{r}')] \times \mathbf{G}(\mathbf{r}, \mathbf{r}') d^3 r' \\ &= \frac{\mu_0}{4\pi} \int_V \mathbf{J}_S(\mathbf{r}') \times \mathbf{G}(\mathbf{r}, \mathbf{r}') d^3 r' - \frac{\mu_0}{4\pi} \sum_{j=1}^{\varpi} \sigma_j \int_{V_j} \nabla \boldsymbol{\Phi}(\mathbf{r}') \times \mathbf{G}(\mathbf{r}, \mathbf{r}') d^3 r'. \end{aligned} \quad (\text{A.49})$$

Using

$$\nabla \boldsymbol{\Phi}(\mathbf{r}') \times \mathbf{G}(\mathbf{r}, \mathbf{r}') = \nabla \times [\boldsymbol{\Phi}(\mathbf{r}') \mathbf{G}(\mathbf{r}, \mathbf{r}')], \quad (\text{A.50})$$

the second term on the right-hand side of Eq. (A.49) can be rewritten as

$$\begin{aligned} \int_{V_j} \nabla \Phi(\mathbf{r}') \times \mathbf{G}(\mathbf{r}, \mathbf{r}') d^3 r' &= \int_{V_j} \nabla \times [\Phi(\mathbf{r}') \mathbf{G}(\mathbf{r}, \mathbf{r}')] d^3 r' \\ &= \int_{S_j} \hat{\mathbf{n}}(\mathbf{r}') dS \times [\Phi(\mathbf{r}') \mathbf{G}(\mathbf{r}, \mathbf{r}')] \\ &= \int_{S_j} \Phi(\mathbf{r}') \hat{\mathbf{n}}(\mathbf{r}') \times \mathbf{G}(\mathbf{r}, \mathbf{r}') dS, \end{aligned} \quad (\text{A.51})$$

where S_j indicates the surface of V_j , and $\hat{\mathbf{n}}(\mathbf{r}')$ is the outward unit normal vector of a surface element on S_j . We use the Gauss theorem to derive the right-hand side of Eq. (A.51). Substituting Eq. (A.51) into (A.49), we can derive the following Geselowitz formula [5]:

$$\mathbf{B}(\mathbf{r}) = \mathbf{B}_0(\mathbf{r}) - \frac{\mu_0}{4\pi} \sum_{j=1}^{\infty} (\sigma_j - \sigma'_j) \int_{S_j} \Phi(\mathbf{r}') \hat{\mathbf{n}}(\mathbf{r}') \times \mathbf{G}(\mathbf{r}, \mathbf{r}') dS, \quad (\text{A.52})$$

where σ'_j is the conductivity just outside V_j . Also, in Eq. (A.52), $\mathbf{B}_0(\mathbf{r})$ is the magnetic field for the infinite homogeneous conductor in Eq. (A.40).

A.2.4 Magnetic Field from a Homogeneous Spherical Conductor

Here, we assume that V_j is spherically symmetric, and we set the coordinate origin at the center of V_j . We have

$$\mathbf{B}(\mathbf{r}) \cdot \mathbf{e}_r = \mathbf{B}_0(\mathbf{r}) \cdot \mathbf{e}_r - \frac{\mu_0}{4\pi} \sum_{j=1}^{\infty} (\sigma_j - \sigma'_j) \int_{S_j} \Phi(\mathbf{r}') \hat{\mathbf{n}}(\mathbf{r}') \times \mathbf{G}(\mathbf{r}, \mathbf{r}') \cdot \mathbf{e}_r dS, \quad (\text{A.53})$$

where \mathbf{e}_r is the unit vector in the radial direction, which is defined as $\mathbf{e}_r = \mathbf{r}/|\mathbf{r}|$. Since $\hat{\mathbf{n}}(\mathbf{r}) = \mathbf{e}_r$ in the case of a spherically symmetric conductor, the second term is equal to zero, and we have the relationship

$$\mathbf{B}(\mathbf{r}) \cdot \mathbf{e}_r = \mathbf{B}_0(\mathbf{r}) \cdot \mathbf{e}_r. \quad (\text{A.54})$$

This equation indicates that the radial component of the magnetic field is not affected by the volume current and that the radial component is determined solely by the primary current.

We next derive a closed-form formula for the magnetic field outside a spherically-symmetric homogeneous conductor. The derivation is according to Sarvas [6]. The

relationship $\nabla \times \mathbf{B}(\mathbf{r}) = 0$ holds outside the volume conductor, because there is no electric current. Thus, $\mathbf{B}(\mathbf{r})$ can be expressed in terms of the magnetic scalar potential $U(\mathbf{r})$, as

$$\mathbf{B}(\mathbf{r}) = -\mu_0 \nabla U(\mathbf{r}). \quad (\text{A.55})$$

This potential function is derived from

$$U(\mathbf{r}) = \frac{1}{\mu_0} \int_0^\infty \mathbf{B}(\mathbf{r} + \tau \mathbf{e}_r) \cdot \mathbf{e}_r d\tau = \frac{1}{\mu_0} \int_0^\infty \mathbf{B}_0(\mathbf{r} + \tau \mathbf{e}_r) \cdot \mathbf{e}_r d\tau, \quad (\text{A.56})$$

where we use the relationship in Eq. (A.54). Assuming that a dipole source exists at \mathbf{r}_0 , by substituting Eq. (A.47) into Eq. (A.56) and performing the integral, we obtain

$$\begin{aligned} U(\mathbf{r}) &= \frac{1}{\mu_0} \int_0^\infty \mathbf{B}_0(\mathbf{r} + \tau \mathbf{e}_r) \cdot \mathbf{e}_r d\tau \\ &= \frac{1}{4\pi} \mathbf{Q} \times (\mathbf{r} - \mathbf{r}_0) \cdot \mathbf{e}_r \int_0^\infty \frac{d\tau}{|\mathbf{r} + \tau \mathbf{e}_r - \mathbf{r}_0|^3} \\ &= -\frac{1}{4\pi} \frac{(\mathbf{Q} \times \mathbf{r}_0) \cdot \mathbf{r}}{\Lambda}, \end{aligned} \quad (\text{A.57})$$

where

$$\Lambda = |\mathbf{r} - \mathbf{r}_0|(|\mathbf{r} - \mathbf{r}_0| |\mathbf{r}|^2 - \mathbf{r}_0 \cdot \mathbf{r}). \quad (\text{A.58})$$

The well-known Sarvas formula [6] for $\mathbf{B}(\mathbf{r})$ is then obtained by substituting Eq. (A.57) into Eq. (A.55) and performing the gradient operation. The results are expressed as

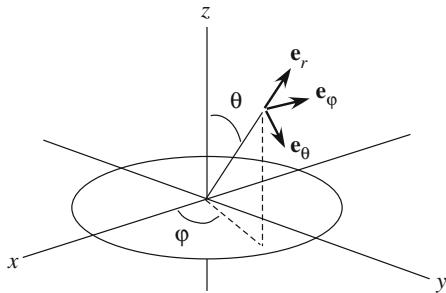
$$\begin{aligned} \mathbf{B}(\mathbf{r}) &= \frac{\mu_0}{4\pi} \nabla \frac{(\mathbf{Q} \times \mathbf{r}_0) \cdot \mathbf{r}}{\Lambda} = \frac{\mu_0}{4\pi} \left[\frac{\mathbf{Q} \times \mathbf{r}_0}{\Lambda} - \frac{1}{\Lambda^2} (\mathbf{Q} \times \mathbf{r}_0) \cdot \mathbf{r} \nabla \Lambda \right] \\ &= \frac{\mu_0}{4\pi \Lambda^2} [\Lambda \mathbf{Q} \times \mathbf{r}_0 - [(\mathbf{Q} \times \mathbf{r}_0) \cdot \mathbf{r}] \nabla \Lambda], \end{aligned} \quad (\text{A.59})$$

where

$$\begin{aligned} \nabla \Lambda &= \left[\frac{|\mathbf{r} - \mathbf{r}_0|^2}{|\mathbf{r}|} + \frac{(\mathbf{r} - \mathbf{r}_0) \cdot \mathbf{r}}{|\mathbf{r} - \mathbf{r}_0|} + 2|\mathbf{r} - \mathbf{r}_0| + 2|\mathbf{r}| \right] \mathbf{r} \\ &\quad - \left[|\mathbf{r} - \mathbf{r}_0| + 2|\mathbf{r}| + \frac{(\mathbf{r} - \mathbf{r}_0) \cdot \mathbf{r}}{|\mathbf{r} - \mathbf{r}_0|} \right] \mathbf{r}_0. \end{aligned} \quad (\text{A.60})$$

We can see that when \mathbf{r}_0 approaches the center of the sphere, $\mathbf{B}(\mathbf{r})$ becomes zero, and no magnetic field is generated outside the conductor from a source at the origin. Also, if the source vector \mathbf{Q} and the location vector \mathbf{r}_0 are parallel, i.e., if the primary current source is oriented in the radial direction, no magnetic fields are generated outside the spherical conductor from such a radial source. This is because the two

Fig. A.1 The three orthogonal directions (e_r, e_ϕ, e_θ) used to express the source vector when the spherically-symmetric conductor model is used for the forward calculation



terms on the left-hand side of Eq. (A.59) contain the vector product $\mathbf{Q} \times \mathbf{r}_0$, which is equal to zero when \mathbf{Q} and \mathbf{r}_0 are parallel.

Therefore, when using the spherically-homogeneous conductor model, instead of the x, y, z directions, we usually use the three orthogonal directions (e_r, e_ϕ, e_θ) to express the source vector. These directions are illustrated in Fig. A.1. Because the e_r component of a source never creates a measurable magnetic field outside the spherical conductor, we can disregard this component and only deal with the e_ϕ and e_θ components of the source vector.

References

1. R. Plonsey, *Bioelectric phenomena*. Wiley Online Library, 1969.
2. H. Tuckwell, “Introduction to Theoretical Neurobiology vols. 1 and 2, 1988,” *Cambridge University Press, Cambridge, Longtin A., Bulsara A., Moss F., Phys. Rev. Lett*, vol. 65, p. 656, 1991.
3. B. Hille, *Ion channels of excitable membranes*, vol. 507. Sinauer Sunderland, MA, 2001.
4. R. M. Gulrajani, *Bioelectricity and biomagnetism*. John Wiley and Sons, 1998.
5. D. B. Geselowitz, “On the magnetic field generated outside an inhomogeneous volume conductor by internal current sources,” *IEEE Trans. Biomed. Eng.*, vol. 2, pp. 346–347, 1970.
6. J. Sarvas, “Basic mathematical and electromagnetic concepts of the biomagnetic inverse problem,” *Phys. Med. Biol.*, vol. 32, pp. 11–22, 1987.

Appendix B

Basics of Bayesian Inference

B.1 Linear Model and Bayesian Inference

This appendix explains the basics of Bayesian inference. Here, we consider the general problem of estimating the vector \mathbf{x} , containing N unknowns: $\mathbf{x} = [x_1, \dots, x_N]^T$, from the observation (the sensor data) \mathbf{y} , containing M elements: $\mathbf{y} = [y_1, \dots, y_M]^T$. We assume a linear model between the observation \mathbf{y} and the unknown vector \mathbf{x} , such that

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\varepsilon}, \quad (\text{B.1})$$

where \mathbf{H} is an $M \times N$ matrix that expresses the forward relationship between \mathbf{y} and \mathbf{x} , and $\boldsymbol{\varepsilon}$ is an additive noise overlapped to the observation \mathbf{y} .

To solve this estimation problem based on the Bayesian inference, we consider \mathbf{x} a vector random variable, and use the following three kinds of probability distributions.

- (1) $p(\mathbf{x})$: The probability distribution on the unknown \mathbf{x} . This is called the prior probability distribution. It represents our prior knowledge on the unknown \mathbf{x} .
- (2) $p(\mathbf{y}|\mathbf{x})$: The conditional probability of \mathbf{y} given \mathbf{x} . This conditional probability is equal to the likelihood. The maximum likelihood method, mentioned in Chap. 2, estimates the unknown \mathbf{x} as the value of \mathbf{x} that maximizes $p(\mathbf{y}|\mathbf{x})$.
- (3) $p(\mathbf{x}|\mathbf{y})$: The probability of \mathbf{x} given observation \mathbf{y} . This is called the posterior probability. The Bayesian inference estimates the unknown parameter \mathbf{x} based on this posterior probability.

The posterior probability is obtained from the prior probability $p(\mathbf{x})$ and the likelihood $p(\mathbf{y}|\mathbf{x})$ using Bayes' rule,

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{\int p(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x}}. \quad (\text{B.2})$$

On the right-hand side of Eq. (B.2), the denominator is used only for the normalization $\int p(\mathbf{x}|\mathbf{y})d\mathbf{x} = 1$, and often the denominator is not needed to estimate the posterior $p(\mathbf{x}|\mathbf{y})$. Therefore, Bayes' rule can be expressed in a simpler form such as

$$p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x})p(\mathbf{x}), \quad (\text{B.3})$$

where the notation \propto means that both sides are equal, ignoring a multiplicative constant.

Bayesian inference uses Bayes' rule in Eq. (B.2) or (B.3) to estimate \mathbf{x} . A problem here is how to determine the prior probability distribution $p(\mathbf{x})$. A general strategy for this problem is to determine the prior probability by taking what we know about the unknown parameters \mathbf{x} into account. However, if there is no prior knowledge on \mathbf{x} , we must use the uniform prior distribution, i.e.,

$$p(\mathbf{x}) = \text{constant}. \quad (\text{B.4})$$

With Eq. (B.4), Bayes' rule in Eq. (B.3) becomes

$$p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x}). \quad (\text{B.5})$$

In this case, the posterior probability $p(\mathbf{x}|\mathbf{y})$ is equal to the likelihood $p(\mathbf{y}|\mathbf{x})$, resulting in a situation that Bayesian and maximum likelihood methods give the same solution. The prior distribution in Eq. (B.4) is referred to as the non-informative prior.

Even when some prior information on the unknown parameter \mathbf{x} is available, exact probability distributions are generally difficult to determine. Therefore, the probability distribution is usually determined according to the convenience in computing the posterior distribution. Some probability distributions have the same forms as the prior and posterior distributions for given $p(\mathbf{y}|\mathbf{x})$. One representative example of such distributions is the Gaussian distribution. That is, if the noise $\boldsymbol{\varepsilon}$ is Gaussian, a Gaussian prior distribution gives a Gaussian posterior distribution. The derivation of the posterior distribution in the Gaussian model is explained in Sect. B.3.

B.2 Point Estimate of Unknown \mathbf{x}

In Bayesian inference, the unknown parameter \mathbf{x} is estimated based on the posterior probability $p(\mathbf{x}|\mathbf{y})$. Then, how can we obtain the optimum estimate $\hat{\mathbf{x}}$ based on $p(\mathbf{x}|\mathbf{y})$? There are two ways to compute the estimate $\hat{\mathbf{x}}$ based on a given posterior distribution. One way chooses the \mathbf{x} that maximizes the posterior, i.e.,

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmax}} p(\mathbf{x}|\mathbf{y}). \quad (\text{B.6})$$

This $\hat{\mathbf{x}}$ is called the maximum a posteriori (MAP) estimate.

The other way is to choose $\hat{\mathbf{x}}$ that minimizes the squared error between the estimate $\hat{\mathbf{x}}$ and the true value \mathbf{x} . The squared error is expressed as $E[(\hat{\mathbf{x}} - \mathbf{x})^T(\hat{\mathbf{x}} - \mathbf{x})]$, and the estimate $\hat{\mathbf{x}}$ is obtained using

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\hat{\mathbf{x}}} E \left[(\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) \right]. \quad (\text{B.7})$$

This estimate is called the minimum mean squared error (MMSE) estimate. Here, taking the relationship,

$$\begin{aligned} E \left[(\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) \right] &= \iint_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x} \right] p(\mathbf{y}) d\mathbf{y}, \end{aligned} \quad (\text{B.8})$$

and the fact that $p(\mathbf{y}) \geq 0$ into consideration, the $\hat{\mathbf{x}}$ that minimizes $E[(\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x})]$ is equal to

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\hat{\mathbf{x}}} \int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x}. \quad (\text{B.9})$$

Taking the derivative of the above integral with respect to $\hat{\mathbf{x}}$, and setting it to zero, we have

$$\frac{\partial}{\partial \hat{\mathbf{x}}} \int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x} = 2 \int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x} = 0. \quad (\text{B.10})$$

Therefore, the MMSE estimate is equal to

$$\hat{\mathbf{x}} = \int_{-\infty}^{\infty} \mathbf{x} p(\mathbf{x}|\mathbf{y}) d\mathbf{x}. \quad (\text{B.11})$$

That is, the MMSE estimate is equal to the mean of the posterior. Note that, when the posterior distribution is Gaussian, the MAP estimate and the MMSE estimate are equal, because the Gaussian distribution reaches its maximum at the mean.

B.3 Derivation of Posterior Distribution in the Gaussian Model

Let us consider the problem in which the time series of unknown parameters $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$ are estimated using the time series of observation data $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K$ where $\mathbf{y}(t_k)$ and $\mathbf{x}(t_k)$ are denoted \mathbf{y}_k and \mathbf{x}_k . The relationship in Eq. (B.1) holds between \mathbf{x}_k and \mathbf{y}_k , i.e.,

$$\mathbf{y}_k = \mathbf{H}\mathbf{x}_k + \boldsymbol{\varepsilon}. \quad (\text{B.12})$$

In this chapter, the whole time series data $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K$ is collectively denoted \mathbf{y} , and the whole time series data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$ are collectively denoted \mathbf{x} . We assume that the noise $\boldsymbol{\varepsilon}$ is Gaussian and is identically and independently distributed across time, i.e.,

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(\boldsymbol{\varepsilon} | \mathbf{0}, \boldsymbol{\Lambda}^{-1}), \quad (\text{B.13})$$

where we omit the notation of the time index k from $\boldsymbol{\varepsilon}$. In Eq. (B.13), $\boldsymbol{\Lambda}$ is a diagonal precision matrix of which the j th diagonal entry is equal to the noise precision for the j th observation data. Then, using Eqs. (B.13) and (C.3), the conditional probability $p(\mathbf{y}_k | \mathbf{x}_k)$ is obtained as

$$p(\mathbf{y}_k | \mathbf{x}_k) = \mathcal{N}(\mathbf{y}_k | \mathbf{H}\mathbf{x}_k, \boldsymbol{\Lambda}^{-1}). \quad (\text{B.14})$$

The conditional probability of the whole time series of \mathbf{y}_k given the whole time series of \mathbf{x}_k is given by

$$\begin{aligned} p(\mathbf{y} | \mathbf{x}) &= p(\mathbf{y}_1, \dots, \mathbf{y}_K | \mathbf{x}_1, \dots, \mathbf{x}_K) \\ &= \prod_{k=1}^K p(\mathbf{y}_k | \mathbf{x}_k) = \prod_{k=1}^K \mathcal{N}(\mathbf{y}_k | \mathbf{H}\mathbf{x}_k, \boldsymbol{\Lambda}^{-1}). \end{aligned} \quad (\text{B.15})$$

The prior distribution of \mathbf{x}_k is assumed to be Gaussian and independent across time:

$$p(\mathbf{x}_k) = \mathcal{N}(\mathbf{x}_k | \mathbf{0}, \boldsymbol{\Phi}^{-1}). \quad (\text{B.16})$$

The prior distribution for the whole time series of \mathbf{x}_k is expressed as

$$p(\mathbf{x}) = p(\mathbf{x}_1, \dots, \mathbf{x}_K) = \prod_{k=1}^K p(\mathbf{x}_k) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}_k | \mathbf{0}, \boldsymbol{\Phi}^{-1}). \quad (\text{B.17})$$

In this case, the posterior probability is independent across time, and given by

$$p(\mathbf{x} | \mathbf{y}) = p(\mathbf{x}_1, \dots, \mathbf{x}_K | \mathbf{y}_1, \dots, \mathbf{y}_K) = \prod_{k=1}^K p(\mathbf{x}_k | \mathbf{y}_k). \quad (\text{B.18})$$

The posterior probability $p(\mathbf{x}_k | \mathbf{y}_k)$ can be derived by substituting Eqs. (B.16) and (B.14) into Bayes' rule:

$$p(\mathbf{x}_k | \mathbf{y}_k) \propto p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k). \quad (\text{B.19})$$

Actual computation of $p(\mathbf{x}_k | \mathbf{y}_k)$ is performed in the following manner. Since we know that the posterior distribution is also Gaussian, the posterior distribution is assumed to be

$$p(\mathbf{x}_k | \mathbf{y}_k) = \mathcal{N}(\mathbf{x}_k | \bar{\mathbf{x}}_k, \boldsymbol{\Gamma}^{-1}), \quad (\text{B.20})$$

where $\bar{\mathbf{x}}$ is the posterior mean, and $\boldsymbol{\Gamma}$ is the posterior precision matrix. The exponential part of the above Gaussian distribution is given by

$$-\frac{1}{2}(\mathbf{x}_k - \bar{\mathbf{x}}_k)^T \boldsymbol{\Gamma} (\mathbf{x}_k - \bar{\mathbf{x}}_k) = -\frac{1}{2}\mathbf{x}_k^T \boldsymbol{\Gamma} \mathbf{x}_k + \mathbf{x}_k^T \boldsymbol{\Gamma} \bar{\mathbf{x}}_k + \mathcal{C}, \quad (\text{B.21})$$

where \mathcal{C} represents terms that do not contain \mathbf{x}_k . The exponential part of the right-hand side of Eq. (B.19) is given by

$$-\frac{1}{2} \left[\mathbf{x}_k^T \boldsymbol{\Phi} \mathbf{x}_k + (\mathbf{y}_k - \mathbf{H} \mathbf{x}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{H} \mathbf{x}_k) \right]. \quad (\text{B.22})$$

The above expression can be rewritten as

$$-\frac{1}{2} \mathbf{x}_k^T (\boldsymbol{\Phi} + \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{H}) \mathbf{x}_k + \mathbf{x}_k^T \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{y}_k + \mathcal{C}', \quad (\text{B.23})$$

where \mathcal{C}' again represents terms that do not contain \mathbf{x}_k . Comparing the quadratic and linear terms of \mathbf{x}_k between Eqs. (B.21) and (B.23) gives the relationships

$$\boldsymbol{\Gamma} = \boldsymbol{\Phi} + \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{H}, \quad (\text{B.24})$$

$$\bar{\mathbf{x}}_k = \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{y}_k = (\boldsymbol{\Phi} + \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{H})^{-1} \mathbf{H}^T \boldsymbol{\Lambda} \mathbf{y}_k. \quad (\text{B.25})$$

The precision matrix and the mean of the posterior distribution are obtained in Eqs. (B.24) and (B.25), respectively. This $\bar{\mathbf{x}}_k$ is the MMSE (and also MAP) estimate of \mathbf{x}_k . Using the matrix inversion formula in Eq. (C.92), $\bar{\mathbf{x}}_k$ can also be expressed as

$$\bar{\mathbf{x}}_k = \boldsymbol{\Phi}^{-1} \mathbf{H}^T (\mathbf{H} \boldsymbol{\Phi}^{-1} \mathbf{H}^T + \boldsymbol{\Lambda}^{-1})^{-1} \mathbf{y}_k = \boldsymbol{\Upsilon} \mathbf{H}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k, \quad (\text{B.26})$$

where $\boldsymbol{\Upsilon}$, which is equal to $\boldsymbol{\Phi}^{-1}$, is the covariance matrix of the prior distribution, and $\boldsymbol{\Sigma}_y$ is expressed in Eq. (B.30).

B.4 Derivation of the Marginal Distribution $p(\mathbf{y})$

The probability for the observation \mathbf{y} is obtained using

$$p(\mathbf{y}) = \prod_{k=1}^K p(\mathbf{y}_k) = \prod_{k=1}^K \int p(\mathbf{x}_k, \mathbf{y}_k) d\mathbf{x}_k = \prod_{k=1}^K \int p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k) d\mathbf{x}_k. \quad (\text{B.27})$$

Substituting Eqs. (B.16) and (B.14) into the equation above, we get

$$p(\mathbf{y}) = \prod_{k=1}^K \int \mathcal{N}(\mathbf{y}_k | \mathbf{H}\mathbf{x}_k, \boldsymbol{\Lambda}^{-1}) \mathcal{N}(\mathbf{x}_k | \mathbf{0}, \boldsymbol{\Phi}^{-1}) d\mathbf{x}_k. \quad (\text{B.28})$$

This $p(\mathbf{y})$ is called the marginal likelihood. This integral is computed in Sect. 4.3. Using Eq. (4.30), we have

$$p(\mathbf{y}) \propto \prod_{k=1}^K \frac{1}{|\boldsymbol{\Sigma}_y|^{1/2}} \exp \left[-\frac{1}{2} \mathbf{y}_k^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_k \right], \quad (\text{B.29})$$

where $\boldsymbol{\Sigma}_y$ is expressed as

$$\boldsymbol{\Sigma}_y = \boldsymbol{\Lambda}^{-1} + \mathbf{H}\boldsymbol{\Phi}^{-1}\mathbf{H}^T. \quad (\text{B.30})$$

This $\boldsymbol{\Sigma}_y$ is referred to as the model data covariance. Taking the normalization into account, we get the marginal distribution $p(\mathbf{y})$, such that

$$p(\mathbf{y}) = \prod_{k=1}^K p(\mathbf{y}_k) \quad \text{where} \quad p(\mathbf{y}_k) = \mathcal{N}(\mathbf{y}_k | \mathbf{0}, \boldsymbol{\Sigma}_y). \quad (\text{B.31})$$

Thus, $\boldsymbol{\Sigma}_y$, the model data covariance, is the covariance matrix of the marginal distribution $p(\mathbf{y}_k)$.

B.5 Expectation Maximization (EM) Algorithm

B.5.1 Marginal Likelihood Maximization

The Bayesian estimate of the unknown \mathbf{x} , $\bar{\mathbf{x}}_k$, is computed using Eq. (B.25). However, when computing $\bar{\mathbf{x}}_k$, $\boldsymbol{\Phi}$ and $\boldsymbol{\Lambda}$ are needed. The precision matrix of the prior distribution $\boldsymbol{\Phi}$ and that of the noise $\boldsymbol{\Lambda}$ are called the hyperparameters. Quite often, the hyperparameters are unknown, and should also be estimated from the observation \mathbf{y} . The hyperparameters may be estimated by maximizing $p(\mathbf{y}|\boldsymbol{\Phi}, \boldsymbol{\Lambda})$, which is the likelihood with respect to $\boldsymbol{\Phi}$ and $\boldsymbol{\Lambda}$. This $p(\mathbf{y}|\boldsymbol{\Phi}, \boldsymbol{\Lambda})$ is called the marginal likelihood or the data evidence to discriminate it from the conventional likelihood $p(\mathbf{y}|\mathbf{x})$.

The marginal likelihood can be computed with $p(\mathbf{x}|\boldsymbol{\Phi})$ and $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\Lambda})$ using¹

$$p(\mathbf{y}|\boldsymbol{\Phi}, \boldsymbol{\Lambda}) = \int_{-\infty}^{\infty} p(\mathbf{y}|\mathbf{x}, \boldsymbol{\Lambda}) p(\mathbf{x}|\boldsymbol{\Phi}) d\mathbf{x}. \quad (\text{B.32})$$

However, computation of $\boldsymbol{\Phi}$ and $\boldsymbol{\Lambda}$ that maximize $p(\mathbf{y}|\boldsymbol{\Phi}, \boldsymbol{\Lambda})$ is generally quite troublesome, as is demonstrated in detail in Chap. 4. In the following, instead of directly

¹ The notation $d\mathbf{x}$ indicates $d\mathbf{x}_1 d\mathbf{x}_2 \cdots d\mathbf{x}_K$.

maximizing the marginal likelihood, we describe an algorithm that maximizes a quantity called the average data likelihood to obtain estimates for the hyperparameters. This algorithm is called the expectation maximization (EM) algorithm, and is described in the following.

B.5.2 Average Data Likelihood

The EM algorithm computes the quantity called the average data likelihood. Computing the average data likelihood is much easier than computing the marginal likelihood. To define the average data likelihood, let us first define the complete data likelihood, such that

$$\log p(\mathbf{y}, \mathbf{x} | \boldsymbol{\Phi}, \boldsymbol{\Lambda}) = \log p(\mathbf{y} | \mathbf{x}, \boldsymbol{\Lambda}) + \log p(\mathbf{x} | \boldsymbol{\Phi}). \quad (\text{B.33})$$

If we observed not only \mathbf{y} but also \mathbf{x} , we could have estimated $\boldsymbol{\Phi}$ and $\boldsymbol{\Lambda}$ by maximizing $\log p(\mathbf{y}, \mathbf{x} | \boldsymbol{\Phi}, \boldsymbol{\Lambda})$ with respect to these hyperparameters. However, since we do not observe \mathbf{x} , we must substitute for the unknown \mathbf{x} in $\log p(\mathbf{y}, \mathbf{x} | \boldsymbol{\Phi}, \boldsymbol{\Lambda})$ with some “reasonable” value.

Having observed \mathbf{y} , we actually know which values of \mathbf{x} are reasonable, and our best knowledge on the unknown \mathbf{x} is represented by the posterior distribution $p(\mathbf{x} | \mathbf{y})$. Thus, the “reasonable” value would be the one that maximizes the posterior probability, and one solution would be to use the MAP estimate of \mathbf{x} in $\log p(\mathbf{y}, \mathbf{x} | \boldsymbol{\Phi}, \boldsymbol{\Lambda})$. A better solution would be to use all possible values of \mathbf{x} in the complete data likelihood and average over it with the posterior probability. This results in the average data likelihood, $\Theta(\boldsymbol{\Phi}, \boldsymbol{\Lambda})$:

$$\begin{aligned} \Theta(\boldsymbol{\Phi}, \boldsymbol{\Lambda}) &= \int p(\mathbf{x} | \mathbf{y}) \log p(\mathbf{y}, \mathbf{x} | \boldsymbol{\Phi}, \boldsymbol{\Lambda}) d\mathbf{x} \\ &= E[\log p(\mathbf{y}, \mathbf{x} | \boldsymbol{\Phi}, \boldsymbol{\Lambda})] \\ &= E[\log p(\mathbf{y} | \mathbf{x}, \boldsymbol{\Lambda})] + E[\log p(\mathbf{x} | \boldsymbol{\Phi})], \end{aligned} \quad (\text{B.34})$$

where the expectation $E[\cdot]$ is taken with respect to the posterior probability $p(\mathbf{x} | \mathbf{y})$. The estimates of the hyperparameters, $\boldsymbol{\Phi}$ and $\boldsymbol{\Lambda}$ are obtained using

$$\hat{\boldsymbol{\Lambda}} = \underset{\boldsymbol{\Lambda}}{\operatorname{argmax}} \Theta(\boldsymbol{\Phi}, \boldsymbol{\Lambda}), \quad (\text{B.35})$$

$$\hat{\boldsymbol{\Phi}} = \underset{\boldsymbol{\Phi}}{\operatorname{argmax}} \Theta(\boldsymbol{\Phi}, \boldsymbol{\Lambda}). \quad (\text{B.36})$$

In the Gaussian model discussed in Sect. B.3, $p(\mathbf{x} | \boldsymbol{\Phi})$ and $p(\mathbf{y} | \mathbf{x}, \boldsymbol{\Lambda})$ are expressed in Eqs. (B.17) and (B.15), respectively. Substituting Eqs. (B.17) and (B.15) into (B.33), the complete data likelihood is expressed as²

² The constant terms containing 2π are ignored here.

$$\begin{aligned}\log p(\mathbf{y}, \mathbf{x} | \boldsymbol{\Phi}, \boldsymbol{\Lambda}) &= \frac{K}{2} \log |\boldsymbol{\Phi}| - \frac{1}{2} \sum_{k=1}^K \mathbf{x}_k^T \boldsymbol{\Phi} \mathbf{x}_k \\ &\quad + \frac{K}{2} \log |\boldsymbol{\Lambda}| - \frac{1}{2} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{H} \mathbf{x}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{H} \mathbf{x}_k).\end{aligned}\quad (\text{B.37})$$

Thus, the average data likelihood is obtained as

$$\begin{aligned}\Theta(\boldsymbol{\Phi}, \boldsymbol{\Lambda}) &= \frac{K}{2} \log |\boldsymbol{\Phi}| - \frac{1}{2} E \left[\sum_{k=1}^K \mathbf{x}_k^T \boldsymbol{\Phi} \mathbf{x}_k \right] \\ &\quad + \frac{K}{2} \log |\boldsymbol{\Lambda}| - \frac{1}{2} E \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{H} \mathbf{x}_k)^T \boldsymbol{\Lambda} (\mathbf{y}_k - \mathbf{H} \mathbf{x}_k) \right].\end{aligned}\quad (\text{B.38})$$

B.5.3 Update Equation for Prior Precision $\boldsymbol{\Phi}$

Let us first derive the update equation for $\boldsymbol{\Phi}$. The updated value of $\boldsymbol{\Phi}$, $\widehat{\boldsymbol{\Phi}}$, is the one that maximizes the average data likelihood, $\Theta(\boldsymbol{\Phi}, \boldsymbol{\Lambda})$. The derivative of $\Theta(\boldsymbol{\Phi}, \boldsymbol{\Lambda})$ with respect to $\boldsymbol{\Phi}$ is given by

$$\frac{\partial}{\partial \boldsymbol{\Phi}} \Theta(\boldsymbol{\Phi}, \boldsymbol{\Lambda}) = \frac{K}{2} \boldsymbol{\Phi}^{-1} - \frac{1}{2} E \left[\sum_{k=1}^K \mathbf{x}_k \mathbf{x}_k^T \right].\quad (\text{B.39})$$

Here we use Eqs. (C.89) and (C.90). Setting this derivative to zero, we have

$$\widehat{\boldsymbol{\Phi}}^{-1} = \frac{1}{K} E \left[\sum_{k=1}^K \mathbf{x}_k \mathbf{x}_k^T \right].\quad (\text{B.40})$$

Using the posterior precision $\boldsymbol{\Gamma}$, the relationship

$$E \left[\sum_{k=1}^K \mathbf{x}_k \mathbf{x}_k^T \right] = \sum_{k=1}^K \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T + K \boldsymbol{\Gamma}^{-1}\quad (\text{B.41})$$

holds. Therefore, the update equation for $\boldsymbol{\Phi}$ is expressed as

$$\widehat{\boldsymbol{\Phi}}^{-1} = \frac{1}{K} \sum_{k=1}^K \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T + \boldsymbol{\Gamma}^{-1}.\quad (\text{B.42})$$

B.5.4 Update Equation for Noise Precision Λ

We next derive the update equation for Λ . The derivative of the average data likelihood with respect to Λ is given by

$$\frac{\partial}{\partial \Lambda} \Theta(\Phi, \Lambda) = \frac{K}{2} \Lambda^{-1} - \frac{1}{2} E \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{Hx}_k)(\mathbf{y}_k - \mathbf{Hx}_k)^T \right]. \quad (\text{B.43})$$

Setting this derivative to zero, we get

$$\widehat{\Lambda}^{-1} = \frac{1}{K} E \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{Hx}_k)(\mathbf{y}_k - \mathbf{Hx}_k)^T \right]. \quad (\text{B.44})$$

Using the mean and the precision of the posterior distribution, the above equation is rewritten as

$$\widehat{\Lambda}^{-1} = \frac{1}{K} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{H}\bar{x}_k)(\mathbf{y}_k - \mathbf{H}\bar{x}_k)^T + \mathbf{H}\Gamma^{-1}\mathbf{H}^T. \quad (\text{B.45})$$

B.5.5 Summary of the EM Algorithm

The EM algorithm is a recursive algorithm. We first set appropriate initial values to the hyperparameters Φ and Λ , which are then used for computing the posterior distribution, i.e., for computing the mean and the precision of the posterior distribution using

$$\begin{aligned} \Gamma &= \Phi + \mathbf{H}^T \Lambda \mathbf{H}, \\ \bar{x}_k &= \Gamma^{-1} \mathbf{H}^T \Lambda \mathbf{y}_k = (\Phi + \mathbf{H}^T \Lambda \mathbf{H})^{-1} \mathbf{H}^T \Lambda \mathbf{y}_k. \end{aligned}$$

With these parameter values of the posterior distribution, the hyperparameters Φ and Λ are updated using

$$\begin{aligned} \widehat{\Phi}^{-1} &= \frac{1}{K} \sum_{k=1}^K \bar{x}_k \bar{x}_k^T + \Gamma^{-1}, \\ \widehat{\Lambda}^{-1} &= \frac{1}{K} \sum_{k=1}^K (\mathbf{y}_k - \mathbf{H}\bar{x}_k)(\mathbf{y}_k - \mathbf{H}\bar{x}_k)^T + \mathbf{H}\Gamma^{-1}\mathbf{H}^T. \end{aligned}$$

The step that computes the posterior distribution is called the E step, and the step that estimates the hyperparameters is called the M step. The EM algorithm updates

the posterior distribution and the values of hyperparameters by repeating the E and M steps. As a result of this recursive procedure, the marginal likelihood $p(\mathbf{y}|\boldsymbol{\Phi}, \boldsymbol{\Lambda})$ is increased. A proof may be found, for example, in [1].

B.5.6 Hyperparameter Update Equations when $\boldsymbol{\Phi} = \alpha\mathbf{I}$ and $\boldsymbol{\Lambda} = \beta\mathbf{I}$

In Sect. 2.10, we show that the Bayesian estimate of $\mathbf{x}_k, \bar{\mathbf{x}}_k$, becomes equal to the L_2 -norm regularized minimum-norm solution if we use $\boldsymbol{\Phi} = \alpha\mathbf{I}$ and $\boldsymbol{\Lambda} = \beta\mathbf{I}$ in the Gaussian model. Let us derive update equations for the scalar hyperparameters α and β in this case.

The average data likelihood in this case is given by

$$\begin{aligned}\Theta(\alpha, \beta) &= \frac{NK}{2} \log \alpha - \frac{\alpha}{2} E \left[\sum_{k=1}^K \mathbf{x}_k^T \mathbf{x}_k \right] \\ &\quad + \frac{MK}{2} \log \beta - \frac{\beta}{2} E \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{Hx}_k)^T (\mathbf{y}_k - \mathbf{Hx}_k) \right].\end{aligned}$$

Therefore, using

$$\frac{\partial}{\partial \alpha} \Theta(\alpha, \beta) = \frac{NK}{2\alpha} - \frac{1}{2} E \left[\sum_{k=1}^K \mathbf{x}_k^T \mathbf{x}_k \right] = 0$$

and

$$\begin{aligned}E \left[\sum_{k=1}^K \mathbf{x}_k^T \mathbf{x}_k \right] &= \text{tr} \left[\sum_{k=1}^K E(\mathbf{x}_k \mathbf{x}_k^T) \right] \\ &= \text{tr} \left[\sum_{k=1}^K \bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T + K \boldsymbol{\Gamma}^{-1} \right] \\ &= \sum_{k=1}^K \bar{\mathbf{x}}_k^T \bar{\mathbf{x}}_k + K \text{tr}(\boldsymbol{\Gamma}^{-1}),\end{aligned}$$

we get

$$\hat{\alpha}^{-1} = \frac{1}{NK} E \left[\sum_{k=1}^K \mathbf{x}_k^T \mathbf{x}_k \right] = \frac{1}{N} \left[\frac{1}{K} \sum_{k=1}^K \|\bar{\mathbf{x}}_k\|^2 + \text{tr}(\boldsymbol{\Gamma}^{-1}) \right] \quad (\text{B.46})$$

for the update equation of α .

Regarding the hyperparameter β , using a similar derivation, we have

$$\hat{\beta}^{-1} = \frac{1}{MK} E \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{H}\mathbf{x}_k)^T (\mathbf{y}_k - \mathbf{H}\mathbf{x}_k) \right]. \quad (\text{B.47})$$

Here, considering the relationships

$$\begin{aligned} E \left[\sum_{k=1}^K (\mathbf{y}_k - \mathbf{H}\mathbf{x}_k)^T (\mathbf{y}_k - \mathbf{H}\mathbf{x}_k) \right] &= E \left[\sum_{k=1}^K \left(\mathbf{y}_k^T \mathbf{y}_k - \mathbf{x}_k^T \mathbf{H}^T \mathbf{y}_k - \mathbf{y}_k^T \mathbf{H} \mathbf{x}_k + \mathbf{x}_k^T \mathbf{H}^T \mathbf{H} \mathbf{x}_k \right) \right] \\ &= \sum_{k=1}^K \left(\mathbf{y}_k^T \mathbf{y}_k - E[\mathbf{x}_k^T] \mathbf{H}^T \mathbf{y}_k - \mathbf{y}_k^T \mathbf{H} E[\mathbf{x}_k] + E[\mathbf{x}_k^T \mathbf{H}^T \mathbf{H} \mathbf{x}_k] \right) \\ &= \sum_{k=1}^K \left(\mathbf{y}_k^T \mathbf{y}_k - \bar{\mathbf{x}}_k^T \mathbf{H}^T \mathbf{y}_k - \mathbf{y}_k^T \mathbf{H} \bar{\mathbf{x}}_k + E[\mathbf{x}_k^T \mathbf{H}^T \mathbf{H} \mathbf{x}_k] \right) \end{aligned} \quad (\text{B.48})$$

and

$$\begin{aligned} E[\mathbf{x}_k^T \mathbf{H}^T \mathbf{H} \mathbf{x}_k] &= E[\text{tr}(\mathbf{H}^T \mathbf{H} \mathbf{x}_k \mathbf{x}_k^T)] \\ &= \text{tr} \left[\mathbf{H}^T \mathbf{H} E \left(\mathbf{x}_k \mathbf{x}_k^T \right) \right] \\ &= \text{tr} \left[\mathbf{H}^T \mathbf{H} \left(\bar{\mathbf{x}}_k \bar{\mathbf{x}}_k^T + \boldsymbol{\Gamma}^{-1} \right) \right] \\ &= \bar{\mathbf{x}}_k^T \mathbf{H}^T \mathbf{H} \bar{\mathbf{x}}_k + \text{tr} \left(\mathbf{H}^T \mathbf{H} \boldsymbol{\Gamma}^{-1} \right), \end{aligned} \quad (\text{B.49})$$

we finally obtain

$$\hat{\beta}^{-1} = \frac{1}{M} \left[\frac{1}{K} \sum_{k=1}^K \|\mathbf{y}_k - \mathbf{H} \bar{\mathbf{x}}_k\|^2 + \text{tr} \left(\mathbf{H}^T \mathbf{H} \boldsymbol{\Gamma}^{-1} \right) \right]. \quad (\text{B.50})$$

B.6 Variational Bayesian Inference

In the preceding sections, the unknown \mathbf{x} is estimated based on the posterior distribution $p(\mathbf{x}|\mathbf{y})$, and the hyperparameters are estimated by using the EM algorithm. In this section, we introduce a method called variational Bayesian inference, which makes it possible to derive approximate posterior distributions for hyperparameters [2].

B.6.1 Derivation of the EM Algorithm Using the Free Energy

B.6.1.1 Derivation of Posterior Distribution (E-step)

As a preparation for introducing the variational technique, we derive the EM algorithm in a different manner based on an optimization of a functional called the free energy. In this section, the hyperparameters are collectively expressed as θ . We define a functional such that,

$$\mathcal{F}[q(\mathbf{x}), \theta] = \int d\mathbf{x} q(\mathbf{x}) [\log p(\mathbf{x}, \mathbf{y}|\theta) - \log q(\mathbf{x})]. \quad (\text{B.51})$$

This $\mathcal{F}[q(\mathbf{x}), \theta]$ is a function of hyperparameters θ and an arbitrary probability distribution $q(\mathbf{x})$. This $\mathcal{F}[q(\mathbf{x}), \theta]$ is called the free energy using a terminology in statistical physics. We show, in the following, that maximizing the free energy $\mathcal{F}[q(\mathbf{x}), \theta]$ with respect to $q(\mathbf{x})$ results in the E step, and maximizing it with respect to the hyperparameters results in the M step of the EM algorithm.

When maximizing $\mathcal{F}[q(\mathbf{x}), \theta]$ with respect to $q(\mathbf{x})$, since $q(\mathbf{x})$ is a probability distribution, the constraint $\int_{-\infty}^{\infty} q(\mathbf{x}) d\mathbf{x} = 1$ must be imposed. Therefore, this maximization problem can be formulated such that,

$$\hat{q}(\mathbf{x}) = \underset{q(\mathbf{x})}{\operatorname{argmax}} \mathcal{F}[q(\mathbf{x}), \theta], \quad \text{subject to} \quad \int_{-\infty}^{\infty} q(\mathbf{x}) d\mathbf{x} = 1. \quad (\text{B.52})$$

Such a constrained optimization problem can be solved by using the method of Lagrange multipliers, in which defining the Lagrange multiplier as γ , the Lagrangian is defined as

$$\begin{aligned} \mathbb{L}[q, \gamma] &= \mathcal{F}[q, \theta] + \gamma \left[\int_{-\infty}^{\infty} q(\mathbf{x}) d\mathbf{x} - 1 \right] \\ &= \int_{-\infty}^{\infty} d\mathbf{x} q(\mathbf{x}) [\log p(\mathbf{x}, \mathbf{y}|\theta) - \log q(\mathbf{x})] + \gamma \left[\int_{-\infty}^{\infty} q(\mathbf{x}) d\mathbf{x} - 1 \right]. \end{aligned} \quad (\text{B.53})$$

The constrained optimization problem in Eq. (B.52) is now rewritten as the unconstrained optimization problem in Eq. (B.53). The probability distribution $q(\mathbf{x})$ that maximizes the Lagrangian $\mathbb{L}[q, \gamma]$ is the solution of the constrained optimization problem in Eq. (B.52).

Differentiating $\mathbb{L}[q, \gamma]$ with respect to $q(\mathbf{x})$, and setting the derivative to zero, we have

$$\frac{\delta \mathbb{L}[q(\mathbf{x}), \gamma]}{\delta q(\mathbf{x})} = \log p(\mathbf{x}, \mathbf{y}|\theta) - \log q(\mathbf{x}) - 1 + \gamma = 0. \quad (\text{B.54})$$

A brief explanation on the differentiation of a functional, as well as the derivation of Eq. (B.54), is presented in Sect. C.5 in the Appendix. Differentiating $\mathbb{L}[q, \gamma]$ with respect to γ gives

$$\frac{\partial \mathbb{L}[q(\mathbf{x}), \gamma]}{\partial \gamma} = \int_{-\infty}^{\infty} q(\mathbf{x}) d\mathbf{x} - 1 = 0. \quad (\text{B.55})$$

Thus, using Eq. (B.54), we have

$$\hat{q}(\mathbf{x}) = e^{\gamma-1} p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}), \quad (\text{B.56})$$

and using Eq. (B.55), we also have

$$\int_{-\infty}^{\infty} q(\mathbf{x}) d\mathbf{x} = e^{\gamma-1} \int_{-\infty}^{\infty} p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) d\mathbf{x} = e^{\gamma-1} p(\mathbf{y}|\boldsymbol{\theta}) = 1. \quad (\text{B.57})$$

We thereby get

$$e^{\gamma-1} = \frac{1}{p(\mathbf{y}|\boldsymbol{\theta})}. \quad (\text{B.58})$$

Substitution of Eq. (B.58) into (B.56) results in the relationship

$$\hat{q}(\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})}{p(\mathbf{y}|\boldsymbol{\theta})} = p(\mathbf{x}|\mathbf{y}). \quad (\text{B.59})$$

The above equation shows that the probability distribution that maximizes the free energy $\mathcal{F}[q(\mathbf{x}), \boldsymbol{\theta}]$ is the posterior distribution $p(\mathbf{x}|\mathbf{y})$.

Note that to derive the posterior distribution, we do not explicitly use Bayes' rule. Instead, we use the optimization of the functional called the free energy. This idea is further extended in variational Bayesian inference in the following sections to derive an approximate posterior distribution in a more complicated situation.

B.6.1.2 Derivation of M-step

We next maximize the free energy with respect to the hyperparameter $\boldsymbol{\theta}$. Once $\mathcal{F}[q, \boldsymbol{\theta}]$ is maximized with respect to $q(\mathbf{x})$, the free energy is written as

$$\begin{aligned} \mathcal{F}[p(\mathbf{x}|\mathbf{y}), \boldsymbol{\theta}] &= \int d\mathbf{x} p(\mathbf{x}|\mathbf{y}) [\log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) - \log p(\mathbf{x}|\mathbf{y})] \\ &= \Theta(\boldsymbol{\theta}) + \mathcal{H}[p(\mathbf{x}|\mathbf{y})], \end{aligned} \quad (\text{B.60})$$

where

$$\Theta(\boldsymbol{\theta}) = \int d\mathbf{x} p(\mathbf{x}|\mathbf{y}) \log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}). \quad (\text{B.61})$$

This $\Theta(\boldsymbol{\theta})$ is equal to the average data likelihood, and

$$\mathcal{H}[p(\mathbf{x}|\mathbf{y})] = - \int d\mathbf{x} p(\mathbf{x}|\mathbf{y}) \log p(\mathbf{x}|\mathbf{y}) \quad (\text{B.62})$$

is the entropy of the posterior distribution. Since the entropy does not depend on θ , maximizing $\mathcal{F}[p(\mathbf{x}|\mathbf{y}), \theta]$ with respect to θ is equivalent to maximizing $\Theta(\theta)$ with respect to θ . Namely, the maximization of the free energy with respect to the hyperparameters results in the M step of the EM algorithm.

Note also that the free energy after the maximization with respect to $q(\mathbf{x})$ is equal to the marginal likelihood, $\log p(\mathbf{y}|\theta)$. This can be seen by rewriting Eq. (B.60) such that

$$\begin{aligned}\mathcal{F}[p(\mathbf{x}|\mathbf{y}), \theta] &= \int d\mathbf{x} p(\mathbf{x}|\mathbf{y}) [\log p(\mathbf{x}, \mathbf{y}|\theta) - \log p(\mathbf{x}|\mathbf{y})] \\ &= \int d\mathbf{x} p(\mathbf{x}|\mathbf{y}) \log p(\mathbf{y}|\theta) = \log p(\mathbf{y}|\theta).\end{aligned}\quad (\text{B.63})$$

This relationship is used to derive the expressions of the marginal likelihood for the Bayesian factor analysis in Chap. 5.

For an arbitrary probability distribution $q(\mathbf{x})$, the relationship between the free energy and the marginal likelihood is expressed as

$$\mathcal{F}[q(\mathbf{x}), \theta] = \log p(\mathbf{y}|\theta) - \mathcal{K}_L [q(\mathbf{x})||p(\mathbf{x}|\mathbf{y})], \quad (\text{B.64})$$

where $\mathcal{K}_L [q(\mathbf{x})||p(\mathbf{x}|\mathbf{y})]$ is the Kullback-Leibler (KL) distance defined in

$$\mathcal{K}_L [q(\mathbf{x})||p(\mathbf{x}|\mathbf{y})] = \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x}|\mathbf{y})} d\mathbf{x}. \quad (\text{B.65})$$

The KL distance represents a distance between the true posterior distribution $p(\mathbf{x}|\mathbf{y})$ and the arbitrary probability distribution $q(\mathbf{x})$. It always has a nonnegative value, and is equal to zero when the two distributions are identical. Hence, for an arbitrary $q(\mathbf{x})$, the inequality $\mathcal{F}[q(\mathbf{x}), \theta] \leq \log p(\mathbf{y}|\theta)$ holds, and the free energy forms a lower-bound of the marginal likelihood.

B.6.2 Variational Bayesian EM Algorithm

In Bayesian inference, to estimate the unknown parameter \mathbf{x} , we must first derive the posterior distribution $p(\mathbf{x}|\mathbf{y})$, assuming the existence of an appropriate prior distribution $p(\mathbf{x}|\theta)$. We then obtain an optimum estimate of the unknown \mathbf{x} based on the posterior distribution $p(\mathbf{x}|\mathbf{y})$. When the hyperparameter θ is unknown, a truly Bayesian approach is to first derive the joint posterior distribution $p(\mathbf{x}, \theta|\mathbf{y})$, and to estimate \mathbf{x} and θ simultaneously based on this joint posterior distribution. To derive the joint posterior $p(\mathbf{x}, \theta|\mathbf{y})$, we can use Bayes' rule, and obtain,

$$p(\mathbf{x}, \theta|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x}, \theta)p(\mathbf{x}, \theta) = p(\mathbf{y}|\mathbf{x}, \theta)p(\mathbf{x}|\theta)p(\theta). \quad (\text{B.66})$$

In order to compute the posterior distribution $p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$, we must compute the integral

$$Z = \iint_{-\infty}^{\infty} p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\mathbf{x})p(\boldsymbol{\theta})d\mathbf{x}d\boldsymbol{\theta}.$$

However, in general, this integral Z does not have a closed-form solution, and it is hard to compute the posterior distribution $p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ using Bayes' rule, i.e., Eq. (B.66).

In Sect. B.6.1, the posterior distribution is obtained using the optimization of the functional called the free energy, and not using Bayes' rule. We here use this variational technique to obtain the posterior distribution. The free energy is expressed as

$$\mathcal{F}[q(\mathbf{x}, \boldsymbol{\theta}), \boldsymbol{\theta}] = \int d\boldsymbol{\theta} d\mathbf{x} q(\mathbf{x}, \boldsymbol{\theta}) [\log p(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) - \log q(\mathbf{x}, \boldsymbol{\theta})]. \quad (\text{B.67})$$

We use an approximation that the joint posterior distribution is factorized, i.e.,

$$q(\mathbf{x}, \boldsymbol{\theta}) = q(\mathbf{x})q(\boldsymbol{\theta}), \quad (\text{B.68})$$

which is called the variational approximation. Using this approximation, the free energy is given by

$$\mathcal{F}[q(\mathbf{x}), q(\boldsymbol{\theta}), \boldsymbol{\theta}] = \int d\boldsymbol{\theta} d\mathbf{x} q(\mathbf{x})q(\boldsymbol{\theta}) [\log p(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) - \log q(\mathbf{x}) - \log q(\boldsymbol{\theta})]. \quad (\text{B.69})$$

The best estimate of the posterior distributions $p(\mathbf{x}|\mathbf{y})$ and $p(\boldsymbol{\theta}|\mathbf{y})$ are derived by jointly maximizing the free energy in Eq. (B.69). That is, the estimated posterior distributions, $\hat{p}(\mathbf{x}|\mathbf{y})$ and $\hat{p}(\boldsymbol{\theta}|\mathbf{y})$, are given by

$$\begin{aligned} \hat{p}(\mathbf{x}|\mathbf{y}), \hat{p}(\boldsymbol{\theta}|\mathbf{y}) &= \underset{q(\mathbf{x}), q(\boldsymbol{\theta})}{\operatorname{argmax}} \mathcal{F}[q(\mathbf{x}), q(\boldsymbol{\theta}), \boldsymbol{\theta}], \\ \text{subject to } \int_{-\infty}^{\infty} d\mathbf{x} q(\mathbf{x}) &= 1, \text{ and } \int_{-\infty}^{\infty} d\boldsymbol{\theta} q(\boldsymbol{\theta}) = 1. \end{aligned} \quad (\text{B.70})$$

Therefore, according to Eq. (B.64), the estimated posterior distributions $\hat{p}(\mathbf{x}|\mathbf{y})$ and $\hat{p}(\boldsymbol{\theta}|\mathbf{y})$ satisfy the relationship:

$$\begin{aligned} \hat{p}(\mathbf{x}|\mathbf{y}), \hat{p}(\boldsymbol{\theta}|\mathbf{y}) &= \underset{q(\mathbf{x}), q(\boldsymbol{\theta})}{\operatorname{argmin}} \mathcal{K}_L [p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) || q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})], \\ \text{where } q(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) &= q(\mathbf{x}|\mathbf{y})q(\boldsymbol{\theta}|\mathbf{y}). \end{aligned} \quad (\text{B.71})$$

Namely, the estimated posterior distributions $\hat{p}(\mathbf{x}|\mathbf{y})$ and $\hat{p}(\boldsymbol{\theta}|\mathbf{y})$ obtained by maximizing the free energy jointly minimizes the KL distance between the true and approximated posterior distributions. In other words, the estimated joint posterior $\hat{p}(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) = \hat{p}(\mathbf{x}|\mathbf{y})\hat{p}(\boldsymbol{\theta}|\mathbf{y})$ is the “best” estimate in the sense that it minimizes the KL distance from the true posterior distribution.

Let us first derive $\hat{p}(\mathbf{x}|y)$. Using arguments similar to those in Sect. B.6.1, defining the Lagrangian as

$$\begin{aligned}\mathbb{L}[q(\mathbf{x})] &= \int d\mathbf{x} d\boldsymbol{\theta} q(\mathbf{x}) q(\boldsymbol{\theta}) [\log p(\mathbf{x}, y, \boldsymbol{\theta}) - \log q(\mathbf{x}) - \log q(\boldsymbol{\theta})] \\ &\quad + \gamma \left[\int_{-\infty}^{\infty} q(\mathbf{x}) d\mathbf{x} - 1 \right],\end{aligned}\tag{B.72}$$

differentiating it with respect to $q(\mathbf{x})$, and setting the derivative to zero, we have

$$\int d\boldsymbol{\theta} q(\boldsymbol{\theta}) [\log p(\mathbf{x}, y, \boldsymbol{\theta}) - \log q(\mathbf{x})] + \mathcal{C} = 0,\tag{B.73}$$

where

$$\mathcal{C} = - \int d\boldsymbol{\theta} q(\boldsymbol{\theta}) \log q(\boldsymbol{\theta}) - 1 + \gamma.$$

Neglecting \mathcal{C} , we obtain

$$\log \hat{p}(\mathbf{x}|y) = \int d\boldsymbol{\theta} q(\boldsymbol{\theta}) \log p(\mathbf{x}, y, \boldsymbol{\theta}) = E_{\boldsymbol{\theta}} [\log p(\mathbf{x}, y, \boldsymbol{\theta})],\tag{B.74}$$

where $E_{\boldsymbol{\theta}} [\cdot]$ indicates computing the mean with respect to the posterior distribution $q(\boldsymbol{\theta})$. Using exactly the same derivation, we obtain the relationship,

$$\log \hat{p}(\boldsymbol{\theta}|y) = \int d\mathbf{x} q(\mathbf{x}) \log p(\mathbf{x}, y, \boldsymbol{\theta}) = E_{\mathbf{x}} [\log p(\mathbf{x}, y, \boldsymbol{\theta})],\tag{B.75}$$

where $E_{\mathbf{x}} [\cdot]$ indicates computing the mean with respect to the posterior distribution $q(\mathbf{x})$.

Equations (B.74) and (B.75) indicate that, to compute the posterior distribution $\hat{p}(\mathbf{x}|y)$, we need $q(\boldsymbol{\theta})$, namely $\hat{p}(\boldsymbol{\theta}|y)$, and to compute the posterior distribution $\hat{p}(\boldsymbol{\theta}|y)$, we need $q(\mathbf{x})$, namely $\hat{p}(\mathbf{x}|y)$. Therefore, in variational Bayesian inference, we need an EM-like recursive procedure. We update $\hat{p}(\mathbf{x}|y)$ and $\hat{p}(\boldsymbol{\theta}|y)$ by repeatedly applying Eqs. (B.74) and (B.75). When computing $\hat{p}(\mathbf{x}|y)$ using Eq. (B.74), we use $\hat{p}(\boldsymbol{\theta}|y)$ obtained in the preceding step. When computing $\hat{p}(\boldsymbol{\theta}|y)$ using Eq. (B.75), we use $\hat{p}(\mathbf{x}|y)$ obtained in the preceding step. This recursive algorithm is referred to as the variational Bayesian EM (VBEM) algorithm.

References

1. C. M. Bishop, *Pattern recognition and machine learning*. Springer, New York, 2006.
2. H. T. Attius, “A variational Bayesian framework for graphical models,” in *Advances in Neural information processing*, pp. 209–215, MIT Press, 2000.

Appendix C

Supplementary Mathematical Arguments

C.1 Multi-dimensional Gaussian Distribution

Let us define a column vector of N random variables as \mathbf{x} . The multi-dimensional Gaussian distribution for \mathbf{x} is expressed as

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{N/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]. \quad (\text{C.1})$$

Here $\boldsymbol{\mu}$ is the mean of \mathbf{x} defined as $\boldsymbol{\mu} = E(\mathbf{x})$ and $\boldsymbol{\Sigma}$ is the covariance matrix of \mathbf{x} defined as $\boldsymbol{\Sigma} = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T]$ where $E(\cdot)$ is the expectation operator. Also, $|\boldsymbol{\Sigma}|$ indicates the determinant of $\boldsymbol{\Sigma}$. The Gaussian distribution in Eq. (C.1) is often written as

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

When the random variable \mathbf{x} follows a Gaussian distribution with mean $\boldsymbol{\mu}$ and the covariance matrix $\boldsymbol{\Sigma}$, the expression

$$\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (\text{C.2})$$

is often used.

Two vector random variables \mathbf{x}_1 and \mathbf{x}_2 have the linear relationship,

$$\mathbf{x}_2 = \mathbf{A}\mathbf{x}_1 + \mathbf{c},$$

where \mathbf{A} is a matrix of deterministic variables, and \mathbf{c} is a vector of deterministic variables. Then, if we have

$$\mathbf{x}_1 \sim \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

we also have

$$\mathbf{x}_2 \sim \mathcal{N}(\mathbf{x}_2 | \mathbf{A}\boldsymbol{\mu} + \mathbf{c}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T). \quad (\text{C.3})$$

Namely, \mathbf{x}_2 obeys the Gaussian distribution with a mean equal to $\mathbf{A}\boldsymbol{\mu} + \mathbf{c}$ and a covariance matrix equal to $\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T$.

Bayesian inference quite often uses the precision, instead of the variance. Let us define the precision matrix $\boldsymbol{\Lambda}$ that corresponds to the covariance matrix $\boldsymbol{\Sigma}$, i.e., $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$. When the precision $\boldsymbol{\Lambda}$ is used, the notation in Eq. (C.2) is used such that

$$\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}). \quad (\text{C.4})$$

Here, since we maintain the notational convenience

$$\mathcal{N}(\text{random variable}|\text{mean, covariance matrix}),$$

we must use $\boldsymbol{\Lambda}^{-1}$ in this notation. Using the precision matrix, $\boldsymbol{\Lambda}$, the explicit form of the Gaussian distribution is given by

$$p(\mathbf{x}) = \frac{|\boldsymbol{\Lambda}|^{1/2}}{(2\pi)^{N/2}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}) \right]. \quad (\text{C.5})$$

Let us compute the entropy when the probability distribution is Gaussian. The definition of the entropy for a continuous random variable is

$$\mathcal{H}[p(\mathbf{x})] = - \int p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x} = -E[\log p(\mathbf{x})]. \quad (\text{C.6})$$

Substituting the probability distribution in Eq. (C.1) into the equation above, we get

$$\begin{aligned} \mathcal{H}[p(\mathbf{x})] &= -E \left[\log \left[\frac{1}{(2\pi)^{N/2} |\boldsymbol{\Sigma}|^{1/2}} \right] \right] \\ &\quad + \frac{1}{2} E \left[(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right], \end{aligned} \quad (\text{C.7})$$

where constant terms are omitted. In the equation above, the first term on the right-hand side is equal to $\frac{1}{2} \log |\boldsymbol{\Sigma}|$. The second term is equal to

$$\begin{aligned} \frac{1}{2} E \left[(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] &= \frac{1}{2} E \left[\text{tr} \left[\boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \right] \right] \\ &= \frac{1}{2} \text{tr} \left[\boldsymbol{\Sigma}^{-1} E \left[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \right] \right] \\ &= \frac{1}{2} \text{tr} \left[\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \right] = \frac{N}{2}. \end{aligned} \quad (\text{C.8})$$

Therefore, omitting constant terms, the entropy is expressed as

$$\mathcal{H}[p(\mathbf{x})] = \mathcal{H}(\mathbf{x}) = \frac{1}{2} \log |\boldsymbol{\Sigma}|. \quad (\text{C.9})$$

Note that we use a simplified notation $\mathcal{H}(\mathbf{x})$ in this book to indicate the entropy of the probability distribution $p(\mathbf{x})$, which should formally be written as $\mathcal{H}[p(\mathbf{x})]$ ³

C.2 Complex Gaussian Distribution

We derive the probability distribution when the random variable is a complex-valued Gaussian. The arguments in this section follows those of Neeser and Massey [1]. Let us define a column vector of N complex random variables as \mathbf{z} , and define $\mathbf{z} = \mathbf{x} + i\mathbf{y}$ where $\mathbf{x} = \Re(\mathbf{z})$, $\mathbf{y} = \Im(\mathbf{z})$ and \mathbf{x} and \mathbf{y} are real-valued N -dimensional Gaussian random vectors. Here, we assume that $E(\mathbf{z}) = 0$, and as a result, $E(\mathbf{x}) = 0$ and $E(\mathbf{y}) = 0$. The covariance matrix of \mathbf{z} is defined as $\boldsymbol{\Sigma}_{zz} = E(\mathbf{z}\mathbf{z}^H)$.

We next define a $(2N \times 1)$ vector $\boldsymbol{\phi}$ such that $\boldsymbol{\phi} = [\mathbf{x}^T, \mathbf{y}^T]^T$. The covariance matrix of $\boldsymbol{\phi}$, $\boldsymbol{\Sigma}_{\phi\phi}$, is given by

$$\boldsymbol{\Sigma}_{\phi\phi} = E\left[\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \begin{bmatrix} \mathbf{x}^T & \mathbf{y}^T \end{bmatrix}\right] = \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{yx}^T \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix}, \quad (\text{C.10})$$

where $\boldsymbol{\Sigma}_{xx} = E(\mathbf{x}\mathbf{x}^T)$, $\boldsymbol{\Sigma}_{yx} = E(\mathbf{y}\mathbf{x}^T)$, and $\boldsymbol{\Sigma}_{yy} = E(\mathbf{y}\mathbf{y}^T)$. We assume that the joint distribution of \mathbf{x} and \mathbf{y} is given by

$$p(\mathbf{x}, \mathbf{y}) = p(\boldsymbol{\phi}) = \frac{1}{(2\pi)^N |\boldsymbol{\Sigma}_{\phi\phi}|^{1/2}} \exp\left[-\frac{1}{2} \boldsymbol{\phi}^T \boldsymbol{\Sigma}_{\phi\phi}^{-1} \boldsymbol{\phi}\right]. \quad (\text{C.11})$$

In this section, we show that this joint distribution is equal to

$$p(\mathbf{z}) = \frac{1}{\pi^N |\boldsymbol{\Sigma}_{zz}|} \exp\left[-\mathbf{z}^H \boldsymbol{\Sigma}_{zz}^{-1} \mathbf{z}\right], \quad (\text{C.12})$$

which is called the complex Gaussian distribution.

To show this equality, we assume a property, called ‘‘proper’’, on the complex random variable \mathbf{z} . The complex random variable \mathbf{z} is proper if its pseudo covariance $\bar{\boldsymbol{\Sigma}}_{zz}$ is equal to zero. The pseudo covariance is defined such that $\bar{\boldsymbol{\Sigma}}_{zz} = E(\mathbf{z}\mathbf{z}^T)$, which is written as

$$\bar{\boldsymbol{\Sigma}}_{zz} = E\left[(\mathbf{x} + i\mathbf{y})(\mathbf{x} + i\mathbf{y})^T\right] = \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{yy} + i\left(\boldsymbol{\Sigma}_{yx} + \boldsymbol{\Sigma}_{yx}^T\right).$$

Therefore, $\bar{\boldsymbol{\Sigma}}_{zz} = \mathbf{0}$ is equivalent to the relationships

$$\boldsymbol{\Sigma}_{xx} = \boldsymbol{\Sigma}_{yy} \quad \text{and} \quad -\boldsymbol{\Sigma}_{yx} = \boldsymbol{\Sigma}_{yx}^T. \quad (\text{C.13})$$

³ The notation $\mathcal{H}(\mathbf{x})$ may look as if the entropy $\mathcal{H}(\mathbf{x})$ is a function of \mathbf{x} , but the entropy is a functional of the probability distribution $p(\mathbf{x})$, not a function of \mathbf{x} .

Note that proper complex Gaussian random vectors with zero mean are called circular. In this section, we assume that the circularity holds for \mathbf{z} . Under the assumption that \mathbf{z} is a proper complex random vector, we derive

$$\boldsymbol{\Sigma}_{zz} = E \left[(\mathbf{x} + i\mathbf{y})(\mathbf{x} - i\mathbf{y})^T \right] = 2 \left(\boldsymbol{\Sigma}_{xx} + i\boldsymbol{\Sigma}_{yx} \right), \quad (\text{C.14})$$

and $\boldsymbol{\Sigma}_{zz}^{-1}$ is found to be⁴

$$\boldsymbol{\Sigma}_{zz}^{-1} = \frac{1}{2} \boldsymbol{\Delta}^{-1} \left(\mathbf{I} - i\boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \right), \quad (\text{C.15})$$

where

$$\boldsymbol{\Delta} = \boldsymbol{\Sigma}_{xx} + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx}. \quad (\text{C.16})$$

We first derive the relationship between the determinants $|\boldsymbol{\Sigma}_{zz}|$ and $|\boldsymbol{\Sigma}_{\phi\phi}|$. Using Eq. (C.14), we have

$$\boldsymbol{\Sigma}_{zz}^T = 2 \left(\boldsymbol{\Sigma}_{xx} + i\boldsymbol{\Sigma}_{yx}^T \right) = 2 \left(\boldsymbol{\Sigma}_{xx} - i\boldsymbol{\Sigma}_{yx} \right),$$

and $\boldsymbol{\Sigma}_{zz}^{-1}$ is expressed as

$$\boldsymbol{\Sigma}_{zz}^{-1} = \frac{1}{2} \boldsymbol{\Delta}^{-1} \left(\boldsymbol{\Sigma}_{xx} - i\boldsymbol{\Sigma}_{yx} \right) \boldsymbol{\Sigma}_{xx}^{-1} = \frac{1}{4} \boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma}_{zz}^T \boldsymbol{\Sigma}_{xx}^{-1}.$$

Taking the determinant of both sides of the equation above, we get

$$|\boldsymbol{\Sigma}_{zz}|^{-1} = \frac{1}{4^N} |\boldsymbol{\Delta}|^{-1} |\boldsymbol{\Sigma}_{zz}| |\boldsymbol{\Sigma}_{xx}|^{-1},$$

and finally,

$$|\boldsymbol{\Sigma}_{zz}| = 2^N \sqrt{|\boldsymbol{\Sigma}_{xx}| |\boldsymbol{\Delta}|}. \quad (\text{C.17})$$

On the other hand, using the determinant identity in Eq. (C.94), $|\boldsymbol{\Sigma}_{\phi\phi}|$ is expressed as

$$|\boldsymbol{\Sigma}_{\phi\phi}| = \begin{vmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{yx}^T \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{vmatrix} = |\boldsymbol{\Sigma}_{xx}| |\boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx}^T| = |\boldsymbol{\Sigma}_{xx}| |\boldsymbol{\Delta}|, \quad (\text{C.18})$$

where the relationships in Eq. (C.13) are used to obtain the right-most expression. Comparing Eqs. (C.17) and (C.18), we finally obtain the relationship,

$$|\boldsymbol{\Sigma}_{zz}| = 2^N \sqrt{|\boldsymbol{\Sigma}_{\phi\phi}|}. \quad (\text{C.19})$$

⁴ It is easy to see that $\boldsymbol{\Sigma}_{zz}^{-1} \boldsymbol{\Sigma}_{zz} = \mathbf{I}$ holds.

We next show the equality of the two quadratic forms, i.e.,

$$\mathbf{z}^H \boldsymbol{\Sigma}_{zz}^{-1} \mathbf{z} = \frac{1}{2} \boldsymbol{\phi}^T \boldsymbol{\Sigma}_{\phi\phi}^{-1} \boldsymbol{\phi}.$$

To show this relationship, we use the matrix inversion formula in Eq. (C.93), and rewrite $\boldsymbol{\Sigma}_{\phi\phi}^{-1}$ such that⁵

$$\boldsymbol{\Sigma}_{\phi\phi}^{-1} = \begin{bmatrix} \Delta^{-1} & \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx} \Delta^{-1} \\ -\Delta^{-1} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} & \Delta^{-1} \end{bmatrix}, \quad (\text{C.20})$$

where Δ is defined in Eq. (C.16). Let us define:

$$\mathbf{M}_c = \Delta^{-1}, \quad (\text{C.21})$$

$$\mathbf{M}_s = -\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx} \Delta^{-1}. \quad (\text{C.22})$$

Then, \mathbf{M}_c is symmetric because Δ is symmetric. Using

$$\Delta \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx} = \left(\boldsymbol{\Sigma}_{xx} + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx} \right) \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx} = \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \Delta,$$

we also get the relationship:

$$\begin{aligned} \mathbf{M}_s &= -\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx} \Delta^{-1} = -\Delta^{-1} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \Delta \Delta^{-1} = -\Delta^{-1} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \\ &= -\left(\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx}^T \Delta^{-1} \right)^T = \left(\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx} \Delta^{-1} \right)^T = -\mathbf{M}_s^T, \end{aligned} \quad (\text{C.23})$$

where $\boldsymbol{\Sigma}_{yx}^T = -\boldsymbol{\Sigma}_{yx}$ is used. Thus, \mathbf{M}_s is skew-symmetric.

We define a matrix \mathbf{M} as

$$\mathbf{M} = (\mathbf{M}_c + i\mathbf{M}_s), \quad (\text{C.24})$$

and compute the quadratic form $\mathbf{z}^H \mathbf{M} \mathbf{z}$. Since we have the relationship,

$$\begin{aligned} \left(\mathbf{z}^H \mathbf{M} \mathbf{z} \right)^H &= \mathbf{z}^H \mathbf{M}^H \mathbf{z} = \mathbf{z}^H (\mathbf{M}_c + i\mathbf{M}_s)^H \mathbf{z} \\ &= \mathbf{z}^H \left(\mathbf{M}_c^T - i\mathbf{M}_s^T \right) \mathbf{z} = \mathbf{z}^H (\mathbf{M}_c + i\mathbf{M}_s) \mathbf{z} = \mathbf{z}^H \mathbf{M} \mathbf{z}, \end{aligned} \quad (\text{C.25})$$

⁵ In Eq. (C.20), the (2, 2)th element is Δ^{-1} , which can be shown as follows. According to Eq. (C.93), this element is equal to

$$\boldsymbol{\Sigma}_{yy}^{-1} + \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} \left(\boldsymbol{\Sigma}_{xx} + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx} \right) \boldsymbol{\Sigma}_{yx}^T \boldsymbol{\Sigma}_{yy}^{-1}.$$

Using Eqs. (C.13) and (C.91), this is equal to $(\boldsymbol{\Sigma}_{xx} + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{yx})^{-1} = \Delta^{-1}$.

the quadratic form $\mathbf{z}^H \mathbf{M} \mathbf{z}$ is real-valued. Therefore, we have

$$\begin{aligned}\mathbf{z}^H \mathbf{M} \mathbf{z} &= \Re \left[\mathbf{z}^H \mathbf{M} \mathbf{z} \right] = \Re \left[(\mathbf{x}^T - i\mathbf{y}^T) (\mathbf{M}_c + i\mathbf{M}_s) (\mathbf{x} + i\mathbf{y}) \right] \\ &= \mathbf{x}^T \mathbf{M}_c \mathbf{x} - \mathbf{x}^T \mathbf{M}_s \mathbf{y} + \mathbf{y}^T \mathbf{M}_s \mathbf{x} + \mathbf{y}^T \mathbf{M}_c \mathbf{y} \\ &= [\mathbf{x}^T, \mathbf{y}^T] \begin{bmatrix} \mathbf{M}_c & -\mathbf{M}_s \\ \mathbf{M}_s & \mathbf{M}_c \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \boldsymbol{\phi}^T \boldsymbol{\Sigma}_{\phi\phi}^{-1} \boldsymbol{\phi}.\end{aligned}\quad (\text{C.26})$$

We now compute \mathbf{M} , such that

$$\mathbf{M} = \boldsymbol{\Delta}^{-1} - i\boldsymbol{\Delta}^{-1} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} = \boldsymbol{\Delta}^{-1} \left(\mathbf{I} - i\boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \right).$$

By comparing the above \mathbf{M} with Eq. (C.15) we can see that $\boldsymbol{\Sigma}_{zz}^{-1} = \frac{1}{2}\mathbf{M}$, and therefore, from Eq. (C.26), we can prove the relationship

$$\mathbf{z}^H \boldsymbol{\Sigma}_{zz}^{-1} \mathbf{z} = \frac{1}{2} \boldsymbol{\phi}^T \boldsymbol{\Sigma}_{\phi\phi}^{-1} \boldsymbol{\phi}. \quad (\text{C.27})$$

On the basis of Eqs. (C.17) and (C.27), it is now clear that the real-valued joint Gaussian distribution in Eq. (C.11) and the complex Gaussian distribution in Eq. (C.12) are equivalent. Using exactly the same derivation for Eq. (C.9) and ignoring the constants, the entropy for the complex-valued Gaussian is obtained as

$$\mathcal{H}(\mathbf{z}) = \log |\boldsymbol{\Sigma}_{zz}|. \quad (\text{C.28})$$

C.3 Canonical Correlation and Mutual Information

C.3.1 Canonical Correlation

This appendix provides a concise explanation on canonical correlation. Let us define real-valued random vectors \mathbf{x} and \mathbf{y} such that

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_q \end{bmatrix}, \quad (\text{C.29})$$

and consider computing the correlation between \mathbf{x} and \mathbf{y} . One way is to compute correlation coefficients between all combinations of (x_i, y_j) . However, this gives total $p \times q$ correlation coefficients and the interpretation of these results may not be easy. The canonical correlation method first projects column vectors \mathbf{x} and \mathbf{y} onto the

direction of \mathbf{a} and \mathbf{b} such that $\hat{x} = \mathbf{a}^T \mathbf{x}$ and $\hat{y} = \mathbf{b}^T \mathbf{y}$, where \mathbf{a} and \mathbf{b} are real-valued vectors. It then computes the correlation between \hat{x} and \hat{y} . This correlation depends on the choices of \mathbf{a} and \mathbf{b} . The maximum correlation between \hat{x} and \hat{y} is considered to represent the correlation between \mathbf{x} and \mathbf{y} , and this maximum correlation is called the canonical correlation.

The correlation between \hat{x} and \hat{y} , ρ , is expressed as

$$\begin{aligned}\rho &= \frac{E(\hat{x}\hat{y})}{\sqrt{E(\hat{x}^2)}\sqrt{E(\hat{y}^2)}} \\ &= \frac{\mathbf{a}^T E(\mathbf{x}\mathbf{y}^T) \mathbf{b}}{\sqrt{[\mathbf{a}^T E(\mathbf{x}\mathbf{x}^T) \mathbf{a}] [\mathbf{b}^T E(\mathbf{y}\mathbf{y}^T) \mathbf{b}]}} = \frac{\mathbf{a}^T \boldsymbol{\Sigma}_{xy} \mathbf{b}}{\sqrt{[\mathbf{a}^T \boldsymbol{\Sigma}_{xx} \mathbf{a}] [\mathbf{b}^T \boldsymbol{\Sigma}_{yy} \mathbf{b}]}}.\end{aligned}\quad (\text{C.30})$$

where $\boldsymbol{\Sigma}_{xy} = E(\mathbf{x}\mathbf{y}^T)$, $\boldsymbol{\Sigma}_{xx} = E(\mathbf{x}\mathbf{x}^T)$, and $\boldsymbol{\Sigma}_{yy} = E(\mathbf{y}\mathbf{y}^T)$ and the notation $E(\cdot)$ indicates the expectation. Here, $\boldsymbol{\Sigma}_{xx}$ and $\boldsymbol{\Sigma}_{yy}$ are symmetric matrices, and the relationship $\boldsymbol{\Sigma}_{xy}^T = \boldsymbol{\Sigma}_{yx}$ holds. The canonical correlation ρ_c is obtained by solving the following maximization problem:

$$\rho_c = \max_{\mathbf{a}, \mathbf{b}} \mathbf{a}^T \boldsymbol{\Sigma}_{xy} \mathbf{b} \quad \text{subject to} \quad \mathbf{a}^T \boldsymbol{\Sigma}_{xx} \mathbf{a} = 1 \quad \text{and} \quad \mathbf{b}^T \boldsymbol{\Sigma}_{yy} \mathbf{b} = 1. \quad (\text{C.31})$$

To solve this maximization problem, we first compute whitened vectors $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$, such that

$$\boldsymbol{\alpha} = \boldsymbol{\Sigma}_{xx}^{1/2} \mathbf{a}, \quad (\text{C.32})$$

$$\boldsymbol{\beta} = \boldsymbol{\Sigma}_{yy}^{1/2} \mathbf{b}. \quad (\text{C.33})$$

Using these vectors, the optimization problem in Eq. (C.31) is rewritten as

$$\rho_c = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \boldsymbol{\alpha}^T \boldsymbol{\Sigma}_{xx}^{-1/2} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1/2} \boldsymbol{\beta} \quad \text{subject to} \quad \boldsymbol{\alpha}^T \boldsymbol{\alpha} = 1 \quad \text{and} \quad \boldsymbol{\beta}^T \boldsymbol{\beta} = 1. \quad (\text{C.34})$$

The constrained maximization problem in Eq. (C.34) can be solved using the Lagrange multiplier method. Denoting the Lagrange multipliers v_1 and v_2 , the Lagrangian is defined as

$$\mathbb{L} = \boldsymbol{\alpha}^T \boldsymbol{\Sigma}_{xx}^{-1/2} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1/2} \boldsymbol{\beta} - \frac{v_1}{2} (\boldsymbol{\alpha}^T \boldsymbol{\alpha} - 1) - \frac{v_2}{2} (\boldsymbol{\beta}^T \boldsymbol{\beta} - 1). \quad (\text{C.35})$$

The optimum $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are obtained by maximizing \mathbb{L} with no constraints. The derivatives of \mathbb{L} with respect to $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are expressed as

$$\frac{\partial \mathbb{L}}{\partial \alpha} = \boldsymbol{\Pi}_{xy}\beta - v_1\alpha = 0, \quad (\text{C.36})$$

$$\frac{\partial \mathbb{L}}{\partial \beta} = \boldsymbol{\Pi}_{xy}^T\alpha - v_2\beta = 0, \quad (\text{C.37})$$

where $\boldsymbol{\Pi}_{xy} = \boldsymbol{\Sigma}_{xx}^{-1/2} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1/2}$. Using the constraint $\alpha^T \alpha = 1$, and left multiplying α^T to Eq. (C.36) gives

$$\alpha^T \boldsymbol{\Pi}_{xy}\beta = v_1.$$

Using the constraint $\beta^T \beta = 1$, and left multiplying β^T to Eq. (C.37) gives

$$\beta^T \boldsymbol{\Pi}_{xy}^T\alpha = v_2.$$

Since the left-hand sides of the two equations above are equal, we have $v_1 = v_2 = \rho_c$.

Also, using Eq. (C.37), we get

$$\beta = \frac{1}{\rho_c} \boldsymbol{\Pi}_{xy}^T\alpha,$$

and substituting the equation above into Eq. (C.36), we get

$$(\boldsymbol{\Pi}_{xy}\boldsymbol{\Pi}_{xy}^T)\alpha = \rho_c^2\alpha. \quad (\text{C.38})$$

This equation indicates that the squared canonical correlation ρ_c^2 is obtained as the eigenvalues of a matrix $\boldsymbol{\Pi}_{xy}\boldsymbol{\Pi}_{xy}^T$, and the corresponding eigenvector is the solution for the vector α . Note that $\boldsymbol{\Pi}_{xy}\boldsymbol{\Pi}_{xy}^T$ is a real symmetric matrix, the eigenvectors α and β are real-valued, and therefore α and β are real-valued because $\boldsymbol{\Sigma}_{xx}$ and $\boldsymbol{\Sigma}_{yy}$ are real-valued matrices.

Denoting the eigenvalues in Eq. (C.38) as μ_j where $j = 1, \dots, d$ and $d = \min\{p, q\}$, the canonical correlation between the two sets of random variables x_1, \dots, x_p and y_1, \dots, y_q is obtained as the largest eigenvalue μ_1 , which is the best overall measure of the association between x and y . However, other eigenvalues may provide complementary information on the linear relationship between those two sets of random variables. The mutual information described in Sect. C.3.2 is a measure that can take all the eigenvalues into account.

Also, it is worth mentioning that the matrices,

$$\boldsymbol{\Pi}_{xy}\boldsymbol{\Pi}_{xy}^T = \boldsymbol{\Sigma}_{xx}^{-1/2} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T \boldsymbol{\Sigma}_{xx}^{-1/2}$$

and

$$\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T \quad (\text{C.39})$$

have the same eigenvalues, according to Property No. 9 in Sect. C.8. Thus, the canonical squared correlation can be obtained as the largest eigenvalue of the matrix $\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T$.

C.3.2 Mutual Information

Next we introduce mutual information, which has a close relationship with the canonical correlation under the Gaussianity assumption. We also assume real-valued random vectors \mathbf{x} and \mathbf{y} , and their probability distribution as $p(\mathbf{x})$ and $p(\mathbf{y})$. The entropy is defined for \mathbf{x} and \mathbf{y} such that,

$$\mathcal{H}(\mathbf{x}) = - \int p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}, \quad (\text{C.40})$$

$$\mathcal{H}(\mathbf{y}) = - \int p(\mathbf{y}) \log p(\mathbf{y}) d\mathbf{y}. \quad (\text{C.41})$$

The entropy is a measure of uncertainty; $\mathcal{H}(\mathbf{x})$ represents the uncertainty when \mathbf{x} is unknown and $\mathcal{H}(\mathbf{y})$ represents the uncertainty when \mathbf{y} is unknown. The joint entropy is defined as

$$\mathcal{H}(\mathbf{x}, \mathbf{y}) = - \iint p(\mathbf{x}, \mathbf{y}) \log p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}, \quad (\text{C.42})$$

which represents the uncertainty when both \mathbf{x} and \mathbf{y} are unknown. When \mathbf{x} and \mathbf{y} are independent, we have the relationship

$$\begin{aligned} \mathcal{H}(\mathbf{x}, \mathbf{y}) &= - \iint p(\mathbf{x}, \mathbf{y}) \log p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= - \iint p(\mathbf{x})p(\mathbf{y}) (\log p(\mathbf{x}) + \log p(\mathbf{y})) d\mathbf{x} d\mathbf{y} = \mathcal{H}(\mathbf{x}) + \mathcal{H}(\mathbf{y}). \end{aligned} \quad (\text{C.43})$$

The conditional entropy is defined as

$$\mathcal{H}(\mathbf{x}|\mathbf{y}) = - \iint p(\mathbf{x}, \mathbf{y}) \log p(\mathbf{x}|\mathbf{y}) d\mathbf{x} d\mathbf{y}, \quad (\text{C.44})$$

which represents the uncertainty when \mathbf{x} is unknown, once \mathbf{y} is given. We then have the relationship,

$$\mathcal{H}(\mathbf{x}, \mathbf{y}) = \mathcal{H}(\mathbf{x}|\mathbf{y}) + \mathcal{H}(\mathbf{y}). \quad (\text{C.45})$$

The above indicates that the uncertainty when both \mathbf{x} and \mathbf{y} are unknown is equal to the uncertainty on \mathbf{x} when \mathbf{y} is given plus the uncertainty when \mathbf{y} is unknown.

The mutual information between \mathbf{x} and \mathbf{y} is defined as

$$\mathcal{I}(\mathbf{x}, \mathbf{y}) = \mathcal{H}(\mathbf{x}) + \mathcal{H}(\mathbf{y}) - \mathcal{H}(\mathbf{x}, \mathbf{y}). \quad (\text{C.46})$$

When \mathbf{x} and \mathbf{y} are independent, we have $\mathcal{I}(\mathbf{x}, \mathbf{y}) = 0$ according to Eq. (C.43). Let us define \mathbf{z} as $\mathbf{z} = [\mathbf{x}^T, \mathbf{y}^T]^T$. Assuming the Gaussian processes for \mathbf{x} and \mathbf{y} , the mutual information is expressed as

$$\begin{aligned} \mathcal{I}(\mathbf{x}, \mathbf{y}) &= \mathcal{H}(\mathbf{x}) + \mathcal{H}(\mathbf{y}) - \mathcal{H}(\mathbf{x}, \mathbf{y}) \\ &= \frac{1}{2} \log |\boldsymbol{\Sigma}_{xx}| + \frac{1}{2} \log |\boldsymbol{\Sigma}_{yy}| - \frac{1}{2} \log |\boldsymbol{\Sigma}_{zz}|. \end{aligned} \quad (\text{C.47})$$

Here, $|\boldsymbol{\Sigma}_{zz}|$ is rewritten as

$$|\boldsymbol{\Sigma}_{zz}| = \left| \begin{array}{cc} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{xy}^T & \boldsymbol{\Sigma}_{yy} \end{array} \right| = |\boldsymbol{\Sigma}_{yy}| \left| \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T \right|, \quad (\text{C.48})$$

where the determinant identity in Eq. (C.94) is used. Substituting Eq. (C.48) into (C.47), we have

$$\begin{aligned} \mathcal{I}(\mathbf{x}, \mathbf{y}) &= \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_{xx}|}{\left| \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T \right|} \\ &= \frac{1}{2} \log \frac{1}{\left| \mathbf{I} - \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T \right|}. \end{aligned} \quad (\text{C.49})$$

Let us define the eigenvalues of $\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T$ as γ_j where $j = 1, \dots, d$ and $d = \min\{p, q\}$.⁶ Using these eigenvalues, we can derive

$$\begin{aligned} \mathcal{I}(\mathbf{x}, \mathbf{y}) &= \frac{1}{2} \log \frac{1}{\left| \mathbf{I} - \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T \right|} \\ &= \frac{1}{2} \log \frac{1}{\prod_{j=1}^d (1 - \gamma_j)} = \frac{1}{2} \sum_{j=1}^d \log \frac{1}{1 - \gamma_j}. \end{aligned} \quad (\text{C.50})$$

Note that γ_1 is equal to the canonical squared correlation ρ_c^2 , according to the arguments in Sect. C.3.1.

When the random vectors \mathbf{x} and \mathbf{y} are complex-valued, the mutual information for the complex random vectors is expressed as

$$\mathcal{I}(\mathbf{x}, \mathbf{y}) = \mathcal{H}(\mathbf{x}) + \mathcal{H}(\mathbf{y}) - \mathcal{H}(\mathbf{x}, \mathbf{y}) = \log |\boldsymbol{\Sigma}_{xx}| + \log |\boldsymbol{\Sigma}_{yy}| - \log |\boldsymbol{\Sigma}_{zz}|. \quad (\text{C.51})$$

⁶ Here, remember that p and q are the sizes of the column vector \mathbf{x} and \mathbf{y}

Here, $\boldsymbol{\Sigma}_{xx}$, $\boldsymbol{\Sigma}_{yy}$, and $\boldsymbol{\Sigma}_{zz}$ are covariance matrices of the complex Gaussian distribution. Using exactly the same derivation, we can obtain

$$\mathcal{I}(\mathbf{x}, \mathbf{y}) = \log \frac{1}{|\mathbf{I} - \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T|} = \sum_{j=1}^d \log \frac{1}{1 - \gamma_j}, \quad (\text{C.52})$$

where γ_j is the j th eigenvalue of $\boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{xy}^T$.

C.3.3 Covariance of Residual Signal and Conditional Entropy

Let us next consider the regression problem:

$$\mathbf{y} = \mathbf{Ax} + \mathbf{e}, \quad (\text{C.53})$$

where the vector \mathbf{e} represents the residual of this regression. The $p \times q$ coefficient matrix \mathbf{A} can be estimated using the least-squares fit,

$$\widehat{\mathbf{A}} = \underset{\mathbf{A}}{\operatorname{argmin}} E \left[\|\mathbf{e}\|^2 \right] = \underset{\mathbf{A}}{\operatorname{argmin}} E \left[\|\mathbf{y} - \mathbf{Ax}\|^2 \right]. \quad (\text{C.54})$$

Here, $E \left[\|\mathbf{e}\|^2 \right]$ is expressed as

$$E \left[\|\mathbf{e}\|^2 \right] = E \left[\mathbf{y}\mathbf{y}^T - \mathbf{x}^T \mathbf{A}^T \mathbf{y} - \mathbf{y}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} \right].$$

Differentiating $E \left[\|\mathbf{e}\|^2 \right]$ with respect to \mathbf{A} gives

$$\frac{\partial}{\partial \mathbf{A}} E \left[\|\mathbf{e}\|^2 \right] = -2E(\mathbf{y}\mathbf{x}^T) + 2AE(\mathbf{x}\mathbf{x}^T).$$

Setting this derivative to zero, we obtain

$$\widehat{\mathbf{A}} = E(\mathbf{y}\mathbf{x}^T)E(\mathbf{x}\mathbf{x}^T)^{-1} = \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}, \quad (\text{C.55})$$

where $\boldsymbol{\Sigma}_{yx} = E(\mathbf{y}\mathbf{x}^T)$ and $\boldsymbol{\Sigma}_{xx} = E(\mathbf{x}\mathbf{x}^T)$. The covariance matrix of the residual \mathbf{e} is defined as

$$\begin{aligned}
\boldsymbol{\Sigma}_{ee} &= E \left[(\mathbf{y} - \widehat{\mathbf{A}}\mathbf{x}) (\mathbf{y} - \widehat{\mathbf{A}}\mathbf{x})^T \right] \\
&= E \left[\mathbf{y}\mathbf{y}^T - \mathbf{y}\mathbf{x}^T \widehat{\mathbf{A}}^T - \widehat{\mathbf{A}}\mathbf{x}\mathbf{y}^T + \widehat{\mathbf{A}}\mathbf{x}\mathbf{x}^T \widehat{\mathbf{A}}^T \right] \\
&= E(\mathbf{y}\mathbf{y}^T) - E(\mathbf{y}\mathbf{x}^T)\widehat{\mathbf{A}}^T - \widehat{\mathbf{A}}E(\mathbf{x}\mathbf{y}^T) + \widehat{\mathbf{A}}E(\mathbf{x}\mathbf{x}^T)\widehat{\mathbf{A}}^T \\
&= \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx}\widehat{\mathbf{A}}^T - \widehat{\mathbf{A}}\boldsymbol{\Sigma}_{xy} + \widehat{\mathbf{A}}\boldsymbol{\Sigma}_{xx}\widehat{\mathbf{A}}^T,
\end{aligned} \tag{C.56}$$

where $\boldsymbol{\Sigma}_{yy} = E(\mathbf{y}\mathbf{y}^T)$ and $\boldsymbol{\Sigma}_{xy} = E(\mathbf{x}\mathbf{y}^T)$. Substituting Eq. (C.55) into the equation above gives,

$$\boldsymbol{\Sigma}_{ee} = \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx}\boldsymbol{\Sigma}_{xx}^{-1}\boldsymbol{\Sigma}_{yx}^T. \tag{C.57}$$

Let us assume that the random variables \mathbf{x} and \mathbf{y} follow the Gaussian distribution. According to Sect. C.1, we have the relationship,

$$\mathcal{H}(\mathbf{x}) = \frac{1}{2} \log |\boldsymbol{\Sigma}_{xx}| \quad \text{and} \quad \mathcal{H}(\mathbf{y}) = \frac{1}{2} \log |\boldsymbol{\Sigma}_{yy}|, \tag{C.58}$$

where we ignore constants that are not related to the current arguments. We also have

$$\begin{aligned}
\mathcal{H}(\mathbf{x}, \mathbf{y}) &= \frac{1}{2} \log \left| E \left[\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \begin{bmatrix} \mathbf{x}^T, \mathbf{y}^T \end{bmatrix} \right] \right| \\
&= \frac{1}{2} \log \left| \begin{array}{cc} E(\mathbf{x}\mathbf{x}^T) & E(\mathbf{x}\mathbf{y}^T) \\ E(\mathbf{y}\mathbf{x}^T) & E(\mathbf{y}\mathbf{y}^T) \end{array} \right| = \frac{1}{2} \log \left| \begin{array}{cc} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{yx}^T \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{array} \right|. \tag{C.59}
\end{aligned}$$

Thus, the conditional entropy is expressed as

$$\mathcal{H}(\mathbf{y}|\mathbf{x}) = \mathcal{H}(\mathbf{x}, \mathbf{y}) - \mathcal{H}(\mathbf{x}) = \frac{1}{2} \log \left| \begin{array}{cc} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{yx}^T \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{array} \right| - \frac{1}{2} \log |\boldsymbol{\Sigma}_{xx}|. \tag{C.60}$$

Using the determinant identity in Eq. (C.94), we finally obtain the formula to compute the conditional entropy

$$\begin{aligned}
\mathcal{H}(\mathbf{y}|\mathbf{x}) &= \frac{1}{2} \log \left| \begin{array}{cc} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{yx}^T \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{array} \right| - \frac{1}{2} \log |\boldsymbol{\Sigma}_{xx}| \\
&= \frac{1}{2} \log \left| \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx}\boldsymbol{\Sigma}_{xx}^{-1}\boldsymbol{\Sigma}_{yx}^T \right| = \frac{1}{2} \log |\boldsymbol{\Sigma}_{ee}|. \tag{C.61}
\end{aligned}$$

The conditional entropy is expressed by the covariance of the residual signal e obtained by regressing \mathbf{y} with \mathbf{x} .

C.4 Definitions of Several Vector Norms

The p th order norm ($p \geq 0$) of an N -dimensional vector, $\mathbf{x} = [x_1, \dots, x_N]^T$ is defined such that

$$\|\mathbf{x}\|_p = (|x_1|^p + |x_2|^p + \dots + |x_N|^p)^{1/p}. \quad (\text{C.62})$$

When $p = 2$, the norm is called the L_2 -norm, which is equal to the conventional Euclidean norm,

$$\|\mathbf{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_N^2}. \quad (\text{C.63})$$

The L_2 -norm is usually denoted $\|\mathbf{x}\|$. When $p = 1$, the norm is expressed as

$$\|\mathbf{x}\|_1 = |x_1| + |x_2| + \dots + |x_N|, \quad (\text{C.64})$$

which is called the L_1 -norm. When $p = 0$, we can define $\mathcal{T}(x)$, which is called the indicator function, such that

$$\mathcal{T}(x) = \begin{cases} 0 & x = 0 \\ 1 & x \neq 0, \end{cases} \quad (\text{C.65})$$

with the norm expressed as

$$\|\mathbf{x}\|_0 = \sum_{i=1}^N \mathcal{T}(x_i). \quad (\text{C.66})$$

Namely, $\|\mathbf{x}\|_0$ is equal to the number of non-zero components of \mathbf{x} . When $p = \infty$, the norm becomes

$$\|\mathbf{x}\|_\infty = \max\{x_1, x_2, \dots, x_N\}. \quad (\text{C.67})$$

Namely, the $\|\mathbf{x}\|_\infty$ is equal to the maximum component of \mathbf{x} .

C.5 Derivative of Functionals

This section provides a brief explanation of the derivative of functionals. A conventional function relates input variables to output variables. A functional relates input functions to output variables. A simple example of a functional is:

$$I = \int_{\alpha}^{\beta} f(x) dx, \quad (\text{C.68})$$

where $f(x)$ is defined in $[\alpha, \beta]$. This I is a functional of the function $f(x)$. The functional is often denoted as $I[f(x)]$.

Let us consider a function $F(x)$ that contains another function $f(x)$. One example is $F(x) = f(x)^2 + 2x$ where $F(x)$ consists of $f(x)$ and x . The integral of $F(x)$,

$$I[f(x)] = \int_{\alpha}^{\beta} F(x) dx, \quad (\text{C.69})$$

is a functional of $f(x)$. Let us define a small change in $f(x)$ as $\delta f(x)$, and a change in $I[f(x)]$ due to $\delta f(x)$ as $\delta I[f(x)]$. For simplicity, $\delta I[f(x)]$ is denoted δI in the following arguments. To derive the relationship between δI and $\delta f(x)$, the region $[\alpha, \beta]$ is divided into small regions Δx . A contribution from the j th small region onto δI , δI_j , is equal to

$$\delta I_j = A_j \delta f(x_j) \Delta x. \quad (\text{C.70})$$

Here we assume that δI_j is proportional to $\delta f(x_j)$ and A_j is a proportional constant. Then, δI is derived as a sum of contributions from all small regions such that

$$\delta I = \sum_j A_j \delta f(x_j) \Delta x. \quad (\text{C.71})$$

Accordingly, when $\Delta x \rightarrow 0$, the relationship between $\delta f(x)$ and δI becomes

$$\delta I = \int_{\alpha}^{\beta} A(x) \delta f(x) dx. \quad (\text{C.72})$$

In Eq. (C.72), $A(x)$ is defined as the derivative of the functional $I[f(x)]$ with respect to $f(x)$, and it is denoted $\delta I/\delta f$.

Let us compute this derivative of a functional. When a functional is given in Eq. (C.68), i.e., $F(x) = f(x)$, the amount of change δI due to the change from $f(x)$ to $f(x) + \delta f(x)$ is expressed as

$$\delta I = \int_{\alpha}^{\beta} [f(x) + \delta f(x)] dx - \int_{\alpha}^{\beta} f(x) dx = \int_{\alpha}^{\beta} \delta f(x) dx. \quad (\text{C.73})$$

Therefore, we obtain $A(x) = \delta I/\delta f = 1$. In the general case where $F(x)$ consists of $f(x)$, the relationship

$$\delta I = \int_{\alpha}^{\beta} [F(x) + \delta F(x)] dx - \int_{\alpha}^{\beta} F(x) dx = \int_{\alpha}^{\beta} \delta F(x) dx \quad (\text{C.74})$$

holds. Using the relationship

$$\delta F(x) = \frac{\partial F}{\partial f} \delta f(x),$$

we have

$$\delta I = \int_{\alpha}^{\beta} \delta F(x) dx = \int_{\alpha}^{\beta} \frac{\partial F}{\partial f} \delta f(x) dx, \quad (\text{C.75})$$

and thus,

$$\frac{\delta I}{\delta f} = \frac{\partial F}{\partial f}. \quad (\text{C.76})$$

As an example of the use of Eq. (C.76), let us derive the derivative of the functional $\mathbb{L}[q(\mathbf{x})]$ in Eq. (B.53). Ignoring the last term, which does not contain $q(\mathbf{x})$, Eq. (B.53) is rewritten as

$$\begin{aligned} \mathbb{L}[q(\mathbf{x})] &= \int_{-\infty}^{\infty} d\mathbf{x} q(\mathbf{x}) \log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) \\ &\quad - \int_{-\infty}^{\infty} d\mathbf{x} q(\mathbf{x}) \log q(\mathbf{x}) + \gamma \int_{-\infty}^{\infty} q(\mathbf{x}) d\mathbf{x}. \end{aligned} \quad (\text{C.77})$$

In the first term, since $F(\mathbf{x}) = q(\mathbf{x}) \log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})$, we have

$$\frac{\partial F}{\partial q} = \log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}). \quad (\text{C.78})$$

In the second term, since $F(\mathbf{x}) = q(\mathbf{x}) \log q(\mathbf{x})$, we have

$$\frac{\partial F}{\partial q} = \frac{1}{\partial q(\mathbf{x})} [q(\mathbf{x}) \log q(\mathbf{x})] = \log q(\mathbf{x}) + q(\mathbf{x}) \frac{1}{q(\mathbf{x})} = \log q(\mathbf{x}) + 1. \quad (\text{C.79})$$

In the third term, since $F(\mathbf{x}) = \gamma q(\mathbf{x})$, we have

$$\frac{\partial F}{\partial q} = \gamma. \quad (\text{C.80})$$

Therefore, we have

$$\frac{\partial \mathbb{L}[q(\mathbf{x})]}{\partial q(\mathbf{x})} = \log p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) - \log q(\mathbf{x}) - 1 + \gamma, \quad (\text{C.81})$$

which is equal to Eq. (B.54).

C.6 Vector and Matrix Derivatives

Differentiating a scalar F with a column vector \mathbf{x} is defined as creating a column vector whose j th element is equal to $\partial F / \partial x_j$. Assuming that \mathbf{a} is a column vector and A is a matrix, The following relationships hold,

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}, \quad (\text{C.82})$$

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}, \quad (\text{C.83})$$

$$\frac{\partial \text{tr}(\mathbf{x} \mathbf{a}^T)}{\partial \mathbf{x}} = \frac{\partial \text{tr}(\mathbf{a} \mathbf{x}^T)}{\partial \mathbf{x}} = \mathbf{a}. \quad (\text{C.84})$$

Let us denote the (i, j) th element of a matrix \mathbf{A} as $A_{i,j}$. Differentiating a scalar F with a matrix \mathbf{A} is defined as creating a matrix whose (i, j) th element is equal to $\partial F / \partial A_{i,j}$. Representative identities are the following, where \mathbf{x} and \mathbf{y} are column vectors and \mathbf{A} and \mathbf{B} are matrices.

$$\frac{\partial \text{tr}(\mathbf{A})}{\partial \mathbf{A}} = \mathbf{I}, \quad (\text{C.85})$$

$$\frac{\partial \text{tr}(\mathbf{AB})}{\partial \mathbf{A}} = \mathbf{B}^T, \quad (\text{C.86})$$

$$\frac{\partial \text{tr}(\mathbf{A}^T \mathbf{B})}{\partial \mathbf{A}} = \mathbf{B}, \quad (\text{C.87})$$

$$\frac{\partial \text{tr}(\mathbf{ABB}^T)}{\partial \mathbf{A}} = \mathbf{A}(\mathbf{B} + \mathbf{B}^T), \quad (\text{C.88})$$

$$\frac{\partial \mathbf{x}^T \mathbf{Ay}}{\partial \mathbf{A}} = \mathbf{xy}^T, \quad (\text{C.89})$$

$$\frac{\partial \log |\mathbf{A}|}{\partial \mathbf{A}} = (\mathbf{A}^{-1})^T. \quad (\text{C.90})$$

C.7 Several Formulae for Matrix Computations Used in this Book

The following are representative formulae for the matrix inversion.

$$(\mathbf{A} + \mathbf{BD}^{-1}\mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{CA}^{-1}\mathbf{B})^{-1}\mathbf{CA}^{-1}, \quad (\text{C.91})$$

$$(\mathbf{A}^{-1} + \mathbf{B}^T \mathbf{C}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{C}^{-1} = \mathbf{AB}^T (\mathbf{BAB}^T + \mathbf{C})^{-1}. \quad (\text{C.92})$$

Also, we have

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{M} & -\mathbf{MBD}^{-1} \\ -\mathbf{D}^{-1}\mathbf{CM} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{CM}\mathbf{BD}^{-1} \end{bmatrix}, \quad (\text{C.93})$$

where

$$\mathbf{M} = (\mathbf{A} - \mathbf{BD}^{-1}\mathbf{C})^{-1}.$$

Regarding the matrix determinant, we have the following identity,

$$\begin{vmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{vmatrix} = |\mathbf{A}| |\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B}| = |\mathbf{D}| |\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}|. \quad (\text{C.94})$$

When a matrix \mathbf{A} is an invertible $(n \times n)$ matrix, and \mathbf{B} and \mathbf{C} are $(n \times m)$ matrices, the following matrix determinant lemma holds:

$$|\mathbf{A}| |\mathbf{I} + \mathbf{B}^T \mathbf{A}^{-1} \mathbf{C}| = |\mathbf{A} + \mathbf{C}\mathbf{B}^T|. \quad (\text{C.95})$$

C.8 Properties of Eigenvalues

Representative properties of eigenvalues that may be used in this book are listed below.

1. If \mathbf{A} is a Hermitian matrix and positive definite, all eigenvalues are real and greater than 0.
2. If \mathbf{A} is a Hermitian matrix and positive semidefinite, all eigenvalues are real and greater than or equal to 0.
3. If \mathbf{A} is a real symmetric matrix and positive definite, all eigenvalues are real and greater than 0. Eigenvectors are also real.
4. If \mathbf{A} is a real symmetric matrix and positive semidefinite, all eigenvalues are real and greater than or equal to 0. Eigenvectors are also real.
5. If \mathbf{A} and \mathbf{B} are square matrices and \mathbf{B} is nonsingular, eigenvalues of \mathbf{A} are also eigenvalues of $\mathbf{B}^{-1}\mathbf{AB}$.
6. If \mathbf{A} and \mathbf{B} are square matrices and \mathbf{B} is unitary, eigenvalues of \mathbf{A} are also eigenvalues of $\mathbf{B}^H\mathbf{AB}$.
7. If \mathbf{A} and \mathbf{B} are square matrices and \mathbf{B} is orthogonal, eigenvalues of \mathbf{A} are also eigenvalues of $\mathbf{B}^T\mathbf{AB}$.
8. If \mathbf{A} and \mathbf{B} are square matrices and \mathbf{B} is positive definite, eigenvalues of \mathbf{BA} are also eigenvalues of $\mathbf{B}^{1/2}\mathbf{AB}^{1/2}$.
9. If \mathbf{A} and \mathbf{B} are square matrices and \mathbf{B} is positive definite, eigenvalues of $\mathbf{B}^{-1}\mathbf{A}$ are also eigenvalues of $\mathbf{B}^{-1/2}\mathbf{AB}^{-1/2}$.
10. Let us assume that \mathbf{A} is an $(m \times n)$ matrix, \mathbf{B} is an $(n \times m)$ matrix, and $n \geq m$. If $\lambda_1, \dots, \lambda_m$ are eigenvalues of \mathbf{AB} , $\lambda_1, \dots, \lambda_m, 0, \dots, 0$ are eigenvalues of \mathbf{BA} .

C.9 Rayleigh-Ritz Formula

This section provides a proof of the Rayleigh-Ritz formula, which is according to [2]. We define \mathbf{A} and \mathbf{B} as positive definite matrices of the same dimension. We introduce the following notations and use them throughout the book:

- The minimum and maximum eigenvalues of a matrix \mathbf{A} are denoted $\mathcal{S}_{\min}\{\mathbf{A}\}$ and $\mathcal{S}_{\max}\{\mathbf{A}\}$.
- The eigenvectors corresponding to the minimum and maximum eigenvalues of a matrix \mathbf{A} are denoted $\vartheta_{\min}\{\mathbf{A}\}$ and $\vartheta_{\max}\{\mathbf{A}\}$.
- The minimum and maximum generalized eigenvalues of a matrix \mathbf{A} with a metric \mathbf{B} are denoted $\mathcal{S}_{\min}\{\mathbf{A}, \mathbf{B}\}$ and $\mathcal{S}_{\max}\{\mathbf{A}, \mathbf{B}\}$, and the corresponding eigenvectors are denoted $\vartheta_{\min}\{\mathbf{A}, \mathbf{B}\}$ and $\vartheta_{\max}\{\mathbf{A}, \mathbf{B}\}$.

Here, if the matrix \mathbf{B} is nonsingular, the following relationships hold:

$$\begin{aligned}\mathcal{S}_{\max}\{\mathbf{A}, \mathbf{B}\} &= \mathcal{S}_{\max}\{\mathbf{B}^{-1}\mathbf{A}\}, \\ \vartheta_{\max}\{\mathbf{A}, \mathbf{B}\} &= \vartheta_{\max}\{\mathbf{B}^{-1}\mathbf{A}\}, \\ \mathcal{S}_{\min}\{\mathbf{A}, \mathbf{B}\} &= \mathcal{S}_{\min}\{\mathbf{B}^{-1}\mathbf{A}\}, \\ \vartheta_{\min}\{\mathbf{A}, \mathbf{B}\} &= \vartheta_{\min}\{\mathbf{B}^{-1}\mathbf{A}\}.\end{aligned}$$

Using \mathbf{x} to denote a column vector with its dimension commensurate with the size of the matrices, this appendix shows that

$$\max_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{x}} = \mathcal{S}_{\max}\{\mathbf{A}, \mathbf{B}\}, \quad (\text{C.96})$$

and

$$\operatorname{argmax}_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{x}} = \vartheta_{\max}\{\mathbf{A}, \mathbf{B}\}. \quad (\text{C.97})$$

Since the value of the ratio $(\mathbf{x}^T \mathbf{A} \mathbf{x}) / (\mathbf{x}^T \mathbf{B} \mathbf{x})$ is not affected by the norm of \mathbf{x} , we set the norm of \mathbf{x} so as to satisfy the relationship $\mathbf{x}^T \mathbf{B} \mathbf{x} = 1$. Then, the maximization problem in Eq. (C.96) is rewritten as

$$\max_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} \quad \text{subject to} \quad \mathbf{x}^T \mathbf{B} \mathbf{x} = 1. \quad (\text{C.98})$$

We change this constrained maximization problem to an unconstrained maximization problem by introducing the Lagrange multiplier κ . We define the Lagrangian $\mathbb{L}(\mathbf{x}, \kappa)$ such that

$$\mathbb{L}(\mathbf{x}, \kappa) = \mathbf{x}^T \mathbf{A} \mathbf{x} - \kappa (\mathbf{x}^T \mathbf{B} \mathbf{x} - 1). \quad (\text{C.99})$$

The maximization in Eq. (C.98) is equivalent to maximizing $\mathbb{L}(\mathbf{x}, \kappa)$ with no constraints.

To obtain the maximum of $\mathbb{L}(\mathbf{x}, \kappa)$, we calculate the derivatives

$$\frac{\partial \mathbb{L}(\mathbf{x}, \kappa)}{\partial \mathbf{x}} = 2(\mathbf{A} \mathbf{x} - \kappa \mathbf{B} \mathbf{x}), \quad (\text{C.100})$$

$$\frac{\partial \mathbb{L}(\mathbf{x}, \kappa)}{\partial \kappa} = -(\mathbf{x}^T \mathbf{B} \mathbf{x} - 1). \quad (\text{C.101})$$

By setting these derivatives to zero, we can derive the relationships, $\mathbf{Ax} = \kappa \mathbf{Bx}$ and $\kappa = \mathbf{x}^T \mathbf{A} \mathbf{x}$. Therefore, the maximum value of $\mathbf{x}^T \mathbf{A} \mathbf{x}$ is equal to the maximum eigenvalue of $\mathbf{Ax} = \kappa \mathbf{Bx}$, and the \mathbf{x} that attains this maximum value is equal to the eigenvector corresponding to this maximum eigenvalue. Namely, we have

$$\max_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{x}} = \mathcal{S}_{\max}\{\mathbf{A}, \mathbf{B}\} = \mathcal{S}_{\max}\{\mathbf{B}^{-1} \mathbf{A}\},$$

and

$$\operatorname{argmax}_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{x}} = \vartheta_{\max}\{\mathbf{A}, \mathbf{B}\} = \vartheta_{\max}\{\mathbf{B}^{-1} \mathbf{A}\}.$$

Using exactly the same derivation, it is easy to show that

$$\min_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{x}} = \mathcal{S}_{\min}\{\mathbf{A}, \mathbf{B}\}, \quad (\text{C.102})$$

and

$$\operatorname{argmin}_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{x}} = \vartheta_{\min}\{\mathbf{A}, \mathbf{B}\}. \quad (\text{C.103})$$

References

1. F. D. Neeser and J. L. Massey, “Proper complex random processes with applications to information theory,” *IEEE Transactions on Information Theory*, vol. 39, pp. 1293–1302, 1993.
2. F. R. Gantmacher, *The Theory of Matrices*. New York, NY: Chelsea Publishing Company, 1960.

Index

Symbols

- L_0 -norm minimization, 20
- L_1 -norm, 259
- L_1 -norm regularization, 26
- L_1 -regularized minimum-norm solution, 19
- L_2 -norm, 259
- L_2 -norm regularized minimum-norm method, 25
- L_2 -regularized minimum-norm solution, 17
- L_p -norm regularization, 22

A

- Action potential, 215
- Adaptive beamformer, 29
- Alternative cost function, 59
- Amplitude–amplitude coupling (AAC), 200
- Amplitude–phase diagram, 202, 206, 211
- Analytic signal, 202
- Array measurement, 9
- Array-gain constraint, 31, 35
- Augmented factor vector, 102
- Automatic relevance determination (ARD), 123
- Auxiliary variables , 60
- Average data likelihood, 237
- Average log likelihood, 77
- Axon, 215

B

- Bayes’ rule, 24, 231
- Bayesian beamformer, 36
- Bayesian factor analysis, 75
- Bayesian inference, 231
- Bayesian minimum-norm method, 25
- Bayesian modeling framework, 121

Beam response, 47

- Beamformer, 29
- Bioelectromagnetic forward modeling, 215
- Biot-Savart law, 225
- Bounded conductor with piecewise-constant conductivity, 226
- Brain oscillation, 199
- Brain rhythm, 199

C

- Canonical correlation, 252
- Canonical imaginary coherence, 154
- Canonical magnitude coherence, 151
- Capacitive effect, 223
- Champagne algorithm, 51
- Charge conservation law, 219
- Circular, 250
- Coherence, 183
- Coherence of the MVAR process, 173
- Complete data likelihood, 237
- Complex Gaussian distribution, 249
- Concave function, 59
- Concentration gradient, 217
- Conditional entropy, 258
- Convexity-based algorithm, 59
- Convexity-based update, 126
- Corrected imaginary coherence, 145
- Cortical anatomy, 216
- Cortico-cortical connection, 216
- Covariance component, 122, 133
- Covariance of residual, 257
- Cross-frequency coupling, 200
- Cross-location PAC, 201, 209

D

- Data evidence, 54, 236
- Data likelihood, 231
- Data vector, 9
- Derivative of functional, 259
- Determinant identity, 263
- Determinant lemma, 263
- Diagonal-loading, 35
- Dielectric permittivity, 223
- Dipole model, 225
- Directed transfer function (DTF), 182
- Dual-state beamformer, 42

E

- EM algorithm, 25, 77, 191
- EM update, 58, 126
- Empirical Bayesian, 123
- Empirical Bayesian schema, 51
- Entropy, 248
- Envelope coherence, 160
- Envelope correlation, 159
- ERD/ERS activities, 209
- Event-related desynchronization (ERD), 42
- Event-related power change, 209
- Event-related spectral power change, 42
- Event-related synchronization (ERS), 42

F

- F-ratio image, 44
- Factor analysis model, 75
- Five-dimensional brain imaging, 44
- FOCUSS, 128
- Free energy, 88, 92, 105, 242

G

- Gamma distribution, 27
- Gamma-band oscillation, 199
- Gauge Transformations, 220
- Gauss theorem, 225, 227
- Gaussian prior, 51
- Gaussian scale mixtures, 122
- Generative model, 122
- Geselowitz formula, 227
- Geweke measure, 178
- Global interaction measure (GIM), 157
- Granger causality for a bivariate process, 174
- Granger-causality measures, 171
- Gulrajani gauge, 220

H

- Head tissue conductivity, 224
- High-frequency oscillator, 200
- Hilbert transform, 202
- Hyperparameter, 69, 236
- Hyperparameter MAP, 123

I

- Ideal gas constant, 217
- Imaginary coherence, 139
- Impressed current, 219
- Inhibitory synaptic connections, 216
- Instantaneous amplitude, 202
- Instantaneous dependence, 178
- Instantaneous interaction, 142
- Instantaneous phase, 202
- Intensity bias, 145
- Interference mixing matrix, 96
- Interference suppression, 75
- Intracortical connection, 216
- Inverse-Gamma, 124

K

- Kullback–Leibler distance, 203, 244
- Kullback-Leibler divergence, 132

L

- Lagrange multiplier, 15, 30, 242
- Lagrangian, 242, 264
- Laplace approximation, 132
- Laplace distribution, 26
- Lead-field matrix, 10
- Lead-field vector, 10
- Leakage effects, 143
- Linearly constrained minimumvariance (LCMV) beamformer, 39

Local PAC, 201

Local PAC analysis, 206

Log-likelihood function, 14

Low-frequency oscillator, 200

M

- MacKay update, 58, 126
- Magnetic scalar potential, 228
- Magnetic vector potential, 219
- Marginal distribution, 235
- Marginal likelihood, 54, 79, 236
- Maximum a posteriori (MAP) estimate, 232
- Maximum likelihood principle, 13
- Maximum statistics, 161

Maxwell's equations, 218
MCE, 128
Mean imaginary coherence (MIC) mapping, 162
Mean-Field approximation, 131
Method of least-squares, 14
Minimum mean squared error (MMSE) estimate, 233
Minimum-norm filter, 44
Minimum-norm solution, 15
Minimum-variance beamformer, 30
Mixing matrix, 76
Model data covariance, 236
Model order, 75
Modulation index (MI), 203
Modulation index analysis, 205
Multi-dimensional Gaussian distribution, 247
Multiple comparison problem, 161
Multivariate Granger causality, 175
Multivariate interaction measure (MIM), 156
Multivariate vector auto-regressive (MVAR) process, 171, 190
Mutual information, 146, 188, 255

N
Narrow-band beamformer, 42
Neuronal current, 215
Neurotransmitter, 216
Noise precision, 239
Non-informative prior, 232
Non-zero-lag correlation, 142
Nonadaptive spatial filter, 44

O
Objective function, 22
Ohm's law, 217
Ohmic current, 219

P
Parseval's theorem, 142
Partial coherence, 173
Partial directed coherence, 184
Partitioned factor analysis (PFA), 96
Penalized likelihood method, 127
Phase of firing, 199
Phase–amplitude coupling (PAC), 200
Phase-informed TF map, 206, 212
Phase-informed time-frequency map, 204
Phase–phase coupling (PPC), 200

Posterior probability, 231
Postsynaptic terminal, 216
Potential equation, 221
Preferred phases, 212
Prior precision, 238
Prior probability distribution, 231
Proper, 249
Pyramidal cells, 216

R

Rayleigh-Ritz formula, 263
Recursive null-Steering (RENS) beamformer, 47
Residual canonical coherence, 157
Residual coherence, 148
Residual envelope correlation, 160
Resolution matrix, 16
Restricted maximum likelihood (ReML), 125

S

Saketini algorithm, 101
Same-frequency coupling (SFC), 200
Sarvas formula, 228
Scalar potential, 220
Scalar-type adaptive beamformer, 36
Seed point, 140
Seed voxel, 141
Semi-Bayesian derivation, 33
Semi-Bayesian formulation, 40
Singular value decomposition, 16
Singular-value spectrum, 17
sLORETA filter, 46
Somatic voltage, 215
Source current, 219
Source MAP, 127
Source space, 11
Source space causality analysis, 171
Source vector, 10
Source-space connectivity analysis, 139
Sparse Bayesian learning, 26, 51
Sparse Bayesian (Champagne) algorithm, 191
Sparse solution, 20, 22
Sparsity, 63
Spatial filter, 29
Spectral Granger causality, 178
Spherically-symmetric homogeneous conductor, 227
Student t -distribution, 27
Surrogate data method, 161, 206
Synapses, 215

Synaptic connection, 216

T

Target location, 140

Target voxel, 141

Temporal encoding, 199

Time-Domain Granger causality, 174

Total interdependence, 177

Transfer entropy, 185

Transmembrane voltage, 215, 217

Type-II maximum likelihood, 125

U

Unbounded homogeneous medium, 224

Unit-gain constraint, 30, 34

Unit-noise-gain constraint, 32

V

Variational Bayes EM algorithm (VBEM),
84

Variational Bayes Factor Analysis (VBFA),
82

Variational Bayesian approximation, 131

Variational Bayesian inference, 241

VBEM algorithm, 97, 103

Vector norm, 259

Vector-type beamformer, 38

Virtual sensor, 29

Voxel, 11

Voxel-by-voxel statistical threshold, 161

W

Weight-normalized minimum-norm filter,
46

Weight-normalized minimum-variance
beamformer, 33

Y

Yule-Walker equation, 191

Z

Zero-time-lag correlation, 142