



Carmen Sippel

Turbo and Message Passing Approaches for Iterative Signal Recovery in Compressed Sensing





Carmen Sippel

Turbo and Message Passing Approaches for Iterative Signal Recovery in Compressed Sensing

Bibliografische Information der Deutschen Nationalbibliothek

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über http://dnb.dnb.de abrufbar.

Dissertation, Universität Ulm, Fakultät für Ingenieurwissenschaften, Informatik und Psychologie, 2023

Impressum

Universität Ulm Institut für Nachrichtentechnik Prof. Dr.-Ing. Robert Fischer Albert-Einstein-Allee 43 89081 Ulm

Eine Übersicht über alle Bände der Schriftenreihe finden Sie unter https://nt.uni-ulm.de/schriften

Diese Veröffentlichung ist im Internet auf dem Open Access Repositorium der Universität Ulm (https://oparu.uni-ulm.de) verfügbar und dort unter der Lizenz "A" publiziert. Details zur Lizenz sind unter http://www.uni-ulm.de/index.php?id=50392 zu finden.

Institut für Nachrichtentechnik der Universität Ulm 2023

e-ISBN 978-3-948303-38-9



Turbo and Message Passing Approaches for Iterative Signal Recovery in Compressed Sensing

DISSERTATION

zur Erlangung des akademischen Grades eines

DOKTOR-INGENIEURS

(Dr.-Ing.)

der Fakultät für Ingenieurwissenschaften, Informatik und Psychologie der Universität Ulm

von

Carmen Sippel

aus Crailsheim

Gutachter: Prof. Dr.-Ing. Robert Fischer Prof. Dr.-Ing. Volker Kühn Amtierende Dekanin: Prof. Dr. Anke Huckauf

Ulm, 16. Juni 2023

Acknowledgment

On the way to creating this thesis more than just one person was involved. Foremost, I would like to thank my supervisor, Prof. Dr.-Ing. Robert Fischer, for giving me the chance to work on the topic of compressed sensing and for the exhaustive discussions about it, where I could learn so much that all of it did not fit into this thesis. Moreover, I would like to thank him for enabling the cooperation with the institute for digital communications in Erlangen, as well as the colleagues in Erlangen for the inclusion to the institute. I also would like to thank the Deutsche Forschungsgemeinschaft (DFG) for funding most parts of my position. My thanks also go to Prof. Dr.-Ing. Martin Bossert, former head of the institute, for the opportunity to broaden the scope of my knowledge in a collaboration with coding theorists previous to the start of my actual PhD topic. Furthermore, I would like to thank Prof. Dr.-Ing. Volker Kühn for agreeing to examine my thesis on short notice and giving me the chance to present my work to his institute. Further thanks go to Prof. Dr.-Ing. Klaus Dietmayer for being part of the colloquium.

I want to thank my current and former scientific colleagues for the fruitful discussions and the administrative staff for their aid in all organizational and technical matters, as well as the institute as a whole for enjoyable breaks, which contributed to fulfilling working days. Among them, but not restricted to my colleagues, I have to thank the people that proofread this thesis and helped enhancing it. Last but not least, I want to thank my family for encouraging me to make my own decisions and my partner for motivating me at all times and his support in every aspect of life.

Abstract

The thesis at hand deals with recovery algorithms for the reconstruction of signals from noisy and linearly compressed measurements. The procedure of linearly compressing a signal with known structure is known as compressed sensing (CS); the recovery is possible under certain conditions.

Meanwhile, there is a number of algorithms that are designed for this kind of task; each with own benefits and drawbacks. The thesis focuses on algorithms that can be derived from a suitable processing in a factor graph, which embodies the structural dependencies of the CS problem. This approach already leads to a variety of algorithms.

The most compartmentalized factorization leads to a form of the factor graph from which the so-called message passing (MP) algorithm can be derived. The opposite case, a factorization with fewest possible factors results in a graph with two factors, yielding Turbo-type algorithms that iterate between those two factors. In the thesis, both types of algorithms are derived from the respective factor graph by constrained optimization of the Kullback–Leibler divergence between the posterior and a suitably adapted substitute distribution, i.e., the respective algorithms are reconciled as instances of the same basic idea.

The thesis also comprises sequential algorithms, which are derived from factorizations of the problem that lie in between the two extreme cases explained above. The resulting algorithms can be seen as sequentialized versions of the Turbo-type algorithms. There are two possibilities. One way to sequentialize the algorithm is by iterating through the measurements; the other is to iterate through the signal components (additionally to iterating between the two main parts). The sequentialization allows for various schedules, which are examined in numerical simulations.

Furthermore, an estimation-theoretic view is cast on the processing between the factors. The result reveals possibilities for improvement of the algorithms. Several variants of the algorithms are proposed and compared in numerical simulations.

Contents

Abstract					
1	Intro	oduction and Outline	1		
2	Compressed Sensing				
	2.1	System Model	6		
	2.2	Estimation Criteria	7		
		2.2.1 Minimum Mean-Squared Error Estimation	8		
		2.2.1.1 LMMSE Estimator	8		
		2.2.1.2 LMMSE Estimator for Scalar Observation	8		
		2.2.1.3 Individual Non-Linear MMSE Estimator	9		
		2.2.2 Approximate Inference	10		
		2.2.2.1 Moment Constraint	10		
		2.2.2.2 Entropy Maximization	11		
		2.2.2.3 Constrained Optimization	12		
		2.2.2.4 Specifying the Moment Constraint	14		
		2.2.2.5 Connection to Statistical Mechanics	14		
	2.3	Structure of the Compressed Sensing Problem	15		
		2.3.1 Factor Graph Representations	16		
		2.3.2 Representatives for the Factors	16		
3	Infe	rence on Factor Graphs	21		
	3.1	Structure of the Factor Graph	21		
		3.1.1 Factorization of the Substitute Distribution	22		
		3.1.2 The Moment-Matching Constraint	24		
	3.2	Constrained Optimization	25		
		3.2.1 Cost Function for Compressed Sensing	26		
		3.2.2 Stationary Points of the Lagrangian	26		
	3.3	Algorithms	28		
		3.3.1 Expectations	28		
		3.3.2 Projection to Exponential Family	29		
		3.3.3 Processing for Stationarity	30		
		3.3.4 Message Passing	35		
		3.3.5 Related Algorithms	37		
4	Turb	po-type Inference	41		
	4.1	Structure of the Factor Graph	42		

		4.1.1	Factorization of the Substitute Distribution	. 43
		4.1.2	The Moment-Matching Constraint	. 43
	4.2	Constr	rained Optimization	. 44
		4.2.1	Cost Function for Compressed Sensing	. 44
		4.2.2	Stationary Points of the Lagrangian	. 44
		4.2.3	The Lagrange Dual Function	. 45
	4.3	Turbo-	Type Algorithms—Optimization-Based Derivation	. 46
		4.3.1	Average Variances—VAMP	. 47
		4.3.2	Individual Variances	. 49
	4.4	Estima	tion-Theoretic Bias Compensation	. 50
		4.4.1	Channel-Constrained Estimation	. 50
			4.4.1.1 Average Unbiasing	. 50
			4.4.1.2 Individual Unbiasing	. 52
		4.4.2	Signal-Constrained Estimation	. 54
			4.4.2.1 Individual Unbiasing	. 57
			4.4.2.2 Average Unbiasing	. 58
	4.5	Turbo-	Type Algorithms—Estimation-Theoretic Adaption	. 59
		4.5.1	Individual Variances and Estimation-Theoretic Bias Compensation	. 60
		4.5.2	VAMP with Individually Unbiased LMMSE Estimate	. 61
	4.6	Block	Diagrams of the Turbo Algorithms	. 64
		4.6.1	VAMP	. 64
		4.6.2	VAMPind	. 64
		4.6.3	VAMPire	. 64
		4.6.4	VAMPia	. 65
		4.6.5	VAMPII	. 65
5	Sequ	4.6.5 Jential	VAMPII	. 65 67
5	Sequ 5.1	4.6.5 J ential Seque	VAMPII	. 65 67 . 67
5	Տeqւ 5.1	4.6.5 Jential Seque 5.1.1	VAMPII	. 65 67 . 67 . 68
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1	VAMPII	. 65 67 . 67 . 68 . 68
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1	VAMPII	. 65 67 . 67 . 68 . 68 . 68 . 69
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2	VAMPII	. 65 67 . 67 . 68 . 68 . 68 . 69 . 70
5	Seq ı 5.1	4.6.5 Jential Seque 5.1.1 5.1.2	VAMPII	. 65 67 . 67 . 68 . 68 . 68 . 69 . 70 . 70
5	Seq ı 5.1	4.6.5 Jential Seque 5.1.1 5.1.2	VAMPII	. 65 67 . 67 . 68 . 68 . 68 . 69 . 70 . 70 . 70
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2	VAMPII	. 65 67 . 67 . 68 . 68 . 68 . 69 . 70 . 70 . 70 . 71
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2 5.1.3	VAMPII	. 65 67 . 67 . 68 . 68 . 69 . 70 . 70 . 70 . 71 . 73
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2 5.1.3 5.1.4	VAMPII	. 65 67 . 67 . 68 . 68 . 69 . 70 . 70 . 70 . 70 . 71 . 73 . 75
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5	VAMPII Inference ntial Processing of Variables Structure of the Factor Graph 5.1.1.1 Factorization of the Substitute Distribution 5.1.1.2 The Moment-Matching Constraint Constrained Optimization Structure of the Function for Compressed Sensing 5.1.2.1 Cost Function for Compressed Sensing 5.1.2.3 Connection to Expectation Propagation Sequential Version of VAMPind Sequential Version with Estimation-Theoretic Bias Compensation	. 65 67 . 67 . 68 . 68 . 69 . 70 . 70 . 70 . 70 . 71 . 73 . 75 . 77
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5 Seque	VAMPII	. 65 67 . 67 . 68 . 68 . 69 . 70 . 70 . 70 . 70 . 71 . 73 . 75 . 77 . 78
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5 Seque 5.2.1	VAMPII Inference ntial Processing of Variables Structure of the Factor Graph 5.1.1.1 Factorization of the Substitute Distribution 5.1.1.2 The Moment-Matching Constraint Constrained Optimization Stationary Points of the Lagrangian 5.1.2.1 Cost Function for Compressed Sensing 5.1.2.2 Stationary Points of the Lagrangian 5.1.2.3 Connection to Expectation Propagation Sequential Version of VAMPind Sequential Version with Estimation-Theoretic Bias Compensation ntial Processing of Observations Structure of the Factor Graph	. 65 67 . 67 . 68 . 68 . 69 . 70 . 70 . 70 . 70 . 71 . 73 . 75 . 77 . 78 . 79
5	Sequ 5.1	4.6.5 Jential 1 Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5 Seque 5.2.1	VAMPII Inference ntial Processing of Variables Structure of the Factor Graph 5.1.1.1 Factorization of the Substitute Distribution 5.1.1.2 The Moment-Matching Constraint Constrained Optimization Structure 5.1.2.1 Cost Function for Compressed Sensing 5.1.2.2 Stationary Points of the Lagrangian 5.1.2.3 Connection to Expectation Propagation Sequential Version of VAMPind Sequential Version with Estimation-Theoretic Bias Compensation ntial Processing of Observations Structure of the Factor Graph 5.2.1.1 Factorization of the Substitute Distribution	 . 65 . 67 . 68 . 68 . 69 . 70 . 70 . 70 . 70 . 71 . 73 . 75 . 77 . 78 . 79 . 79
5	Sequ 5.1	4.6.5 Jential 1 Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5 Seque 5.2.1	VAMPit Inference ntial Processing of Variables Structure of the Factor Graph 5.1.1 Factorization of the Substitute Distribution 5.1.1.1 Factorization of the Substitute Distribution 5.1.1.2 The Moment-Matching Constraint Constrained Optimization Constraint 5.1.2.1 Cost Function for Compressed Sensing 5.1.2.2 Stationary Points of the Lagrangian 5.1.2.3 Connection to Expectation Propagation Sequential Version of VAMPind Sequential Version with Estimation-Theoretic Bias Compensation ntial Processing of Observations Structure of the Factor Graph 5.2.1.1 Factorization of the Substitute Distribution	 . 65 . 67 . 68 . 68 . 69 . 70 . 70 . 70 . 71 . 73 . 75 . 77 . 78 . 79 . 80
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5 Seque 5.2.1 5.2.2	VAMPII	. 65 67 . 67 . 68 . 68 . 69 . 70 . 70 . 70 . 70 . 70 . 70 . 70 . 71 . 73 . 75 . 77 . 78 . 79 . 80 . 80
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5 Seque 5.2.1 5.2.2	VAMPI	. 65 67 . 67 . 68 . 68 . 69 . 70 . 71 . 73 . 75 . 79 . 79 . 80 . 80 . 79 . 79 . 80 . 80 . 79 . 79 . 79 . 80 . 80 . 80 . 80 . 70 . 71 . 78 . 79 . 80 . 80
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5 Seque 5.2.1 5.2.2	VAMPII	. 65 67 . 67 . 68 . 68 . 69 . 70 . 75 . 79 . 80 . 81 . 81
5	Sequ 5.1	4.6.5 Jential I Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5 Seque 5.2.1 5.2.2 5.2.3	VAMPI	 . 65 . 67 . 68 . 68 . 69 . 70 . 70 . 70 . 70 . 71 . 73 . 75 . 77 . 78 . 79 . 80 . 81 . 82
5	Sequ 5.1	4.6.5 Jential Seque 5.1.1 5.1.2 5.1.3 5.1.4 5.1.5 Seque 5.2.1 5.2.2 5.2.3	VAMPI	 . 65 . 67 . 67 . 68 . 69 . 70 . 70 . 70 . 70 . 71 . 73 . 75 . 77 . 78 . 79 . 80 . 81 . 81 . 82 . 82

	5.2.3.3ExtrinsicExtrinsic5.2.3.4Procedure	83 85			
6	Numerical Results6.1Overview over Recovery Algorithms6.2Simulation Setup6.3I.i.d. Gaussian Sensing Matrix6.4Model for Wireless Sensor Networks6.5Schedules for the Sequential Processing of Variables6.6Variable Noise Scenario				
7	Conclusions and Future Work17.1Conclusions	15 115 117			
A	Vector Spaces 1 A.1 Vectors	19 119 121 121 121			
B	Convex Optimization 1	23			
С	Functional Derivatives 1	27			
D	Exponential Families1D.1Projection Property of the Kullback–Leibler Divergence1D.2Sufficient Statistics1D.3Derivative of the Log-Partition Function1	29 129 130			
		132			
E	Reconstruction Algorithms1E.1Gaussian Message Passing (GMP)1E.2GMP with Post-processed Extrinsic1E.3Approximate Message Passing (AMP)1E.4Generalized Approximate Message Passing (GAMP)1E.5Vector Approximate Message Passing (VAMP)1E.6VAMP with Individual Variances (VAMPind)1E.7VAMPire1E.8VAMP with Individually Unbiased LMMSE Estimator1E.9VAMPind with Sequentially Processed Variables1E.10VAMPire with Sequentially Processed Variables1E.11Sequential Processing of Observations1	132 134 135 136 137 138 139 140 141 142 143 144			

F.2.4 F.2.5 F.2.6 F.2.7 F.2.8	Other Variables	149 149 150 150 151						
Bibliography 1								
List of Scientific Publications and Presentations								
Curriculum Vitæ								

1. Introduction and Outline

Compressed sensing (a.k.a. compressive sampling) [CRT06, Don06] denotes the strategy of simplifying the compression procedure at the point where data is acquired, at the cost of higher effort for reconstruction of the signal at a receiving point. Hence, the strategy is in need for suitable algorithms that recover the signal as exact as possible with least possible expense.

Starting with ℓ_1 -minimization¹ [CRT06], which is a convex relaxation of the problem but still has relatively high complexity, the journey of possible reconstruction algorithms went over so-called *greedy* algorithms such as *orthogonal matching pursuit* (OMP) [PRK93] and *compressive sampling matching pursuit* (CoSaMP) [NT09], which *successively* build a support set for the recovered signal, to *iterative* algorithms. Both, greedy and iterative algorithms reduce the computational complexity compared to ℓ_1 -minimization. For the iterative algorithms, there is to mention iterative hard thresholding (IHT) [BD09] and iterative soft thresholding (IST) [DDD04], which both iterate between the solution of the linear system of equations obtained from the (noisy) compressed measurements—we denote this as the *channel-constrained* part—and a solution that follows a given signal structure—denoted as *signal-constrained* part the difference being the way of adjusting to the signal structure. While being much simpler than ℓ_1 -minimization, the recovery performance of these initial iterative algorithms suffered.

A huge step in terms of recovery performance was made with the invention of *approximate message passing* (AMP) [DMM09, Mal11] without increasing the complexity significantly compared to the other iterative algorithms. One step more complex and more robust to various compression scenarios is the so-called *vector approximate message passing* (VAMP) [RSF19] algorithm. The eponymous message-passing approach is an inference technique that is widely used in communications engineering, for example in the decoding of error-correcting codes such as low-density parity-check (LDPC) codes [Gal62] and Turbo codes [BGT93, HOP96]. Message passing (MP) usually works on graphical models, such as *factor graphs* [KFL01, Loe04], that represent the structure of the problem in detail, leading to messages being passed on the many edges in the graph, thereby creating the name. On the contrary, both, the AMP and the

¹Usually performed by the simplex algorithm [DT97]. For the definition of the ℓ_1 -norm, see Appendix A.



Figure 1.1: Visualization of the structure of the dissertation.

VAMP algorithm, combine so much of the detailed processing that in fact a Turbo-type processing results; similarly to IHT and IST, the iteration is between two parts but the recovery performance is much better.

In this thesis, the message-passing and Turbo-type approach are contrasted and both are derived from the framework of constrained optimization for the compressed sensing problem. Moreover, sequential approaches that stand in between both strategies are discussed. Furthermore, an estimation-theoretic point of view is cast on the exchange of the variables between the different solutions; leading to improved variants of the algorithms. The outline of the thesis is illustrated in Fig. 1.1 and described more detailed in the following.

Chapter 2 contains the theoretical basis to understand the reconstruction problem behind compressed sensing upon which the rest of the thesis is built. It defines the basic terms, especially for the compressed sensing model and states the estimations that are necessary to obtain solutions for the subproblems (signal- and channel-constrained part). Furthermore, we motivate why it is sensible to use the Kullback–Leibler divergence as cost function for the constrained optimization and state the different forms of the factor graph, that will be used in the subsequent chapters to derive the algorithms.

The next two chapters are depicted next to each other in Fig. 1.1, because they contain two opposing approaches to solve the recovery problem.

Chapter 3 focuses on the message-passing-type inference for CS with a detailed view on the connections in the factor graph. A message passing (MP) algorithm as in [KMS⁺12] is derived as result from minimizing the Kullback–Leibler divergence for a substitute distribution that is tailored to the factor graph for CS under expectation constraints. Similar derivations can be found in [YFW05] and [HOW⁺05] for general inference problems; here, we focus on the CS problem. The result stresses the importance of the computation of *extrinsic* [BG96] for the exchange of parameters between solutions of the subproblems.

In Chapter 4, the same derivation is used to obtain Turbo-type algorithms; the difference being how the factor graph is considered. In contrast to Chapter 3, the factor graph is combined to only contain two major parts between which the iterations take place. The results lead among others to the famous VAMP [RSF19] algorithm; with the *extrinsic* being the crucial point making VAMP able to stand out from the other iterative algorithms IHT, IST, and, e.g., iterative soft-feedback (ISF) [SF16], which is processed as IHT and IST, but uses the same signal-constrained estimation as VAMP. Additionally, several adaptions of the derived algorithms are discussed that result from an estimation-theoretic view on the extrinsic calculations.

Chapter 5 discusses two sequential recovery approaches that lie in between the two extreme cases of Turbo-type and message-passing-type algorithms. The derivation utilizes the same constrained optimization strategy as before for yet different factor graph representations. The chapter is divided into two parts; the first part (Sec. 5.1) considers the sequential processing of the signal components, while the second part (Sec. 5.2) focuses on sequentially processing the measurements.

In Chapter 6, the results of numerical simulations comparing the derived algorithms in different scenarios are shown. We depict the compressibility capability of the derived algorithms in the standard scenario for CS in Sec. 6.3, show in Sec. 6.4 a scenario, where the estimationtheoretic adaptions from Chapter 4 can beat the standard VAMP algorithm, discuss schedules for the sequential processing of signal components in Sec. 6.5, and present a scenario, where the sequential processing of measurements can improve over VAMP in Sec. 6.6.

Chapter 7 contains the conclusions and future work.

The Appendices A–D contain basic mathematical notions that are used throughout the thesis, Appendix E collects the recovery algorithms that are presented in the thesis and in Appendix F the notation and abbreviations used in the thesis can be found.

2. Compressed Sensing

The idea for compressed sensing (CS) dates back to 2006, where two groups of scientists independently realized that it is possible to reconstruct information from a linearly compressed signal under certain conditions [CRT06, Don06]. This is astonishing in the sense that the linear compression does not make use of actual features of the signal. So far, compression techniques had considered the signal itself to find out how to suitably compress it. Considering the signal, let's call it x, as vector¹ of an N-dimensional vector space (here \mathbb{R}^N), the new method requires only a vector-matrix-multiplication (we denote the matrix by $A \in \mathbb{R}^{M \times N}$ with M < N) and is therefore completely independent of signal \boldsymbol{x} . The result of the compression is called $y \in \mathbb{R}^M$ here. In engineering applications it usually inherits noise in the form of an additional random vector ² $\mathbf{n} \in \mathbb{R}^{M}$. The benefit of the linear compression technique lies in the fact that the computational effort is shifted to the recovery of signal x from the measurements y together with known sensing matrix A, whereas the deterministic technique splits the processing effort equally between compression and decompression. Especially the compression part is tremendously simplified by the linear acquisition of the data, since no deeper knowledge of x is leveraged. Usual decompression obtains what it believes to be the original signal x directly from y, which carries the recipe to the reconstruction. In the compressed sensing approach, y does not carry this recipe and therefore requires powerful algorithms for the reconstruction of the signal. The two strategies are summarized and compared in Fig. 2.1. Double strike boxes indicate non-linear functions.

Naturally, the question arises, how reconstruction of the signal x is possible at all, if the compression process is not tailored the signal.

The answer lies in the structure of the signal.

Common sources for signals, such as photographs, speech, or binary transmission streams, although providing different realizations, usually inherit a certain statistic that applies to all

¹The notion of vectors and vector spaces is briefly introduced in Appendix A; the general notation is explained in Appendix F.2

²Note that variables in serif font, e.g., x or x, denote realizations of random variables or vectors in sans-serif font, e.g., \times or **x**.



Figure 2.1: Usual transmission scheme (upper row) and compressed sensing idea (lower row).

realizations. Making use of this statistic is key to signal recovery in CS. A basic notion that specifies the structure of the signal is its *sparsity* [EK12]. A signal is said to be *sparse* if it has few non-zero entries; the sparsity being the number of non-zero entries. If the signals of a source can be represented such that they inherit a certain sparsity, this sparsity gives a natural bound on how much the signals can be compressed, i.e., the observation dimension M [EK12].

With this knowledge it is possible to describe our model for CS in detail.

2.1 System Model

Following an engineering perspective, a noisy model for the measurements in CS is considered

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n} \in \mathbb{R}^M , \qquad (2.1)$$

where the sensing matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ is known and $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma_{\mathbf{n}}^{2} \mathbf{I}_{M})$ is independent and identically distributed (i.i.d.) Gaussian. Entries and rows of the sensing matrix are defined as follows

$$\boldsymbol{A} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1N} \\ a_{21} & a_{22} & \dots & a_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ a_{M1} & a_{M2} & \dots & a_{MN} \end{bmatrix} = \begin{bmatrix} \boldsymbol{a}_1^{\top} \\ \boldsymbol{a}_2^{\top} \\ \vdots \\ \boldsymbol{a}_M^{\top} \end{bmatrix}, \quad \boldsymbol{a}_i = [a_{i1}, a_{i2}, \dots, a_{iN}]^{\top}.$$
(2.2)

So far, it has been withhold that the sensing matrix also needs to fulfill some properties for recovery to be possible [EK12]. The basic point is that columns of the sensing matrix should be as incoherent³ as possible to each other, so that the observation vector \boldsymbol{y} carries as much (different) information as possible about the individual entries of \boldsymbol{x} . A counterexample for a coherent matrix is one with two equal columns [EK12]; the corresponding elements of \boldsymbol{x} influence \boldsymbol{y} in the same way and can therefore not be distinguished anymore. Another astonishing

 $^{^{3}}$ For the definition of the *coherence*, see Appendix A.

fact figured out in [Don06] is that, when the elements a_{ij} ($i \in \{1, ..., M\}$, $j \in \{1, ..., N\}$) are drawn i.i.d. for example from a Gaussian distribution, the resulting matrix is likely to have good properties. This fact will be used in the simulations in Chapter 6.

Let us turn to the statistic of the signal **x**. For simplicity, it is assumed that the elements x_j of $\mathbf{x} = [x_1, \dots, x_N]^\top$ are i.i.d. with marginal probability density function (pdf) $f_x(x_j)$, meaning the pdf of **x** is separable,

$$\mathbf{f}_{\mathbf{x}}(\boldsymbol{x}) = \prod_{j=1}^{N} \mathbf{f}_{\mathbf{x}}(x_j) .$$
(2.3)

A zero-mean signal **x** is assumed throughout the thesis. The fact that **x** is sparse is considered via a Dirac delta function $\delta(x)$ at zero in the marginal pdf. A possible pdf is given by [SF14]

$$f_{x}(x) = \frac{s}{2N}\delta(x+1) + \frac{N-s}{N}\delta(x) + \frac{s}{2N}\delta(x-1) , \qquad (2.4)$$

which we call the discrete ternary (DT) prior. A prior that is very commonly used in CS is the so-called Bernoulli-Gaussian (BG) prior, whose non-zero elements are Gaussian distributed. The definition used in this thesis is given by

$$\mathbf{f}_{\mathsf{x}}(x) = \left(1 - \frac{s}{N}\right)\delta(x) + \frac{s/N}{\sqrt{2\pi}}\exp\left(-\frac{x^2}{2}\right) \ . \tag{2.5}$$

The task in CS is to recover the signal x for given observation y, sensing matrix A and known noise variance σ_n^2 .

2.2 Estimation Criteria

The problem of recovering the signal x from the compressed measurements y under a probabilistic view can be solved by finding the vector \hat{x}_{MAP} that maximizes the probability that under the condition that y was measured, x is the vector which has been compressed, i.e.,

$$\hat{\boldsymbol{x}}_{\text{MAP}} = \operatorname*{argmax}_{\boldsymbol{x}} \mathsf{f}_{\mathsf{x}|\mathsf{y}}(\boldsymbol{x}) , \qquad (2.6)$$

which is the so-called maximum a-posteriori (MAP) criterion [Kay93].

Using Bayes' rule, the conditional pdf or posterior, corresponding to the CS problem, is given by

$$f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) = \frac{1}{f_{\mathbf{y}}(\boldsymbol{y})} f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) \cdot f_{\mathbf{x}}(\boldsymbol{x})$$
(2.7)

with channel part or likelihood function ($\|\cdot\|_2$ denotes the ℓ_2 -norm, i.e., for a real vector \boldsymbol{x} it is $\|\boldsymbol{x}\|_2^2 = \boldsymbol{x}^\top \boldsymbol{x}$, cf. Appendix A)

$$f_{\mathbf{y}|\mathbf{x}}(\mathbf{x}) = f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{x}) = \frac{1}{\sqrt{(2\pi\sigma_{\mathbf{n}}^2)^M}} \exp\left(-\frac{1}{2\sigma_{\mathbf{n}}^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2\right) .$$
(2.8)

We call $f_x(x)$ the signal part or a-priori distribution or simply *prior*.

2.2.1 Minimum Mean-Squared Error Estimation

Another approach is the <u>minimum mean-squared error</u> (MMSE) estimation criterion, which leads to the conditional mean⁴ [Kay93]⁵

$$\hat{\boldsymbol{x}}_{\text{MMSE}} \stackrel{\text{def}}{=} E_{\boldsymbol{x}} \{ \boldsymbol{x} \mid \boldsymbol{y} \} = E_{\boldsymbol{x} \sim f_{\boldsymbol{x}|\boldsymbol{y}}} \{ \boldsymbol{x} \} = \int \boldsymbol{x} f_{\boldsymbol{x}|\boldsymbol{y}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} \; . \tag{2.9}$$

MAP and MMSE estimate coincide if x and y are jointly Gaussian (which is not the case for the given problem). Obtaining these solutions directly is infeasible due to the high dimensionality of the signal x. Hence, there is a need for tractable procedures to recover the signal x. First, two separate MMSE solutions are considered, which solve the two main parts of the problem by approximating the other, respectively.

2.2.1.1 LMMSE Estimator

Considering the *channel part* $f_{y|x}(x)$ of the conditional pdf (2.7) and assuming that the signal vector is Gaussian distributed with mean m_x and covariance matrix Φ_x , i.e., $\mathbf{x} \sim \mathcal{N}(m_x, \Phi_x)$, (which is not the case in CS) linear estimation is optimal in the MMSE sense [Kay93]. The corresponding conditional mean estimator reads [Kay93]

$$\boldsymbol{m}_{c} = \mathrm{E}_{\mathbf{x}}\{\mathbf{x} \mid \boldsymbol{y}\} = \frac{\int \boldsymbol{x} \mathrm{f}_{\mathbf{y} \mid \mathbf{x}}(\boldsymbol{x}) \mathrm{f}_{\mathbf{x}}(\boldsymbol{x}) \mathrm{d}\boldsymbol{x}}{\int \mathrm{f}_{\mathbf{y} \mid \mathbf{x}}(\boldsymbol{x}) \mathrm{f}_{\mathbf{x}}(\boldsymbol{x}) \mathrm{d}\boldsymbol{x}}$$

$$= \boldsymbol{m}_{\mathbf{x}} + \boldsymbol{\Phi}_{\mathbf{x}} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{\Phi}_{\mathbf{x}} \boldsymbol{A}^{\top} + \sigma_{n}^{2} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \boldsymbol{m}_{\mathbf{x}})$$

$$= \boldsymbol{m}_{\mathbf{x}} + (\boldsymbol{A}^{\top} \boldsymbol{A} + \sigma_{n}^{2} \boldsymbol{\Phi}_{\mathbf{x}}^{-1})^{-1} \boldsymbol{A}^{\top} (\boldsymbol{y} - \boldsymbol{A} \boldsymbol{m}_{\mathbf{x}}) , \qquad (2.10)$$

the conditional covariance matrix of the estimation error is obtained as [Kay93]

$$\boldsymbol{\Phi}_{c} = E_{\mathbf{x}} \left\{ (\mathbf{x} - \boldsymbol{m}_{c}) (\mathbf{x} - \boldsymbol{m}_{c})^{\top} \mid \boldsymbol{y} \right\} = \sigma_{n}^{2} \left(\boldsymbol{A}^{\top} \boldsymbol{A} + \sigma_{n}^{2} \boldsymbol{\Phi}_{x}^{-1} \right)^{-1} .$$
(2.11)

The diagonal entries of the covariance matrix are denoted by $\sigma_{c,j}^2 \stackrel{\text{def}}{=} [\boldsymbol{\Phi}_c]_{jj}$. The estimates \boldsymbol{m}_c and $\boldsymbol{\Phi}_c$ are called the <u>channel</u>-constrained estimates or LMMSE (linear minimum mean-squared error) estimates.

Note that the second line in (2.10) results in a less complex implementation, since an $M \times M$ matrix has to be inverted, which is easier than inverting the $N \times N$ matrix in the third line, because in the CS scenario typically M < N.

2.2.1.2 LMMSE Estimator for Scalar Observation

Now, consider that only the *i*th row $\mathbf{a}_i^{\top} = [a_{i1}, \ldots, a_{iN}]$ of the sensing matrix \mathbf{A} is given and, thus, only a scalar observation y_i is measured, i.e., $(\mathbf{n}_i \sim \mathcal{N}(0, \sigma_{\mathbf{n}}^2) \text{ for } i \in \{1, \ldots, M\})$

$$\mathsf{y}_i = \boldsymbol{a}_i^\top \mathsf{x} + \mathsf{n}_i \;. \tag{2.12}$$

⁴For convenience two notations are used to denote the estimation w.r.t. a certain pdf, either by indicating which pdf is used inside the integral or by conditioning the argument of the expectation.

⁵Whenever the integration boundaries are missing, the integral is restricted to the domain of the integrand.

However, as above a Gaussian signal vector **x** with mean m_x and covariance matrix Φ_x is assumed. Then, the optimal conditional estimator for the entire signal vector x is again linear and looks as follows

$$\boldsymbol{m}_{\mathrm{c},i} = \mathrm{E}_{\mathbf{x}}\{\mathbf{x} \mid y_i\} = \boldsymbol{m}_{\mathrm{x}} + \boldsymbol{\Phi}_{\mathrm{x}}\boldsymbol{a}_i(\boldsymbol{a}_i^{\top}\boldsymbol{\Phi}_{\mathrm{x}}\boldsymbol{a}_i + \sigma_{\mathrm{n}}^2)^{-1}(y_i - \boldsymbol{a}_i^{\top}\boldsymbol{m}_{\mathrm{x}}), \qquad (2.13)$$

with (conditional) covariance matrix of the estimation error

$$\boldsymbol{\Phi}_{\mathrm{c},i} = \mathrm{E}_{\mathbf{x}} \left\{ (\mathbf{x} - \boldsymbol{m}_{\mathrm{c},i}) (\mathbf{x} - \boldsymbol{m}_{\mathrm{c},i})^{\top} \mid y_i \right\} = \sigma_{\mathsf{n}}^2 \left(\boldsymbol{a}_i \boldsymbol{a}_i^{\top} + \sigma_{\mathsf{n}}^2 \boldsymbol{\Phi}_{\mathsf{x}}^{-1} \right)^{-1} .$$
(2.14)

Note that the inverse in (2.13) is scalar. Via the Sherman-Morrison-Woodbury identity [GV96], the computation of the conditional variance $\boldsymbol{\Phi}_{c,i}$ can as well be written with a scalar inversion, instead of inverting the $N \times N$ matrix as stated above, which results in

$$\boldsymbol{\Phi}_{\mathrm{c},i} = \boldsymbol{\Phi}_{\mathrm{x}} - \frac{1}{\sigma_{\mathrm{n}}^{2} + \boldsymbol{a}_{i}^{\top} \boldsymbol{\Phi}_{\mathrm{x}} \boldsymbol{a}_{i}} \boldsymbol{\Phi}_{\mathrm{x}} \boldsymbol{a}_{i} \boldsymbol{a}_{i}^{\top} \boldsymbol{\Phi}_{\mathrm{x}} .$$
(2.15)

With $\boldsymbol{m}_{c,i} = [m_{c,i,1}, \ldots, m_{c,i,N}]^{\top}$ and $\sigma_{c,i,j}^2 = [\boldsymbol{\Phi}_{c,i}]_{jj}$ as well as $\boldsymbol{m}_{\mathsf{x}} = [m_{\mathsf{x},1}, \ldots, m_{\mathsf{x},N}]^{\top}$ and $\sigma_{\mathsf{x},j}^2 = [\boldsymbol{\Phi}_{\mathsf{x}}]_{jj}$ the respective *j*th components can be written as

$$m_{\mathrm{c},i,j} = m_{\mathrm{x},j} + \frac{\sigma_{\mathrm{x},j}^{2} a_{ij} \left(y_{i} - \sum_{j'=1}^{N} a_{ij'} m_{\mathrm{x},j'} \right)}{\sigma_{\mathrm{n}}^{2} + \sum_{j'=1}^{N} \sigma_{\mathrm{x},j'}^{2} a_{ij'}^{2}} \\ = \frac{m_{\mathrm{x},j} \sigma_{\mathrm{n}}^{2} + m_{\mathrm{x},j} \sum_{j'=1}^{N} \sigma_{\mathrm{x},j'}^{2} a_{ij'}^{2} + \sigma_{\mathrm{x},j}^{2} a_{ij} \left(y_{i} - \sum_{j'=1}^{N} a_{ij'} m_{\mathrm{x},j'} \right)}{\sigma_{\mathrm{n}}^{2} + \sum_{j'=1}^{N} \sigma_{\mathrm{x},j'}^{2} a_{ij'}^{2}} , \qquad (2.16)$$

$$\sigma_{\mathrm{n}}^{2} = \sigma_{\mathrm{x},j}^{2} \left(1 - \frac{\sigma_{\mathrm{x},j}^{2} a_{ij}^{2}}{\sigma_{\mathrm{x},j}^{2} a_{ij}^{2}} \right)$$

$$\sigma_{\mathbf{c},i,j}^{2} = \sigma_{\mathbf{x},j}^{2} \left(1 - \frac{1 - \frac{1}{\sigma_{\mathbf{n}}^{2} + \sum_{j'=1}^{N} \sigma_{\mathbf{x},j'}^{2} a_{ij'}^{2}}}{\sigma_{\mathbf{n}}^{2} + \sum_{j'=1}^{N} \sigma_{\mathbf{x},j'}^{2} a_{ij'}^{2}} \right) = \sigma_{\mathbf{x},j}^{2} \frac{\sigma_{\mathbf{n}}^{2} + \sum_{j'=1}^{N} \sigma_{\mathbf{x},j'}^{2} a_{ij'}^{2}}{\sigma_{\mathbf{n}}^{2} + \sum_{j'=1}^{N} \sigma_{\mathbf{x},j'}^{2} a_{ij'}^{2}} .$$
(2.17)

2.2.1.3 Individual Non-Linear MMSE Estimator

The opposite view considers the *signal part* $f_x(x)$ and assumes Gaussian noise. Indeed, Gaussian noise is given in the model (2.1), however, the form of $f_{y|x}(x)$ in (2.8) with couplings between x and y given by the sensing matrix makes the problem (2.9) intractable. Since separability is assumed, i.e., independent random variables x_j , the channel-constrained part of the posterior is approximated by ignoring the couplings and assuming independent measurements $\tilde{x}_{s,j}$, which represent the signal x_j in zero-mean Gaussian noise, i.e., the observation is a realization of the random variable

$$\tilde{\mathbf{x}}_{s,j} = \mathbf{x}_j + \mathbf{n}_{s,j}, \quad j \in \{1, \dots, N\},$$
(2.18)

where $x_j \sim f_x(x_j)$ and $n_{s,j} \sim \mathcal{N}(0, \tilde{\sigma}_{s,j}^2)$. The corresponding estimator is non-linear but can be computed per component, i.e., the conditional mean $\boldsymbol{m}_s = [m_{s,1}, \ldots, m_{s,N}]^{\top}$ can be

individually computed via

$$m_{\mathbf{s},j} = \mathcal{E}_{\mathsf{x}_j} \{ \mathsf{x}_j \mid \tilde{x}_{\mathbf{s},j} \} = \frac{\int x_j \mathsf{f}_\mathsf{x}(x_j) \exp\left(-\frac{(\tilde{x}_{\mathbf{s},j} - x_j)^2}{2\tilde{\sigma}_{\mathbf{s},j}^2}\right) \, \mathrm{d}x_j}{\int \mathsf{f}_\mathsf{x}(x_j) \exp\left(-\frac{(\tilde{x}_{\mathbf{s},j} - x_j)^2}{2\tilde{\sigma}_{\mathbf{s},j}^2}\right) \, \mathrm{d}x_j} ,$$
(2.19)

respectively, the conditional variance of estimation error

$$\sigma_{s,j}^{2} = E_{x_{j}} \left\{ (x_{j} - m_{s,j})^{2} \mid \tilde{x}_{s,j} \right\} = \frac{\int (x_{j} - m_{s,j})^{2} f_{x}(x_{j}) \exp\left(-\frac{(\tilde{x}_{s,j} - x_{j})^{2}}{2\tilde{\sigma}_{s,j}^{2}}\right) dx_{j}}{\int f_{x}(x_{j}) \exp\left(-\frac{(\tilde{x}_{s,j} - x_{j})^{2}}{2\tilde{\sigma}_{s,j}^{2}}\right) dx_{j}} , \qquad (2.20)$$

for $j \in \{1, \ldots, N\}$. The exact computation depends of course on the prior $f_x(x_j)$ itself. Note that the estimation requires both variables, $\tilde{x}_{s,j}$ and $\tilde{\sigma}_{s,j}^2$, as inputs, however, we will only denote this in the conditional expectation in case that the respective dependency needs to be emphasized. In the CS literature, these computations are usually referred to as *denoising* and the respective estimator is the *denoiser*. Here, $m_{s,j}$ and $\sigma_{s,j}^2$ ($j \in \{1, \ldots, N\}$) may be called signal-constrained estimates.

2.2.2 Approximate Inference

In Sec. 2.2 (cf. (2.6) and (2.9)) it was shown that $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$ is the function of interest in the CS scenario. One way to approach the CS problem is to replace the "complicated" posterior $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$ by a simpler substitute distribution that is called $q(\mathbf{x})$. A natural question that arises is how to choose $q(\mathbf{x})$ suitably in order to approximate the posterior adequately. In the following, we derive why it is suitable to approach the target distribution $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$ by minimizing the Kullback–Leibler divergence or cross entropy [KL51] (log(·) refers to the natural logarithm)

$$D_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) \mid\mid f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})) \stackrel{\text{def}}{=} \int \mathbf{q}(\boldsymbol{x}) \log \frac{\mathbf{q}(\boldsymbol{x})}{\mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})} \, \mathrm{d}\boldsymbol{x} \,. \tag{2.21}$$

To that end, it is motivated why the (differential) entropy of the substitute distribution q(x) should be maximized, just to find that the entropy itself should also be used to connect the substitute distribution q(x) with the target distribution $f_{x|y}(x)$. The result is also motivated in [SJ80, SJ81] under the name *cross entropy minimization*. It is related to the so-called *maximum entropy method* [Jay82, WJ08, Wu12], but utilizes additionally prior knowledge. It is furthermore used in [Moh97a] for the design of iterative decoders and detectors in digital communication scenarios.

2.2.2.1 Moment Constraint

The substitute distribution q(x) should relate to the target $f_{x|y}(x)$ in some way. This is obtained by requiring the equality of expectations based on the two distributions of some yet to be specified moment. Let U(x) denote the statistic of interest, then it is required that

$$\mathbf{E}_{\mathbf{x}\sim q}\{U(\mathbf{x})\} = \mathbf{E}_{\mathbf{x}\sim f_{\mathbf{x}\mid \mathbf{y}}}\{U(\mathbf{x})\} .$$
(2.22)

Since there are many distributions q(x) that fulfill this condition, a criterion that specifies how to choose a "suitable" one is needed.

2.2.2.2 Entropy Maximization

In order to answer the question how to approximate the posterior suitably, we make use of the information theoretic quantity of uncertainty, as introduced by Shannon in 1948 [Sha48]. Since the problem under consideration deals with continuous distributions, indeed, the *differential entropy* is required. For a distribution q(x), it is defined as

$$h(\mathbf{q}(\boldsymbol{x})) = -\int \mathbf{q}(\boldsymbol{x}) \log \mathbf{q}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} \,. \tag{2.23}$$

The notion of entropy offers a way to measure the uncertainty of a random variable, which gives us a criterion for the adjustment of q(x). Indeed, maximizing the (differential) entropy (2.23) prevents us from preferring certain possibilities over others when there is no reason to do so [Jay57], i.e., a distribution with maximum entropy allows the most flexibility in choosing x. This philosophy is known as *maximum entropy method* [Jay82, WJ08, Wu12].

For example, in the discrete case with finite set \mathcal{X} , the uniform distribution is the one that maximizes entropy [Sha48], which means that every possible $x \in \mathcal{X}$ has the same probability of being the solution; no element is preferred over another. Under a moment constraint the distribution may change. With specified second order moment, the Gaussian distribution is the distribution that maximizes (differential) entropy (in the continuous case), which means that the Gaussian distribution is the "most random" one in this case [HF14]. In order to derive this for the scalar case in our notation, the maximization is formulated as minimization task with certain equality constraints⁶

minimize
$$-h(q(x))$$

subject to $E_{x\sim q} \{x^2\} = \sigma_x^2$, (2.24)
 $\int q(x) dx = 1$.

i.e., the differential entropy is maximized while restricting the second order moment to some value σ_x^2 and the integral over q(x) to 1, in order to obtain a valid pdf. The usual approach to bring together the optimization task with the constraints is the combination in a so-called Lagrangian, which reads in this case

$$\mathcal{L}(\mathbf{q}(x),\,\mu,\,\nu) = -\mathbf{h}(\mathbf{q}(x)) + \mu\left(\mathbf{E}_{\mathsf{x}\sim\mathsf{q}}\left\{\mathsf{x}^{2}\right\} - \sigma_{\mathsf{x}}^{2}\right) + \nu\left(\int \mathbf{q}(x)\,\mathrm{d}x - 1\right) \,. \tag{2.25}$$

In order to find (arguments of) extreme points in this functional, the functional derivative (for a detailed explanation, see Appendix C) of the Lagrangian is necessary. With

$$\frac{\partial}{\partial \mathbf{q}} \int x^2 \mathbf{q}(x) \, \mathrm{d}x = x^2 \tag{2.26}$$

⁶For a summary of the notions from optimization theory, which are used within this thesis, the reader is referred to Appendix B.

and (C.8) it is

$$\frac{\partial \mathcal{L}}{\partial \mathbf{q}} = 1 + \log \mathbf{q}(x) + \mu x^2 + \nu \stackrel{!}{=} 0 , \qquad (2.27)$$

which can be solved for q(x) to obtain

$$q(x) = \exp(-\mu x^2 - 1 - \nu)$$
 (2.28)

The second constraint in (2.24) causes normalization via the multiplier ν , i.e., it forces $\exp(1 + \nu) = \int \exp(-\mu x^2) \, dx = \sqrt{\pi/\mu}$. Hence, the result is a Gaussian distribution with variance $\sigma^2 = 1/(2\mu)$. Since the mean was not specified, it ended up to be zero. Additionally specifying the mean would similarly have resulted in a Gaussian distribution with respective mean. Of course, for non-zero mean, the value of $E_{x\sim q}\{x^2\}$ would not be the variance but the second-order non-centralized moment instead.

2.2.2.3 Constrained Optimization

On the way to motivating the use of the Kullback–Leibler divergence, a constrained optimization problem for the task of approximating $f_{x|y}(x)$ with a substitute distribution q(x) under a given moment constraint (and a normalization constraint) can now be stated

minimize
$$-h(\mathbf{q}(\mathbf{x}))$$

subject to $E_{\mathbf{x}\sim q}\{U(\mathbf{x})\} = E_{\mathbf{x}\sim f_{\mathbf{x}|\mathbf{y}}}\{U(\mathbf{x})\}, \qquad (2.29)$
 $\int \mathbf{q}(\mathbf{x}) \, \mathrm{d}\mathbf{x} = 1.$

Since the entropy is concave in its corresponding distribution, the optimization problem is convex and therefore uniquely solvable, as long as there are feasible points. Although this approach seems as infeasible as the ones before, it gives an interesting viewpoint on the problem. Indeed, the point of view on the problem can be changed by considering the so-called *Lagrange dual problem*, that is obtained by minimizing the *Lagrangian*, which additively combines the above minimization problem (2.29) with the constraints, over \boldsymbol{x} . This new point of view yields a representation that leads to the task of minimizing the Kullback–Leibler divergence, if the statistic $U(\boldsymbol{x})$ is chosen accordingly, which will get clear below. Before turning to that, the mathematical derivation of the dual problem is considered.

With Lagrangian multipliers β and ν , the respective Lagrangian reads

$$\mathcal{L}(\mathbf{q}(\boldsymbol{x}), \beta, \nu) = -h(\mathbf{q}(\boldsymbol{x})) + \beta \left(E_{\mathbf{x} \sim \mathbf{q}} \{ U(\mathbf{x}) \} - E_{\mathbf{x} \sim f_{\mathbf{x}|\mathbf{y}}} \{ U(\mathbf{x}) \} \right) + \nu \left(\int \mathbf{q}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 \right) .$$
(2.30)

Differentiation with respect to q(x) yields a functional derivative, see Appendix C

$$\frac{\partial}{\partial \mathbf{q}(\boldsymbol{x})} \mathcal{L}(\mathbf{q}(\boldsymbol{x}), \beta, \nu) = 1 + \log \mathbf{q}(\boldsymbol{x}) + \beta U(\boldsymbol{x}) + \nu \stackrel{!}{=} 0$$
(2.31)

$$\Leftrightarrow \qquad \qquad \mathbf{q}(\boldsymbol{x}) = \frac{1}{Z_U(\beta)} \exp(-\beta U(\boldsymbol{x})) \quad , \qquad (2.32)$$



Figure 2.2: Schematic representation of the inequality (2.35).

where $Z_U(\beta) = \exp(1 + \nu) = \int \exp(-\beta U(\boldsymbol{x})) \, d\boldsymbol{x}$ fulfills the normalization constraint. The Lagrange dual function is then given by

$$\mathcal{L}_{\mathrm{D}}(\beta) = -\log Z_{U}(\beta) - \beta \mathrm{E}_{\mathbf{x} \sim f_{\mathbf{x}|\mathbf{y}}} \{ U(\mathbf{x}) \} .$$
(2.33)

The dual problem is specified by the maximization of the Lagrange dual function, thus, another approach to obtain the best substitute distribution $q^*(x)$ is to maximize the right-hand side of (2.33), i.e.,

$$-h(\mathbf{q}^*(\boldsymbol{x})) = \sup_{\beta} \left(-\beta E_{\mathbf{x} \sim f_{\mathbf{x}|\mathbf{y}}} \{ U(\mathbf{x}) \} - \log Z_U(\beta) \right) , \qquad (2.34)$$

where $q^*(x)$ is the optimizer⁷. This means especially that for feasible q(x) and arbitrary β it is

$$-h(\mathbf{q}(\boldsymbol{x})) \ge -\beta E_{\mathbf{x} \sim f_{\mathbf{x}|\mathbf{y}}} \{ U(\mathbf{x}) \} - \log Z_U(\beta) , \qquad (2.35)$$

To verify this intuitively, consider that, on the one hand, the entropy is concave in its distribution, i.e., $-h(q(\boldsymbol{x}))$ is convex. On the other hand, the supremum on the right-hand side of (2.34) is obviously an upper bound for its argument. Hence, the optimizer $q^*(\boldsymbol{x})$ defines the point, where both functionals touch. Moving away from that point yields the inequality. This is depicted schematically in Fig. 2.2. Noteworthy, the argument of the supremum is concave in β since the integral over the convex but positive $\exp(-\beta U(\boldsymbol{x}))$ is convex, as well as $\beta E_{\mathbf{x} \sim f_{\mathbf{x}|\mathbf{y}}} \{U(\mathbf{x})\}$, which is linear in β , i.e., both concave and convex, so the negative sum of both terms is concave.

Noteworthy, the optimization problem can as well be solved by considering any reformulation of the inequality (2.35) and pulling it towards equality. For feasible $q(\boldsymbol{x})$, which fulfill the moment constraint, $E_{\mathbf{x}\sim f_{\mathbf{x}|\mathbf{y}}}\{U(\mathbf{x})\}$ can be replaced by $E_{\mathbf{x}\sim q}\{U(\mathbf{x})\}$. This way, the dual representation (2.33) of the optimization problem can be reformulated to the inequality

$$-\log Z_U(\beta) \le \beta \mathbb{E}_{\mathbf{x} \sim \mathbf{q}} \{ U(\mathbf{x}) \} - h(\mathbf{q}(\mathbf{x})) , \qquad (2.36)$$

i.e., an equivalent optimization problem is obtained by minimizing the right-hand side of (2.36). Since, the inequality (2.36) holds for any β , one may deliberately specify it to a certain value, say $\beta = 1$. This specification will be kept throughout the thesis.

⁷The optimizer is indeed achieved. Since the optimization problem does not inherit inequality constraints, *Slater's condition* holds whenever the constraints are fulfilled, i.e., *strong duality* [BV04] holds.

2.2.2.4 Specifying the Moment Constraint

So far, an arbitrary statistic $U(\mathbf{x})$ has been considered that should tie the moments of $q(\mathbf{x})$ and $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$ together. Recalling that the differential entropy can be written as

$$h(f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})) = -\int f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) \log f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = E_{\mathbf{x} \sim f_{\mathbf{x}|\mathbf{y}}} \left\{ -\log f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) \right\}, \quad (2.37)$$

a suitable statistic may be given by

$$U(\boldsymbol{x}) = -\log f_{\boldsymbol{x}|\boldsymbol{y}}(\boldsymbol{x}) . \tag{2.38}$$

This specification means that the moments specified by $-\log f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$ of the substitute distribution $q(\mathbf{x})$ should match the entropy of the posterior $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$. With $U(\mathbf{x})$ as defined in (2.38) and $\beta = 1$, the right-hand side of (2.36) can be rewritten to [YFW05]⁸

$$E_{\mathbf{x}\sim q}\{U(\mathbf{x})\} - h(\mathbf{q}(\mathbf{x})) = \int U(\mathbf{x})\mathbf{q}(\mathbf{x}) \,\mathrm{d}\mathbf{x} + \int \mathbf{q}(\mathbf{x})\log \mathbf{q}(\mathbf{x}) \,\mathrm{d}\mathbf{x}$$
$$= \int \mathbf{q}(\mathbf{x})\log \mathbf{q}(\mathbf{x}) \,\mathrm{d}\mathbf{x} - \int \mathbf{q}(\mathbf{x})\log \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x}) \,\mathrm{d}\mathbf{x}$$
$$= \int \mathbf{q}(\mathbf{x})\log \frac{\mathbf{q}(\mathbf{x})}{\mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x})} \,\mathrm{d}\mathbf{x}$$
$$= D_{\mathrm{KL}}(\mathbf{q}(\mathbf{x}) || \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x})), \qquad (2.39)$$

with $D_{KL}(\cdot||\cdot)$ being the *Kullback–Leibler divergence* [KL51] or *cross entropy*, which was defined in (2.21). The path behind clarifies that the Kullback–Leibler divergence measures the "closeness" between two densities on the basis of the differential entropy. This shows again the actual aim in this procedure to approximate the target distribution $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$, since the minimum of the Kullback–Leibler divergence $(D_{KL}(\mathbf{q}(\mathbf{x})||\mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x})) = 0)$ is achieved exactly, if $\mathbf{q}(\mathbf{x}) = \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$. When minimizing the Kullback–Leibler divergence, one should always keep in mind that the basic idea behind the procedure is to maximize the entropy of a substitute distribution $\mathbf{q}(\mathbf{x})$ while staying consistent to the differential entropy of the original density $\mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$, which means the goal is to find the most freely chosen distribution $\mathbf{q}(\mathbf{x})$ that fits to prior knowledge of $\mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$ [SJ80].

2.2.2.5 Connection to Statistical Mechanics

In statistical mechanics, the distribution in (2.32) is known as *Boltzmann-Gibbs distribution*, which describes the distribution of particles of a system in thermal equilibrium, i.e., when there is no energy flow [Ber20]. Indeed, the function $U(\mathbf{x})$ is hereby the *energy* or *Hamiltonian* of the system and β is the inverse temperature. All in all, the result of the derivations coincides with the second law of thermodynamics, stating that under the constraint of conserving energy, the entropy of a system is always maximized, which results in the distribution of particles according to Boltzmann's law (2.32). It can be shown that when the system freezes

⁸When the energy is only defined for the part that depends on \boldsymbol{x} , i.e., by separating the posterior according to $f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) = \frac{1}{Z}f(\boldsymbol{x})$ and defining $U(\boldsymbol{x}) = -\log f(\boldsymbol{x})$, it is $F_{V}(\mathbf{q}(\boldsymbol{x})) = -\log Z + D_{KL}(\mathbf{q}(\boldsymbol{x}))|f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}))$.

 $(\beta \to \infty)$, the distribution bundles up around the state x that minimizes U(x), i.e., the system converges to a state of minimum energy. If the energy represents the underlying system via the specification of U(x) as in (2.38) the minimization of the energy

$$\hat{\boldsymbol{x}} = \operatorname*{argmin}_{\boldsymbol{x}} U(\boldsymbol{x}) \tag{2.40}$$

is equivalent to the MAP estimate (2.6).

Noteworthy, by defining the energy of the problem via

$$U(\boldsymbol{x}) = -\log(f_{\boldsymbol{x}}(\boldsymbol{x}) \cdot f_{\boldsymbol{y}|\boldsymbol{x}}(\boldsymbol{x}))$$
(2.41)

the problem of minimizing the energy becomes

$$\min_{\boldsymbol{x}} \frac{1}{2\sigma_{n}^{2}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_{2}^{2} - \log f_{\boldsymbol{x}}(\boldsymbol{x}) , \qquad (2.42)$$

which itself for Laplacian prior with variance σ^2 , given by $f_x(x) = \frac{1}{\sqrt{2}\sigma} \exp\left(-\frac{\sqrt{2}}{\sigma}|x|\right)$, yields so-called <u>basis pursuit den</u>oising (BPDN) [CDS01], also known as LASSO [Tib96]⁹

$$\min_{\boldsymbol{x}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_{2}^{2} + \lambda \|\boldsymbol{x}\|_{1}$$
(2.43)

with tuning parameter $\lambda = 2\sqrt{2}\frac{\sigma_n^2}{\sigma}$ and $\|\cdot\|_1$ being the so-called ℓ_1 -norm, for the definition, see Appendix A.

2.3 Structure of the Compressed Sensing Problem

Now the conditional pdf (2.7) is considered, which can be factorized in different ways. With

$$\mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x}) = \frac{1}{Z} \mathbf{f}_{\mathbf{x}}(\mathbf{x}) \cdot \mathbf{f}_{\mathbf{y}|\mathbf{x}}(\mathbf{x}) , \qquad (2.44)$$

where

$$\mathbf{f}_{\mathbf{x}}(\boldsymbol{x}) = \prod_{j=1}^{N} \mathbf{f}_{\mathbf{x}}(x_j)$$
(2.45)

is the signal part,

$$f_{\mathbf{y}|\mathbf{x}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi\sigma_{\mathbf{n}}^2)^M}} \exp\left(-\frac{1}{2\sigma_{\mathbf{n}}^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2\right)$$
$$= \prod_{i=1}^M \frac{1}{\sqrt{2\pi\sigma_{\mathbf{n}}^2}} \exp\left(-\frac{1}{2\sigma_{\mathbf{n}}^2} (y_i - \mathbf{a}_i^\top \mathbf{x})^2\right) \stackrel{\text{def}}{=} \prod_{i=1}^M f_{\mathbf{y}_i|\mathbf{x}}(\mathbf{x}) , \qquad (2.46)$$

the channel part, and $Z = f_y(y)$ serves as normalization, the problem can be split into two parts that are connected via the signal x. Note that signal and channel part can itself be factorized as given.

⁹LASSO stands for least absolute shrinkage and selection operator.

2.3.1 Factor Graph Representations

A graphical representation of the structure of the posterior is given as *factor graph* [KFL01] in Fig. 2.3. In the upper left, the detailed structure of the problem with the factors broken down to the most fragmented factorization, representing $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}) = \frac{1}{Z} \prod_{i=1}^{M} f_{\mathbf{y}_i|\mathbf{x}}(\mathbf{x}) \prod_{j=1}^{N} f_{\mathbf{x}}(x_j)$, is shown. Factors are represented by squares, the *variable nodes* by circles. The connections, represented by lines between the nodes, indicate dependencies between the nodes, i.e., variable nodes connected to a factor node mean that the respective variables are arguments of the factor. One can see that each factor of the channel part is connected to all variables. On the other hand, the signal factors only depend on one corresponding variable. The lower right factor graph shows the splitting according to (2.44), which focuses on the fact that there are only two different kinds of factors. In order to emphasize that the factor nodes here are combined versions of the factors in the upper left, a different representation is used. Furthermore, combined variable nodes are represented by a double circle.

Additionally, mixtures of both extreme cases are possible. The lower left shows a representation that combines the channel parts into $f_{\mathbf{y}|\mathbf{x}}(\mathbf{x})$, but keeps the signal parts as factors, thereby causing distinct variable nodes. Whereas, on the upper right, the opposite is shown with combined signal parts and individual channel parts $f_{\mathbf{y}_i|\mathbf{x}}(\mathbf{x})$ $(i \in \{1, \ldots, M\})$. Noteworthy, in this last case, the variable nodes are combined because the channel parts depend on the entire signal vector \mathbf{x} . According to the denomination we may refer to the respective nodes as channel- or signal-constrained (factor) nodes.

2.3.2 Representatives for the Factors

The approximation of $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$ by the substitute distribution $q(\mathbf{x})$ needs to represent the structure of the problem suitably. To that end, distributions are defined that represent the substitute distribution at certain nodes. The representation for the entire channel part $f_{\mathbf{y}|\mathbf{x}}(\mathbf{x})$ is called $q_c(\mathbf{x})$, whereas if only the dependency of \mathbf{x} on a single observation, i.e., $f_{\mathbf{y}_i|\mathbf{x}}(\mathbf{x})$, is considered the representative is denoted by $q_{c,i}(\mathbf{x})$. The signal part $f_{\mathbf{x}}(\mathbf{x})$ is represented by $q_s(\mathbf{x})$ and the combined variable node by $q_v(\mathbf{x})$, respectively. In contrast to $q_c(\mathbf{x})$ and $q_s(\mathbf{x})$, the representative of the combined variable node $q_v(\mathbf{x})$ cannot be found in the formula for the posterior $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$, indeed its existence is only motivated by the factor graph representation. The task of the (combined) variable node \mathbf{x} is to force equality between the arguments of the factors $f_{\mathbf{y}|\mathbf{x}}(\mathbf{x})$ and $f_{\mathbf{x}}(\mathbf{x})$ in (2.44), i.e., if they are treated individually, one has to make sure that the argument is the same. This behavior can be described by the Dirac delta function, i.e., one can write [LDH⁺07]

$$f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) = \text{const.} \cdot \int f_{\mathbf{x}}(\tilde{\boldsymbol{x}}) f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) \delta(\boldsymbol{x} - \tilde{\boldsymbol{x}}) \, \mathrm{d}\tilde{\boldsymbol{x}} \; . \tag{2.47}$$

This means, the Dirac delta function can be seen as the function of the factor behind the variable node. Nevertheless, we will keep the denomination for variable nodes for the distinguishability to the actual factor nodes. The derivations in subsequent chapters will show that the representative for the variable nodes $q_v(x)$ softens the equality constraint by allowing also variations in the factors.



(a) Detailed view on all individual factors.



(c) Detailed view on variable nodes, combined channel-constrained factor nodes.



(b) Combined variable node, detailed view on channel-constrained factor nodes.



mbined (d) Focus on the two different kinds of factors.

Figure 2.3: Several representations of the factor graph for the conditional pdf of the CS problem (2.7).



(a) Detailed factorization.

(b) Two-part factorization.

Figure 2.4: Factor graphs for different factorizations with respective representatives.

Considering the individual variable nodes and individual signal factor nodes as, e.g., in Fig. 2.3(a), requires also representatives for these (sub)nodes. As these representatives do not depend on all variables anymore, one needs to get rid of them. This is done by *marginaliza-tion* [KFL01, LDH⁺07]. A marginal of any distribution is in general obtained for a variable x_j by integrating (or summing in case of discrete random variables) over all elements of x that are not x_j , i.e., to obtain a marginal for a single variable at a variable node, one needs to compute

$$\mathbf{q}_{\mathbf{v},j}(x_j) = \int \cdots \int \mathbf{q}_{\mathbf{v}}(\boldsymbol{x}) \, \mathrm{d}x_1 \dots \, \mathrm{d}x_{j-1} \, \mathrm{d}x_{j+1} \dots \, \mathrm{d}x_N \,, \quad j \in \{1, \dots, N\} \,. \tag{2.48}$$

Analogously, the marginal at a single signal factor $f_x(x_j)$ is

$$\mathbf{q}_{\mathbf{s},j}(x_j) = \int \cdots \int \mathbf{q}_{\mathbf{s}}(\boldsymbol{x}) \, \mathrm{d}x_1 \dots \, \mathrm{d}x_{j-1} \, \mathrm{d}x_{j+1} \dots \, \mathrm{d}x_N \,, \quad j \in \{1, \dots, N\} \,. \tag{2.49}$$

With the marginals $q_{x,j}(x_j)$ and $q_{s,j}(x_j)$ the connection between x_j and $f_x(x_j)$ can be expressed, as will be seen later in Chapter 3.

Together with $q_{v,j}(x_j)$, the marginals of $q_{c,i}(x)$ can be used to specify the connections between x_j and the channel part $f_{y_i|x}(x)$. They are denoted and calculated by

$$\mathbf{q}_{\mathrm{c},i}(x_j) = \int \cdots \int \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \,\mathrm{d}x_1 \dots \,\mathrm{d}x_{j-1} \,\mathrm{d}x_{j+1} \dots \,\mathrm{d}x_N \,, \quad j \in \{1, \dots, N\} \,. \tag{2.50}$$

In this last case, it is clearly visible that the factor marginals are edge-dependent, as one marginal specifies the connecting edge between corresponding factor and variable node.

To get the denomination of the different marginals straight, two of the factor graph figures are repeated in Fig. 2.4 to show the correspondences. Furthermore, the effect of marginalizing a distribution is shown in Ex. 2.1.

Example 2.1:

As an example, $f_{\mathbf{y}|\mathbf{x}}(\mathbf{x})$ from (2.46) of a two-dimensional CS scenario is considered as a function of \mathbf{x} only. To that end, let us define

$$\mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) \stackrel{\text{def}}{=} \frac{1}{\sqrt{(2\pi\sigma_{n}^{2})^{2}}} \exp\left(-\frac{1}{2\sigma_{n}^{2}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_{2}^{2}\right) , \qquad (2.51)$$



Figure 2.5: Contour plot of $p_y(x)$ of Ex. 2.1 and respective marginals.

where $\boldsymbol{y} = [-0.23, 0.76]^{\top}$ is obtained from compressing the signal $\boldsymbol{x} = [0, 1]^{\top}$, with the sensing matrix

$$\boldsymbol{A} = \begin{bmatrix} 0.75 & -0.25\\ -0.5 & 0.75 \end{bmatrix}$$
(2.52)

and adding Gaussian noise with variance $\sigma_n^2 = 0.2$. The two-dimensional pdf $p_y(x)$ is shown as a contour plot in Fig. 2.5, surrounded by the marginals which are here computed by

$$\mathbf{p}_{\boldsymbol{y},1}(x_1) = \int \mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) \, \mathrm{d}x_2 \tag{2.53}$$

and

$$p_{\boldsymbol{y},2}(x_2) = \int p_{\boldsymbol{y}}(\boldsymbol{x}) \, \mathrm{d}x_1 \;, \tag{2.54}$$

respectively. The integration removes the dependency of the variable that is integrated over. Noteworthy, the pdf $p_y(x)$ has a unique maximum, since the sensing matrix is square. In a scenario, where the signal is actually compressed, e.g., by considering only $y_1 = a_1^T x + n_1$ with $y_1 = -0.23$, $a_1 = [0.75, -0.25]^T$ and $n_1 \sim \mathcal{N}(0, \sigma_n^2)$, the contour plot shows a Gaussian distribution around the line defined by

$$y_1 = -0.23 = \boldsymbol{a}_1^\top \boldsymbol{x} = a_{11}x_1 + a_{12}x_2 = 0.75x_1 - 0.25x_2 \tag{2.55}$$

$$\Leftrightarrow \quad x_2 = -4 \cdot (y_1 - 0.75x_1) = 3x_1 + 0.92 . \tag{2.56}$$



Figure 2.6: Contour plot of $p_{y_1}(x)$ of Ex. 2.1, the line $x_2 = 3x_1 + 0.92$, specifying the maximum of the pdf, and two possible sparse solutions.

All points on that line maximize the pdf and may therefore be seen as solution. Setting either of the signal values x_1 or x_2 to zero, recovers a sparse solution, e.g., $x_1 = 0$ leads to $x_2 = 0.92$. The pdf that corresponds to the single measurement scenario

$$\mathbf{p}_{y_1}(\boldsymbol{x}) = \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{1}{2\sigma_n^2}(y_1 - \boldsymbol{a}_1^{\top}\boldsymbol{x})^2\right)$$
(2.57)

is depicted in Fig. 2.6 as contour plot. Additionally the line following the maximum and the two sparse solutions are shown.

3. Inference on Factor Graphs

The use of factor graphs [Loe04, For01] and other graphical models for inference was exploited in several fields of engineering already. There is to mention the decoding of LDPC codes via Tanner graphs [Tan81], or algebraic codes via Ising-type models [For18], other decoding strategies, such as the BCJR or the Viterbi algorithm [Loe04], multiuser detection in CDMA [GW06] and several aspects in signal processing [LDH⁺07], as, e.g., Kalman filtering or recursive leastsquares algorithms. The underlying algorithm is called *belief propagation* (BP) [Pea88] or *message passing* (MP) and may adopt certain forms; for an overview see [KFL01].

In the following, an optimization problem with constraints that are imposed by the structure of the considered factor graph is set up. Based on that, a form of (Gaussian) message passing for the CS problem with factor graph representation in Fig. 2.3(a) as basis is derived. The resulting algorithm is equal to the one stated in [KMS⁺12, Sec. 3.1], however, the derivation differs from the one considered here.

An optimization-based view has been used in [YFW05] to show the connection between MP and the so-called *Bethe free energy* [Bet35]. A similar derivation is given in [HOW⁺05], where approximate inference with expectation constraints is considered for general inference problems. An overview over signal processing strategies encompassing, next to MP and the Bethe approach, also AMP and the EP framework for general inference is given in [PSC⁺15]. Here, we follow the approaches given in [HOW⁺05, YFW05], but focus on the CS problem.

Noteworthy, the optimization-based derivation reveals that the notion of *extrinsic* [BGT93, HOP96] is inherently given in this approach to the solution of the problem.

3.1 Structure of the Factor Graph

In Chapter 2 it has been pointed out that the high dimensionality of x makes it difficult to recover the signal from the observations. As a first approach, the computation is simplified as much as possible by breaking down the structure to its basic components. To that end, the representation for the CS problem that factorizes the problem as much as possible, cf. Fig. 2.3(a), is regarded. With this fine-granular view, the multi-dimensional estimation is split

into several one-dimensional ones, which simplifies the estimations the most. In the following, this structure will be used to derive a suitable representation of the approximation for the factorization.

3.1.1 Factorization of the Substitute Distribution

First, it is important to find a suitable form for the substitute q(x) that should be worked with. The substitute distribution q(x), which is used to approximate the posterior $f_{x|y}(x)$, should inherit the structure of the factor graph in Fig. 2.3(a). To that end, the corresponding factorization is considered as a function of x. In order to allow variations in the factors, the product is expanded by functions of x for each factor, i.e.,

$$f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) = \text{const.} \cdot \prod_{j=1}^{N} f_{\mathbf{x}}(x_j) \prod_{i=1}^{M} f_{\mathbf{y}_i|\mathbf{x}}(\boldsymbol{x})$$

$$= \text{const.} \cdot \prod_{j=1}^{N} f_{\mathbf{x}}(x_j) \prod_{i=1}^{M} f_{\mathbf{y}_i|\mathbf{x}}(\boldsymbol{x}) \cdot \frac{\prod_{j=1}^{N} s_j(x_j) \prod_{i=1}^{M} c_i(\boldsymbol{x})}{\prod_{j=1}^{N} s_j(x_j) \prod_{i=1}^{M} f_{\mathbf{x}_i}(\boldsymbol{x})}$$

$$= \text{const.} \cdot \frac{\prod_{j=1}^{N} f_{\mathbf{x}}(x_j) s_j(x_j) \prod_{i=1}^{M} f_{\mathbf{y}_i|\mathbf{x}}(\boldsymbol{x}) c_i(\boldsymbol{x})}{\prod_{j=1}^{N} s_j(x_j) \prod_{i=1}^{M} c_i(\boldsymbol{x})} .$$
(3.1)

The association of the factors and the respective functions may now be interpreted as the corresponding representatives from Sec. 2.3.2, i.e.,

$$\mathbf{q}_{c,i}(\boldsymbol{x}) = \mathbf{f}_{\mathbf{y}_i | \mathbf{x}}(\boldsymbol{x}) \cdot c_i(\boldsymbol{x}) , \qquad i \in \{1, \ldots, M\} , \qquad (3.2)$$

$$q_{s,j}(x_j) = f_x(x_j) \cdot s_j(x_j) , \qquad j \in \{1, \dots, N\} .$$
 (3.3)

Assuming that the function for the channel-constrained approximation split as follows

$$c_i(\boldsymbol{x}) = \prod_{j=1}^N c_{ij}(x_j) , \qquad i \in \{1, \dots, M\} ,$$
 (3.4)

one can write

$$\mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x}) = \text{const.} \cdot \frac{\prod_{j=1}^{N} \mathbf{q}_{\mathrm{s},j}(x_j) \prod_{i=1}^{M} \mathbf{q}_{\mathrm{c},i}(\mathbf{x})}{\prod_{j=1}^{N} \left(s_j(x_j) \prod_{i=1}^{M} c_{ij}(x_j) \right)} .$$
(3.5)

This way, the sorting in the denominator is chosen such that functions of the same variable x_j ($j \in \{1, ..., N\}$) are packed together. These clustered variable-dependent functions can now be set in relation with the representatives $q_{v,j}(x_j)$ ($j \in \{1, ..., N\}$) that were chosen for the variable nodes.

The expansion of the product by the introduced deviations may be visualized by attaching them as factors (in numerator and denominator) to the corresponding variable nodes, cf. Fig. 3.1. The approximations $q_{c,i}(x)$ and $q_{s,j}(x_j)$ are then obtained by grouping the factors correspondingly.

Recall, that the distributions representing the variable nodes are Dirac delta functions, cf. (2.47), i.e., the equality between the arguments of the different functions is obtained by the


Figure 3.1: Visualization of the grouping of the factors to allow for variations.

integration over respective Dirac functions. The central box in Fig. 3.1 can, thus, be described by

$$\frac{\text{const.}}{s_j(x_j)\prod_{i=1}^M c_{ij}(x_j)} = \text{const.} \cdot \int \cdots \int \frac{\prod_{i=1}^M \delta(x_j - x_{i,j})}{s_j(x_j)\prod_{i=1}^M c_{ij}(x_{i,j})} \, \mathrm{d}x_{1,j} \dots x_{M,j} , \qquad (3.6)$$

where the slack variables $x_{1,j}, \ldots, x_{M,j}$ are introduced to indicate the arguments of the different deviations. The argument of $s_j(x_j)$ is chosen to be x_j for convenience, so that the usual argument is recovered after the integration.

It is important to note that the M-fold integral on the right-hand side of (3.6) does not describe the effect of the variable node on other factors as in (2.47) but rather represents an internal behavior. In Chapter 2, variable nodes have been defined by their (integral) effect of the Dirac delta function (cf. (2.47)) on other factors, i.e., they do not have an internal behavior. Therefore, the internal behavior that is present on the right-hand side in (3.6) needs to be compensated. If the integral were separable, it would constitute an M-fold multiplication, i.e., the multiplication of M functions. However, the marginal $q_{v,j}(x_j)$ for the variable node x_j represents one function. When drawing the connection between marginal $q_{v,j}(x_j)$ and (3.6), this inherent multiplication is therefore compensated by taking the Mth root, i.e.,

$$\mathbf{q}_{\mathbf{v},j}(x_j) = \text{const.} \cdot \sqrt[M]{s_j(x_j) \prod_{i=1}^M c_{ij}(x_j)}, \qquad j \in \{1, \dots, N\}.$$
 (3.7)

Note that the denominator is represented directly for better readability. Interestingly, the approach of allowing variations in the factors led to the effect of the Dirac delta function being replaced by corresponding variations, which means that the equality between the arguments of different factors is not as strict as before. This way, the entire substitute reads [YFW05]

$$\mathbf{q}(\boldsymbol{x}) = \text{const.} \cdot \frac{\prod_{j=1}^{N} \mathbf{q}_{\mathrm{s},j}(x_j) \prod_{i=1}^{M} \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x})}{\left(\prod_{j=1}^{N} \mathbf{q}_{\mathrm{v},j}(x_j)\right)^{M}} .$$
(3.8)

Note that since the interest lies in the argument x of the function, the constant factor in the formula does not harm the processing [LDH⁺07]; it is therefore not stated in subsequent derivations.



Figure 3.2: Factor graph and representatives for Ex. 2.1 with one observation y_1 .

This factorization represents the most fine-granular view possible, as every individual factor and variable node is represented by its own marginal distribution. According to [YFW05], this factorization is the underlying principle of the Bethe free energy, that was shown to be connected to the BP algorithm. The procedure amounts to the simple rule that for each variable x_j ($j \in \{1, ..., N\}$) that occurs "too often" in the numerator, one representative $q_{v,j}(x_j)$ occurs in the denominator as compensation. In the given case, each variable x_j is part of $q_{s,j}(x_j)$ as well as the M channel-constrained representatives, i.e., the node x_j has M + 1 neighbors. Therefore, (M + 1) - 1 = M instances of $q_{v,j}(x_j)$ appear in the denominator.

An example computation is conducted in Ex. 3.1.

Example 3.1: _

For Ex. 2.1, when considering only one observation y_1 for the two signal components x_1 and x_2 , the factor graph, shown in Fig. 3.2, consists of three factors and two variable nodes, which are connected by the single channel-constrained factor, i.e., the observation. Instead of the factors, the representatives of the notation above are placed at the respective positions. Since the multiplication of $q_{s,1}(x_1)$, $q_{s,2}(x_1)$, and $q_{c,1}(x_1, x_2)$ contains each of the variables x_1 and x_2 two times, one needs to divide by $q_{v,1}(x_1)$ and $q_{v,2}(x_2)$, which results in an overall approximation

$$\mathbf{q}(\boldsymbol{x}) = \frac{\mathbf{q}_{\mathrm{s},1}(x_1)\mathbf{q}_{\mathrm{s},2}(x_2)\mathbf{q}_{\mathrm{c},1}(x_1, x_2)}{\mathbf{q}_{\mathrm{v},1}(x_1)\mathbf{q}_{\mathrm{v},2}(x_2)} \,. \tag{3.9}$$

3.1.2 The Moment-Matching Constraint

As already mentioned, a suitable solution q(x) must represent the structure of the problem correctly. In this fine-granular view, this means that the marginals of respective nodes must be consistent, when following the connecting edges. Consistency, here, means that there is no contradiction to neighboring factor nodes, when marginalizing the respective pdfs to the variable under consideration. In [HOW⁺05] this constraint is called *local consistency condition*. For the marginals at hand, this reads

$$\mathbf{q}_{\mathbf{c},i}(x_j) \stackrel{!}{=} \mathbf{q}_{\mathbf{v},j}(x_j) , \quad j \in \{1, \ldots, N\} , \quad i \in \{1, \ldots, M\} ,$$
 (3.10)

$$\mathbf{q}_{\mathrm{s},j}(x_j) \stackrel{!}{=} \mathbf{q}_{\mathrm{v},j}(x_j) , \quad j \in \{1, \dots, N\} .$$
 (3.11)

The resulting constrained optimization problem leads to the famous *message passing* (MP) algorithm, also known as (loopy) *belief propagation* (BP)¹ [Pea88]. The algorithm propagates pdfs as messages to update the so-called *beliefs*, which are approximates for the exact marginals of the signal x [YFW05].

Since this is a rather unhandy representation, this step is skipped and it is directly worked with moments instead of the distributions, which is equal to restricting the marginals to some *exponential family* [Bro86]². This is a class of distributions that is entirely determined by certain moments; the respective moments are specified by the so-called *sufficient statistics* g(x). For a given sufficient statistics, an exponential family is parameterized by the so-called *natural parameters*, which is here denoted by θ . Noteworthy, there is a direct dependency between the moments specified by the sufficient statistics and the natural parameters; examples are given in Appendix D. An exponential family has, consequently, the beneficial property of allowing to transfer the computations w.r.t. distributions into a parameter space, i.e., instead of handling the distribution, one can compute estimates and process them. Restricting to the Gaussian distribution, which requires only first and second order moments for its statistic to be specified, will lead to so-called *Gaussian BP* [KMS⁺12]. For the derivation thereof, the sufficient statistics g(x) of the respective exponential family are fixed to (D.11), i.e.,

$$\boldsymbol{g}(\boldsymbol{x}) = [\boldsymbol{g}_{\boldsymbol{\lambda}}(\boldsymbol{x})^{\top}, \, \boldsymbol{g}_{\boldsymbol{\Lambda}}(\boldsymbol{x})^{\top}]^{\top}$$
(3.12)

with

$$\boldsymbol{g}_{\boldsymbol{\lambda}}(\boldsymbol{x}) = [x_1, \ldots, x_N]^{\top}, \quad \boldsymbol{g}_{\boldsymbol{\Lambda}}(\boldsymbol{x}) = -\frac{1}{2}[x_1^2, \ldots, x_N^2]^{\top},$$
 (3.13)

and respective natural parameter $\boldsymbol{\theta} = [\boldsymbol{\lambda}^{\top}, \boldsymbol{\Lambda}^{\top}]^{\top}$. In order to specify the moment of a single variable, we may use the notation $\boldsymbol{g}(x_j) = [g_{\lambda_j}(x_j), g_{\Lambda_j}(x_j)]^{\top} = [x_j, -x_j^2/2]^{\top}$.

The respective constraints, ensuring the representation of the structure of the CS posterior $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x})$, is called *weak consistency condition* in [HOW⁺05] and can be stated by using the marginals and the specified sufficient statistics for $j \in \{1, ..., N\}$

$$\mathbf{E}_{\mathsf{x}_{j} \sim \mathfrak{q}_{\mathsf{v},j}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} \stackrel{!}{=} \mathbf{E}_{\mathsf{x}_{j} \sim \mathfrak{q}_{\mathsf{c},i}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} , \quad i \in \{1, \ldots, M\} , \qquad (3.14)$$

respectively,

$$\mathbb{E}_{\mathsf{x}_{j} \sim \mathfrak{q}_{\mathsf{v},j}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} \stackrel{!}{=} \mathbb{E}_{\mathsf{x}_{j} \sim \mathfrak{q}_{\mathsf{s},j}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} .$$
(3.15)

These requirements will be referred to as the moment-matching constraints.

3.2 Constrained Optimization

In the following, the task of approximating $f_{x|y}(x)$ suitably is described as constrained optimization problem for the given structure and the stationary points of the respective Lagrangian is sought, which eventually is used to derive a recovery algorithm for the CS problem.

¹In the given case the BP algorithm is loopy because the factor graph is not a tree, i.e., has cycles.

²The necessary knowledge about exponential families is summarized in Appendix D.

To that end, the idea of minimizing the Kullback–Leibler divergence, which was motivated in Sec. 2.2.2, is seized and combined with the moment-matching constraints (3.14) and (3.15) of the fine-granular view on the factor graph. This way, a constrained optimization problem is obtained that reads

$$\begin{array}{ll} \text{minimize} & \mathrm{D}_{\mathrm{KL}} \Big(\mathbf{q}(\boldsymbol{x}) \mid| \, \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) \Big) \\ \text{subject to} & \mathrm{E}_{\mathsf{x}_{j} \sim \mathsf{q}_{\mathsf{v},j}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} - \mathrm{E}_{\mathsf{x}_{j} \sim \mathsf{q}_{\mathsf{c},i}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} = \boldsymbol{0} \,, \qquad j \in \{1, \ldots, N\} \,, \, i \in \{1, \ldots, M\} \,, \\ & \mathrm{E}_{\mathsf{x}_{j} \sim \mathsf{q}_{\mathsf{v},j}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} - \mathrm{E}_{\mathsf{x}_{j} \sim \mathsf{q}_{\mathsf{s},j}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} = \boldsymbol{0} \,, \qquad j \in \{1, \ldots, N\} \,, \\ & \int \mathbf{q}_{\mathsf{v},j}(x_{j}) \, \mathrm{d}x_{j} = 1 \,, \quad \int \mathbf{q}_{\mathsf{s},j}(x_{j}) \, \mathrm{d}x_{j} = 1 \,, \quad j \in \{1, \ldots, N\} \,, \\ & \int \mathbf{q}_{\mathsf{c},i}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = 1 \,, \qquad \qquad i \in \{1, \ldots, M\} \,. \end{array}$$

3.2.1 Cost Function for Compressed Sensing

Before considering the optimization itself, the Kullback–Leibler divergence to be minimized is stated, utilizing the structure imposed by CS. Inserting the factorization (3.8) into the Kullback–Leibler divergence yields a certain form, because the term disintegrates to a sum of parts depending solely on the factors, respectively variable marginals due to the normalization property of pdfs. The result to be minimized in the above optimization problem reads

$$D_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) \mid\mid \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})) = \sum_{j=1}^{N} \int \mathbf{q}_{\mathrm{s},j}(x_{j}) \log \mathbf{q}_{\mathrm{s},j}(x_{j}) \,\mathrm{d}x_{j} + \sum_{i=1}^{M} \int \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \log \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \,\mathrm{d}x_{i} - \sum_{j=1}^{N} \int \mathbf{q}_{\mathrm{s},j}(x_{j}) \log \mathbf{f}_{\mathrm{x}}(x_{j}) \,\mathrm{d}x_{j} - \sum_{i=1}^{M} \int \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \log \mathbf{f}_{\mathrm{y}_{i}|\mathbf{x}}(\boldsymbol{x}) \,\mathrm{d}x_{i} - M \sum_{j=1}^{N} \int \mathbf{q}_{\mathrm{v},j}(x_{j}) \log \mathbf{q}_{\mathrm{v},j}(x_{j}) \log \mathbf{q}_{\mathrm{v},j}(x_{j}) \,\mathrm{d}x_{j} \,.$$
(3.17)

This form of the Kullback–Leibler divergence is called *Bethe free energy* in [YFW05].

3.2.2 Stationary Points of the Lagrangian

Constrained optimization problems can be dealt with by considering the so-called *Lagrangian*, which combines the optimization task and the constraints in one function(al), cf. Appendix B. As motivated in Appendix B, the optimizer of the problem needs to fulfill a stationary condition of the Lagrangian, which means the stationary points of the Lagrangian are of interest for us to be achieved. By introducing the Lagrangian Multipliers $\nu_{c,ij}$ and $\nu_{s,jj}$ for the moment-matching constraints, as well as ν_j , $\nu_{c,i}$, and $\nu_{s,j}$ as normalization constraint to ensure valid marginals, the Lagrangian reads (omitting the Lagrangian multipliers in the argument

for brevity)

$$\mathcal{L}(\mathbf{q}(\boldsymbol{x})) = D_{\mathrm{KL}}\left(\mathbf{q}(\boldsymbol{x}) \mid| \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})\right) + \sum_{i=1}^{M} \sum_{j=1}^{N} \boldsymbol{\nu}_{\mathrm{c},ij}^{\top} \left(\mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathrm{v},j}} \{\boldsymbol{g}(\mathbf{x}_{j})\} - \mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathrm{c},i}} \{\boldsymbol{g}(\mathbf{x}_{j})\}\right) \\ + \sum_{j=1}^{N} \boldsymbol{\nu}_{\mathrm{s},jj}^{\top} \left(\mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathrm{v},j}} \{\boldsymbol{g}(\mathbf{x}_{j})\} - \mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathrm{s},j}} \{\boldsymbol{g}(\mathbf{x}_{j})\}\right) \\ + \sum_{j=1}^{N} \left(\boldsymbol{\nu}_{j} \left(\int \mathbf{q}_{\mathrm{v},j}(x_{j}) \, \mathrm{d}x_{j} - 1 \right) + \boldsymbol{\nu}_{\mathrm{s},j} \left(\int \mathbf{q}_{\mathrm{s},j}(x_{j}) \, \mathrm{d}x_{j} - 1 \right) \right) \\ + \sum_{i=1}^{M} \boldsymbol{\nu}_{\mathrm{c},i} \left(\int \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 \right) .$$
(3.18)

The stationary points can be obtained by differentiating the Lagrangian with respect to its parameters and setting the term to zero. When differentiating w.r.t. one of the Lagrangian multipliers, this leads to the respective constraint, while differentiation w.r.t. the marginals specify the form of the respective marginal, because the stationarity constraint on the derivative fixates the form of the marginal in the optimal point. This form will help us to derive algorithms for the solution to the problem, since the stationary points are the targets to aim for. In the latter case, the differentiation is again a functional derivative, for more details see Appendix C.

The derivatives, that subsequently are set to zero, read

$$\frac{\partial}{\partial \mathbf{q}_{\mathbf{v},j}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = -M \log \mathbf{q}_{\mathbf{v},j}(x_j) + \left(\boldsymbol{\nu}_{\mathbf{s},jj}^\top + \sum_{i=1}^M \boldsymbol{\nu}_{\mathbf{c},ij}^\top\right) \boldsymbol{g}(x_j) + \text{const.} \stackrel{!}{=} 0, \quad (3.19)$$

$$\frac{\partial}{\partial \mathbf{q}_{\mathrm{c},i}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = \log \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) - \log \mathbf{f}_{\mathbf{y}_i|\mathbf{x}}(\boldsymbol{x}) - \sum_{j=1}^{N} \boldsymbol{\nu}_{\mathrm{c},ij}^{\top} \boldsymbol{g}(x_j) + \text{const.} \stackrel{!}{=} 0 , \qquad (3.20)$$

$$\frac{\partial}{\partial \mathbf{q}_{\mathbf{s},j}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = \log \mathbf{q}_{\mathbf{s},j}(x_j) - \log \mathbf{f}_{\mathbf{x}}(x_j) - \boldsymbol{\nu}_{\mathbf{s},jj}^{\top} \boldsymbol{g}(x_j) + \text{const.} \stackrel{!}{=} 0 , \qquad (3.21)$$

where parts that do not depend on any element of x are summarized to a constant addend.

The resulting terms, solved for the marginals, read [HOW⁺05]

$$\mathbf{q}_{\mathbf{v},j}(x_j) = \frac{1}{Z_{\mathbf{v},j}} \exp\left(\frac{1}{M} (\boldsymbol{\nu}_{\mathbf{s},jj} + \sum_{i=1}^{M} \boldsymbol{\nu}_{\mathbf{c},ij})^{\mathsf{T}} \boldsymbol{g}(x_j)\right) , \qquad (3.22)$$

$$\mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) = \frac{1}{Z_{\mathrm{c},i}} \mathbf{f}_{\mathrm{y}_i|\mathbf{x}}(\boldsymbol{x}) \exp\left(\sum_{j=1}^{N} \boldsymbol{\nu}_{\mathrm{c},ij}^{\top} \boldsymbol{g}(x_j)\right) , \qquad (3.23)$$

$$\mathbf{q}_{\mathbf{s},j}(x_j) = \frac{1}{Z_{\mathbf{s},j}} \mathbf{f}_{\mathbf{x}}(x_j) \exp\left(\boldsymbol{\nu}_{\mathbf{s},jj}^{\top} \boldsymbol{g}(x_j)\right) , \qquad (3.24)$$

where $Z_{v,j}$, $Z_{c,i}$, and $Z_{s,j}$ are introduced to fulfill the normalization constraint.

As expected, the marginals at the stationary points are members of exponential families. Noteworthy, the Lagrangian multipliers turn out to work as natural parameters. We define

$$\boldsymbol{\theta}_{j} \stackrel{\text{\tiny def}}{=} \frac{1}{M} (\boldsymbol{\nu}_{\mathrm{s},jj} + \sum_{i=1}^{M} \boldsymbol{\nu}_{\mathrm{c},ij}) , \qquad (3.25)$$



Figure 3.3: Factor graph for compressed sensing and edge-dependent multipliers.

in order to specify the natural parameter at the variable node. The computational connection of the variable node to the factor nodes, i.e., the marginals $q_{c,i}(\boldsymbol{x})$, respectively $q_{s,j}(x_j)$ is given by the Lagrangian multipliers $\boldsymbol{\nu}_{c,ij}$, respectively $\boldsymbol{\nu}_{s,jj}$ since the multipliers appear in these connected marginals as well. Hence, each multiplier belongs to one edge, i.e., they are *edge-dependent*, cf. Fig. 3.3. The variable node marginal $q_{v,j}(x_j)$ is entirely specified by the respective Lagrangian multipliers and specifies a Gaussian distribution with the given sufficient statistics $\boldsymbol{g}(\boldsymbol{x})$. The factor marginals $q_{c,i}(\boldsymbol{x})$ and $q_{s,j}(x_j)$, however, include also the factors $f_{y_i|\mathbf{x}}(\boldsymbol{x})$ and $f_{\mathbf{x}}(x_j)$, respectively.

3.3 Algorithms

The stationary points derived above can be seen as targets to aim for, because they solve the given problem. A recovery algorithm for CS should therefore end up with the forms for the marginals (3.22)–(3.24) above. In order to get there, the algorithm needs to make use of the structural connections incorporated in the Lagrangian multipliers. Since these multipliers are natural parameters of an exponential family, they express the entire knowledge about the corresponding distribution. As they are, furthermore, edge-dependent, they can be seen as "messages" to be broadcasted to other nodes. The connection between natural parameters and moments of exponential families, cf. Appendix D.2, thus, makes it possible to shift the processing from the distributions to the parameter domain, which comes in handy, as the interest lies in the conditional mean (2.9). In the following, it is important to find out, how the processing is done. At first, the moments that are computed at the factor nodes are considered.

3.3.1 Expectations

Since the above specified sufficient statistic represents a Gaussian distribution, the natural parameters of $q_{s,j}(x_j)$ can be identified as $\boldsymbol{\nu}_{s,jj} = [\tilde{x}_{s,j}/\tilde{\sigma}_{s,j}^2, 1/\tilde{\sigma}_{s,j}^2]^{\top}$ with $\tilde{x}_{s,j}$ and $\tilde{\sigma}_{s,j}^2$ being mean and variance of the Gaussian distribution. With this identification, it is obvious that the

moments to be computed relate to the non-linear estimators (2.19) and (2.20), i.e.,

$$\mathbf{E}_{\mathsf{x}_{j}\sim\mathsf{q}_{\mathrm{s},j}}\left\{g_{\lambda_{j}}(\mathsf{x}_{j})\right\} = \mathbf{E}_{\mathsf{x}_{j}}\left\{\mathsf{x}_{j} \mid \tilde{x}_{\mathrm{s},j}\right\} = m_{\mathrm{s},j} , \qquad (3.26)$$

$$\mathbf{E}_{\mathsf{x}_{j}\sim\mathsf{q}_{\mathrm{s},j}}\left\{g_{A_{j}}(\mathsf{x}_{j})\right\} = -\frac{1}{2}\mathbf{E}_{\mathsf{x}_{j}}\left\{\mathsf{x}_{j}^{2} \mid \tilde{x}_{\mathrm{s},j}\right\} = -\frac{1}{2}\left(\sigma_{\mathrm{s},j}^{2} + m_{\mathrm{s},j}^{2}\right) .$$
(3.27)

For the channel-constrained factor, this is similar. First note that by stacking, it is possible to define a new natural parameter, such that

$$\sum_{j=1}^{N} \boldsymbol{\nu}_{\mathrm{c},ij}^{\top} \boldsymbol{g}(x_j) = \boldsymbol{\nu}_{\mathrm{c},i}^{\top} \boldsymbol{g}(\boldsymbol{x}) .$$
(3.28)

Using the identification

$$\boldsymbol{\nu}_{\mathrm{c},i} = [m_{\mathrm{x},1}/\sigma_{\mathrm{x},1}^2, \dots, m_{\mathrm{x},N}/\sigma_{\mathrm{x},N}^2, 1/\sigma_{\mathrm{x},1}^2, \dots, 1/\sigma_{\mathrm{x},N}^2]^\top , \qquad (3.29)$$

the moment computations (for the entire vector) relate to (2.13) and (2.14) by

$$\mathbf{E}_{\mathbf{x} \sim \mathbf{q}_{\mathrm{c},i}} \{ \boldsymbol{g}_{\boldsymbol{\lambda}}(\mathbf{x}) \} = \mathbf{E}_{\mathbf{x}} \{ \mathbf{x} \mid y_i \} = \boldsymbol{m}_{\mathrm{c},i} , \qquad (3.30)$$

$$\mathbf{E}_{\mathbf{x} \sim \mathbf{q}_{\mathrm{c},i}} \{ \boldsymbol{g}_{\boldsymbol{\Lambda}}(\mathbf{x}) \} = -\frac{1}{2} \mathbf{diag} \left(\mathbf{E}_{\mathbf{x}} \{ \mathbf{x} \mathbf{x}^{\top} \mid y_i \} \right) = -\frac{1}{2} \mathbf{diag} \left(\boldsymbol{\varPhi}_{\mathrm{c},i} + \boldsymbol{m}_{\mathrm{c},i} \boldsymbol{m}_{\mathrm{c},i}^{\top} \right).$$
(3.31)

The operator $\operatorname{diag}(\cdot)$ shall denote the vector of diagonal entries of the matrix, which is input to the operator. Otherwise, if a single element with an index, usually j is inserted, it shall denote the diagonal matrix with diagonal entries specified by the indexed input variable, e.g., if $j \in \{1, \ldots, N\}$ then the resulting matrix is a $N \times N$ matrix.

The estimates for a single variable, thus, relate to the respective jth components in (2.16) and (2.17) via

$$\mathbf{E}_{\mathsf{x}_{j}\sim\mathsf{q}_{\mathrm{c},i}}\left\{g_{\lambda_{j}}(\mathsf{x}_{j})\right\} = m_{\mathrm{c},i,j} , \qquad (3.32)$$

$$E_{x_{j} \sim q_{c,i}} \left\{ g_{\Lambda_{j}}(x_{j}) \right\} = -\frac{1}{2} \left(\sigma_{c,i,j}^{2} + m_{c,i,j}^{2} \right) .$$
(3.33)

As intended, all estimations are one-dimensional for the given consideration.

3.3.2 Projection to Exponential Family

Having the moments computed at the factor nodes, it is necessary to find out, how to process them for the exchange between the nodes, i.e., we have to get back to the edges. Since the Lagrangian multipliers $\nu_{c,ij}$ and $\nu_{s,jj}$ specify the connections in the factor graph, they can be interpreted as "messages" on an edge. To that end, the two possible sending directions on the edge have to be distinguished. Consider the variable node with $q_{v,j}(x_j)$ in (3.22). Here, as already mentioned, the distribution is entirely specified by the respective Lagrangian multipliers, i.e., they tell something about the distribution and can therefore be interpreted as *outgoing* messages. For the distinction to the outgoing messages, the *incoming* messages are denoted by $\tilde{\nu}_{c,ij}$ and $\tilde{\nu}_{s,jj}$, the denomination is exemplary shown for one node in Fig. 3.4.



Figure 3.4: Detailed view on the factor graph for CS with edge- and direction-dependent messages.

In order to close the gap between the moments on the nodes and the natural parameters that serve as messages on the edges, we have a look at the moment-matching constraints (3.14) and (3.15). These were enforced between factor nodes and connected variable nodes. Since the marginals at the variable nodes represent members of an exponential family, i.e., a Gaussian distribution for the sufficient statistics g(x) as specified above, the moments can be mapped according to Ex. D.2 from Appendix D to the natural parameter of the marginal at the variable node to obtain for a channel-constrained factor node ($j \in \{1, ..., N\}$)

$$\boldsymbol{\theta}_{\mathrm{c},i,j} = [m_{\mathrm{c},i,j} / \sigma_{\mathrm{c},i,j}^2, \ 1 / \sigma_{\mathrm{c},i,j}^2]^\top, \quad i \in \{1, \dots, M\},$$
(3.34)

respectively,

$$\boldsymbol{\theta}_{\mathrm{s},j} = [m_{\mathrm{s},j} / \sigma_{\mathrm{s},j}^2, \, 1 / \sigma_{\mathrm{s},j}^2]^\top \,, \tag{3.35}$$

for a signal-constrained factor node.

3.3.3 Processing for Stationarity

The mapping in (3.34) and (3.35) above shows a way how to represent the signal component x_j after the corresponding processing. However, it does not state what messages have to be sent on the respective edges. The edges are currently described by messages that pass it in the two possible directions. For example in Fig. 3.4 the edge between variable x_j and (signal-constrained) factor $f_x(x_j)$ transports the messages $\tilde{\nu}_{s,jj}$ and $\nu_{s,jj}$. Since this edge is connected to the signal component x_j , it should provide knowledge about this variable. As already mentioned, this knowledge is (after processing at the factor node) given by $\theta_{s,j}$. In the following, the connection between these three parameters is sought.

To that end, consider the following partitioning of the posterior

$$\mathbf{f}_{\mathbf{x}|\mathbf{y}}(\mathbf{x}) = \text{const.} \cdot \mathbf{f}_{\mathbf{x}}(x_j) \cdot \left(\prod_{\substack{j'=1\\j'\neq j}}^{N} \mathbf{f}_{\mathbf{x}}(x_{j'}) \prod_{i=1}^{M} \mathbf{f}_{\mathbf{y}_i|\mathbf{x}}(\mathbf{x})\right) , \qquad (3.36)$$

i.e., it is split into the factor $f_x(x_j)$ on the one side and everything else on the other side. This extracts the edge under consideration as the combining element between the two parts. It also



Figure 3.5: Depiction of the partitioning of the factor graph with the messages between the parts and the representation of the knowledge $\theta_{s,j}$ on the connecting edge.

shows that a message from $f_x(x_j)$ must represent knowledge from the respective factor, while a message towards $f_x(x_j)$ must carry the knowledge from everywhere else. Obviously, the multiplication of both parts yields (up to a constant factor) the posterior again. The splitting is displayed in Fig. 3.5 with $f_x(x_j)$ on the left, the other part on the right, and the corresponding representation of the knowledge $\theta_{s,j}$, as well as the messages in between.

In this consideration, the interest lies in the variable x_j only, which means that mathematically

$$\mathbf{f}_{\mathbf{x}_{j}|\mathbf{y}}(x_{j}) \stackrel{\text{def}}{=} \int \cdots \int \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) \, \mathrm{d}x_{1} \dots \mathrm{d}x_{j-1} \mathrm{d}x_{j+1} \dots \mathrm{d}x_{N}$$
(3.37)

should be approximated. Since $f_x(x_j)$ does not carry an argument that is used in the integral, it is

$$\mathbf{f}_{\mathsf{x}_{j}|\mathbf{y}}(x_{j}) = \text{const.} \cdot \mathbf{f}_{\mathsf{x}}(x_{j}) \cdot \int \cdots \int \prod_{\substack{j'=1\\j' \neq j}}^{N} \mathbf{f}_{\mathsf{x}}(x_{j'}) \prod_{i=1}^{M} \mathbf{f}_{\mathsf{y}_{i}|\mathbf{x}}(\boldsymbol{x}) \, \mathrm{d}x_{1} \dots \, \mathrm{d}x_{j-1} \mathrm{d}x_{j+1} \dots \, \mathrm{d}x_{N} ,$$
(3.38)

i.e., the same separation as before can be obtained. This can now be set into relation with the messages. The messages used here are natural parameters specifying a member of an exponential family. The distributions corresponding to the messages therefore express approximations for the two separated parts. Following the thought that the message $\tilde{\nu}_{s,jj}$ must represent the factor $f_x(x_j)$ and $\nu_{s,jj}$ the other part leads to an entire approximation, which reads

$$f_{x_j|y}(x_j) \approx \text{const.} \cdot \exp\left(\tilde{\boldsymbol{\nu}}_{s,jj}^{\top} \boldsymbol{g}(x_j)\right) \cdot \exp\left(\boldsymbol{\nu}_{s,jj}^{\top} \boldsymbol{g}(x_j)\right)$$
 (3.39)

The knowledge about x_j that is associated to this edge is, as already mentioned, given by $\theta_{s,j}$, i.e.,

$$\mathbf{f}_{\mathbf{x}_j|\mathbf{y}}(x_j) \approx \text{const.} \cdot \exp\left(\boldsymbol{\theta}_{\mathbf{s},j}^{\top} \boldsymbol{g}(x_j)\right)$$
 (3.40)

Hence, the approximation of $f_{x_j|y}(x_j)$, which should fit to the respective incoming and outgoing messages, is given by the exponential family carrying $\theta_{s,j}$ as natural parameter. The multiplication of the members of this exponential family transforms to an addition of the natural parameters, meaning that the natural parameter $\theta_{s,j}$ representing the entire knowledge is given by the summation of the natural parameters describing the messages, i.e., [LDH⁺07, Eqs. (54), (55)]

$$\boldsymbol{\theta}_{\mathrm{s},j} = \tilde{\boldsymbol{\nu}}_{\mathrm{s},jj} + \boldsymbol{\nu}_{\mathrm{s},jj} , \quad j \in \{1, \dots, N\} .$$
(3.41)

Note that in [LDH⁺07] this connection is derived from the *sum-product rule* [KFL01] for Gaussian distributions, which is one of the instances of a message passing algorithm.

Neglecting the loops in the graph, this holds similarly for the channel-constrained factors $f_{y_i|x}(x)$, i.e.,

$$\boldsymbol{\theta}_{c,i,j} = \tilde{\boldsymbol{\nu}}_{c,ij} + \boldsymbol{\nu}_{c,ij}, \quad j \in \{1, \dots, N\}, \quad i \in \{1, \dots, M\}.$$
 (3.42)

Having these connections, it is now possible to turn to the necessary processing in the algorithm. To that end, recall that the stationary points of the optimization problem are already given in (3.22)–(3.24). These are the goals an iterative algorithm should achieve. In order to derive the processing for the algorithms, it is therefore suitable to stick to updates that fulfill these equations, i.e., stay in a stationary point once reached. In particular (3.25), i.e.,

$$\boldsymbol{\theta}_{j} = \frac{1}{M} \left(\boldsymbol{\nu}_{\mathrm{s},jj} + \sum_{i=1}^{M} \boldsymbol{\nu}_{\mathrm{c},ij} \right) , \quad j \in \{1, \ldots, N\} , \qquad (3.43)$$

is of importance for that matter as it encodes the structure of the factor graph in itself. In a stationary point, the representation of x_j given in θ_j from (3.25) should not change with the processing, which means that eventually

$$\boldsymbol{\theta}_{j} \stackrel{!}{=} \boldsymbol{\theta}_{\mathrm{s},j} \stackrel{!}{=} \boldsymbol{\theta}_{\mathrm{c},1,j} \stackrel{!}{=} \cdots \stackrel{!}{=} \boldsymbol{\theta}_{\mathrm{c},M,j} , \quad j \in \{1, \ldots, N\} , \qquad (3.44)$$

must hold. With this equality, (3.41) and (3.42) can be inserted into the stationary condition (3.25), which yields

$$\boldsymbol{\theta}_{j} \stackrel{(3.25)}{=} \frac{1}{M} \left(\boldsymbol{\nu}_{\mathrm{s},jj} + \sum_{i=1}^{M} \boldsymbol{\nu}_{\mathrm{c},ij} \right)$$
$$= \frac{1}{M} \left(\boldsymbol{\theta}_{j} - \tilde{\boldsymbol{\nu}}_{\mathrm{s},jj} + \sum_{i=1}^{M} (\boldsymbol{\theta}_{j} - \tilde{\boldsymbol{\nu}}_{\mathrm{c},ij}) \right)$$
$$= \frac{M+1}{M} \boldsymbol{\theta}_{j} - \frac{1}{M} \left(\tilde{\boldsymbol{\nu}}_{\mathrm{s},jj} + \sum_{i=1}^{M} \tilde{\boldsymbol{\nu}}_{\mathrm{c},ij} \right)$$
(3.45)

$$\Leftrightarrow \quad \left(\frac{M+1}{M}-1\right)\boldsymbol{\theta}_{j} = \frac{1}{M}\boldsymbol{\theta}_{j} = \frac{1}{M}\left(\tilde{\boldsymbol{\nu}}_{\mathrm{s},jj} + \sum_{i=1}^{M}\tilde{\boldsymbol{\nu}}_{\mathrm{c},ij}\right) , \qquad (3.46)$$

i.e.,

$$\boldsymbol{\theta}_{j} = \tilde{\boldsymbol{\nu}}_{\mathrm{s},jj} + \sum_{i=1}^{M} \tilde{\boldsymbol{\nu}}_{\mathrm{c},ij}$$
 (3.47)



(a) Visualization fo the computation of the message to channel factor $f_{y_i|x}(x)$.



(b) Visualization of the computation of the message to signal factor $f_x(x_i)$.

Figure 3.6: Visualizations for the connections between incoming and outgoing messages, i.e., Eqs. (3.49) on the left and (3.48) on the right.

Comparing again with (3.41) and (3.42), this means that the connection between messages leaving the variable node and adjacent incoming messages must be given by

$$\boldsymbol{\nu}_{\mathrm{s},jj} = \sum_{i=1}^{M} \tilde{\boldsymbol{\nu}}_{\mathrm{c},ij} , \qquad (3.48)$$

respectively

$$\boldsymbol{\nu}_{\mathrm{c},ij} = \tilde{\boldsymbol{\nu}}_{\mathrm{s},jj} + \sum_{\substack{i'=1\\i'\neq i}}^{M} \tilde{\boldsymbol{\nu}}_{\mathrm{c},i'j} .$$
(3.49)

The processing is depicted in Fig. 3.6.

The character of the separation in (3.41) and (3.42) is known as *extrinsic* [BG96], i.e., the a-priori knowledge (from before) is distinguished from the knowledge obtained by the processing. From (3.48) and (3.49) it becomes clear that the procedure is characterized by excluding knowledge of the incoming edge on which the next message should be sent. From an algorithmic point of view, especially for graphs with loops, broadcasting extrinsic is crucial to the success of the algorithm, because sending something on an edge that has been received there before, would lead to feedback loops and overemphasizing of certain aspects, which is precisely what was pursued to prevent by the approach based on the maximum entropy method, which was motivated in Sec. 2.2.2. This concept has been proven worthwhile not only in the context of the message passing algorithm, but also, e.g., in Turbo decoding [HOP96] and detection [GH11].

On the one hand, the connections in (3.41) and (3.42) reveal a simple relationship between outgoing and incoming messages as well as the moments computed at the nodes, i.e., the input has to be subtracted from the result θ_j to obtain the returning message on the respective edge. The Equations (3.48) and (3.49), on the other hand, indicate that when computing the message for a certain edge, it is necessary to leave out the incoming message of that very

edge. Hence, there are two ways of processing extrinsic. In the following, the relation is examined in detail for a channel-constrained node $f_{y_i|\mathbf{x}}(\mathbf{x})$, where the input parameter is given by $\mathbf{\nu}_{c,ij} = [m_{\mathbf{x},j}/\sigma_{\mathbf{x},j}^2, 1/\sigma_{\mathbf{x},j}^2]^{\top}$ as above, Eq. (3.34) is used for the mapping between moments and natural parameters and $\tilde{\mathbf{\nu}}_{c,ij} \stackrel{\text{def}}{=} [m_{c,i,j}^{\setminus j}/\sigma_{c,i,j}^{2\setminus j}, 1/\sigma_{c,i,j}^{2\setminus j}]^{\top}$ is defined for $i \in \{1, \ldots, M\}$ and $j \in \{1, \ldots, N\}$. Then Eq. (3.42) can be rewritten as

$$\tilde{\boldsymbol{\nu}}_{\mathrm{c},ij} = \boldsymbol{\theta}_{\mathrm{c},i,j} - \boldsymbol{\nu}_{\mathrm{c},ij} \tag{3.50}$$

so that the elements compute to

$$\frac{m_{c,i,j}^{(j)}}{\sigma_{c,i,j}^{2,\backslash j}} = \frac{m_{c,i,j}}{\sigma_{c,i,j}^{2}} - \frac{m_{x,j}}{\sigma_{x,j}^{2}} \\
\stackrel{(2.16)}{=} \frac{m_{x,j}\sigma_{n}^{2} + m_{x,j}\sum_{\substack{j'=1\\j'\neq j}}^{N}\sigma_{x,j'}^{2}a_{ij'}^{2} + \sigma_{x,j}^{2}a_{ij}\left(y_{i} - \sum_{\substack{j'=1\\j'\neq j}}^{N}a_{ij'}m_{x,j'}\right)}{\sigma_{x,j}^{2}\left(\sigma_{n}^{2} + \sum_{\substack{j'=1\\j'\neq j}}^{N}\sigma_{x,j'}^{2}a_{ij'}^{2}\right)} - \frac{m_{x,j}}{\sigma_{x,j}^{2}} \\
= \frac{a_{ij}\left(y_{i} - \sum_{\substack{j'=1\\j'\neq j}}^{N}a_{ij'}m_{x,j'}\right)}{\sigma_{n}^{2} + \sum_{\substack{j'=1\\j'\neq j}}^{N}\sigma_{x,j'}^{2}a_{ij'}^{2}} \tag{3.51}$$

and

$$\frac{1}{\sigma_{\mathrm{c},i,j}^{2,\backslash j}} = \frac{1}{\sigma_{\mathrm{c},i,j}^{2}} - \frac{1}{\sigma_{\mathrm{x},j}^{2}} \stackrel{(2.17)}{=} \frac{\sigma_{\mathrm{n}}^{2} + \sum_{j=1}^{N} \sigma_{\mathrm{x},j}^{2} a_{ij}^{2}}{\sigma_{\mathrm{x},j}^{2} \left(\sigma_{\mathrm{n}}^{2} + \sum_{j'\neq j}^{N} \sigma_{\mathrm{x},j'}^{2} a_{ij'}^{2}\right)} - \frac{1}{\sigma_{\mathrm{x},j}^{2}} \\
= \frac{\sigma_{\mathrm{n}}^{2} + \sum_{j=1}^{N} \sigma_{\mathrm{x},j}^{2} a_{ij}^{2} - \sigma_{\mathrm{n}}^{2} - \sum_{j'\neq j}^{N} \sigma_{\mathrm{x},j'}^{2} a_{ij'}^{2}}{\sigma_{\mathrm{x},j}^{2} \left(\sigma_{\mathrm{n}}^{2} + \sum_{j'\neq j}^{N} \sigma_{\mathrm{x},j'}^{2} a_{ij'}^{2}\right)} \\
= \frac{a_{ij}^{2}}{\sigma_{\mathrm{n}}^{2} + \sum_{j'\neq j}^{N} \sigma_{\mathrm{x},j'}^{2} a_{ij'}^{2}}.$$
(3.52)

This reveals that in the calculation for $\tilde{\nu}_{c,ij}$, one can directly leave out $\nu_{c,ij} = [m_{x,j}/\sigma_{x,j}^2, 1/\sigma_{x,j}^2]^\top$ as input parameter, i.e., the calculation of the extrinsic can be done as *post-processing* via the subtraction imposed by (3.42) or as *pre-processing* by omitting the respective input parameter. The respective computations are contrasted in the upper row of Fig. 3.8 as block diagrams, where the notation of the GMP algorithm (stated below) is used. On the left (Fig. 3.8(a)), the pre-processed extrinsic is shown, whereas on the right (Fig. 3.8(b)) the extrinsic is postprocessed. Both philosophies, pre- and post-processed extrinsic, are used in the literature. For example, estimation theoretic extrinsic [GH11] in the context of Gaussian noise uses the subtractions as specified above. The pre-processing is used for example in the decoding procedure for low-density parity-check (LDPC) codes [Gal62]. In the message passing context, usually the pre-processing, i.e., incident parameter omitting strategy is considered [KMS⁺12, Pea88]. In case of a signal factor $f_x(x_j)$, the relation (3.48) indicates to compute the estimates (2.19) and (2.20) based on all messages coming from the channel-constrained nodes. Via the projection (3.35) and connection (3.42), an outgoing message $\nu_{c,ij}$ is computed. Since the estimators (2.19) and (2.20) are usually non-linear, it is not possible to move the effect of the subtraction to the input of the estimator. Thus, the equivalence of post- and pre-processed extrinsic only holds if the factor marginal is Gaussian. This is usually not the case for the signal factors in CS. Hence, both procedures yield different results. The different processings are depicted in the lower row of Fig. 3.8; pre-processed extrinsic can be found on the left side, post-processed extrinsic on the right side.

3.3.4 Message Passing

Following the notion of pre-processing the extrinsic, i.e., omitting the incoming message of the edge for which the message is computed yields the *Gaussian belief propagation* (GBP) or *Gaussian message passing* (GMP) as stated in [KMS⁺12, Sec. 3.1]. In order to get there, recall that an outgoing message $\nu_{c,ij}$ carries the information of the signal-constrained factor stimulated by all incoming messages $\tilde{\nu}_{c,i'j}$ with $i' \in \{1, \ldots, M\} \setminus \{i\}$, i.e., all except the message of the considered edge. This means that the input parameters to the corresponding estimates (2.19) and (2.20) are computed via

$$\frac{\tilde{x}_{s,j}^{\backslash i}}{\tilde{\sigma}_{s,j}^{2,\backslash i}} = \sum_{\substack{i'=1\\i'\neq i}}^{M} \frac{m_{c,i',j}^{\backslash j}}{\sigma_{c,i',j}^{2,\backslash j}}, \quad \text{and} \quad \frac{1}{\tilde{\sigma}_{s,j}^{2,\backslash i}} = \sum_{\substack{i'=1\\i'\neq i}}^{M} \frac{1}{\sigma_{c,i',j}^{2,\backslash j}}.$$
(3.53)

With these inputs, one can define

$$m_{\mathrm{s},j}^{\backslash i} \stackrel{\text{def}}{=} \mathrm{E}_{\mathsf{x}_j} \left\{ \mathsf{x}_j \mid \tilde{x}_{\mathrm{s},j}^{\backslash i} \right\} = \frac{\int x_j \mathsf{f}_\mathsf{x}(x_j) \exp\left(-\frac{(\tilde{x}_{\mathrm{s},j}^{\backslash i} - x_j)^2}{2\tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i}}\right) \, \mathrm{d}x_j}{\int \mathsf{f}_\mathsf{x}(x_j) \exp\left(-\frac{(\tilde{x}_{\mathrm{s},j}^{\backslash i} - x_j)^2}{2\tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i}}\right) \, \mathrm{d}x_j}$$
(3.54)

and

$$\sigma_{\mathrm{s},j}^{2,\backslash i} \stackrel{\text{def}}{=} \mathrm{E}_{\mathsf{x}_{j}} \Big\{ (\mathsf{x}_{j} - m_{\mathrm{s},j}^{\backslash i})^{2} \mid \tilde{x}_{\mathrm{s},j}^{\backslash i} \Big\} = \frac{\int (x_{j} - m_{\mathrm{s},j})^{2} \mathsf{f}_{\mathsf{x}}(x_{j}) \exp\left(-\frac{(\tilde{x}_{\mathrm{s},j}^{\backslash i} - x_{j})^{2}}{2\tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i}}\right) \, \mathrm{d}x_{j}}{\int \mathsf{f}_{\mathsf{x}}(x_{j}) \exp\left(-\frac{(\tilde{x}_{\mathrm{s},j}^{\backslash i} - x_{j})^{2}}{2\tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i}}\right) \, \mathrm{d}x_{j}} \,. \tag{3.55}$$

The resulting estimates, mapped to the natural parameters, serve as input to the computation of the channel-constrained estimate. Therefore, defining $\boldsymbol{\nu}_{\mathrm{c},ij} = [m_{\mathrm{s},j}^{\backslash i}/\sigma_{\mathrm{s},j}^{2,\backslash i}, 1/\sigma_{\mathrm{s},j}^{2,\backslash i}]^{\top}$, the computations at the channel-constrained nodes become

$$\tilde{\boldsymbol{\nu}}_{c,ij} = \left[\frac{m_{c,i,j}^{\backslash j}}{\sigma_{c,i,j}^{2,\backslash j}}, \frac{1}{\sigma_{c,i',j}^{2,\backslash j}}\right]^{\top} = \left[\frac{a_{ij}\left(y_{i} - \sum_{\substack{j'=1\\j'\neq j}}^{N} a_{ij'}m_{s,j'}^{\backslash i}\right)}{\sigma_{n}^{2} + \sum_{\substack{j'=1\\j'\neq j}}^{N} \sigma_{s,j'}^{2,\backslash i}a_{ij'}^{2}}, \frac{a_{ij}^{2}}{\sigma_{n}^{2} + \sum_{\substack{j'=1\\j'\neq j}}^{N} \sigma_{s,j'}^{2,\backslash i}a_{ij'}^{2}}\right]^{\top}.$$
 (3.56)

The resulting algorithm iterates between the nodes by choosing the indices $j \in \{1, ..., N\}$ and $i \in \{1, ..., M\}$ in a suitable manner. The entire procedure is stated in Algorithm 3.1 below and in Algorithm E.1 in Appendix E. The notation uses only means and variances for comprehensibility, however, for the implementation it is beneficial to resort to scaled means and inverse variances (precisions), so that, e.g., the summations in Lines 7 and 8 spare the inverses, and scalings.

Algorithm 3.1: $m{m}_{ m s}={ t GMP}(m{y},m{A},\sigma_{ m n}^2)$	
1 N	$m_{\mathrm{s},j}^{\setminus i} = 0$, $\sigma_{\mathrm{s},j}^{2,\setminus i} = s/N \; \forall \; j \in \{1, \ldots, N\}, \; i \in \{1, \ldots, M\}$ // initialization
2 $m_{c,i,j}^{\langle j \rangle} = 0$, $\sigma_{c,i,j}^{2,\langle j \rangle} = s/N \ \forall \ j \in \{1, \ldots, N\}, \ i \in \{1, \ldots, M\}$	
³ while stopping criterion not met do	
4	Choose $j \in \{1,, N\}$ and $i \in \{1,, M\}$
5	$m_{\mathrm{c},i,j}^{\setminus j} = (y_i - \sum_{\substack{j'=1\\j'\neq i}}^N a_{ij'} m_{\mathrm{s},j'}^{\setminus i})/a_{ij}$ // "upward" messages
6	$\sigma_{\mathbf{c},i,j}^{2,\backslash j} = (\sigma_{\mathbf{n}}^{2} + \sum_{\substack{j'=1\\j'\neq j}}^{N} \sigma_{\mathbf{s},j'}^{2,\backslash i} a_{ij'}^{2}) / a_{ij}^{2}$
7 8	$ \begin{array}{l} \tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i} = \left(\sum_{\substack{i'=1\\i'\neq i}}^{M} 1/\sigma_{\mathrm{c},i',j}^{2,\backslash j}\right)^{-1} & // \text{ input to NLMMSE (edge-dependent)} \\ \tilde{x}_{\mathrm{s},j}^{\backslash i} = \tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i} \sum_{\substack{i'=1\\i'\neq i}}^{M} m_{\mathrm{c},i',j}^{\backslash j}/\sigma_{\mathrm{c},i',j}^{2,\backslash j} \end{array} $
9 10	$ \begin{bmatrix} m_{\mathrm{s},j}^{\backslash i} = \mathrm{E}_{x_j} \{x_j \mid \tilde{x}_{\mathrm{s},j}^{\backslash i}, \tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i} \} \\ \sigma_{\mathrm{s},j}^{2,\backslash i} = \mathrm{E}_{x_j} \{(x_j - m_{\mathrm{s},j}^{\backslash i})^2 \mid \tilde{x}_{\mathrm{s},j}^{\backslash i}, \tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i} \} \end{bmatrix} // \text{ ''downward'' messages} $
11 õ	$\hat{\tau}_{\mathrm{s},j}^2 = \left(\sum_{i=1}^M 1/\sigma_{\mathrm{c},i',j}^{2,\backslash j}\right)^{-1} \qquad \forall j \in \{1, \dots, N\}$
12 \tilde{x}	$\tilde{c}_{\mathbf{s},j} = \tilde{\sigma}_{\mathbf{s},j}^2 \sum_{i=1}^M m_{\mathbf{c},i',j}^{\setminus j} / \sigma_{\mathbf{c},i',j}^{2\setminus j} \qquad \forall j \in \{1, \dots, N\}$
13 N	$m{n}_{ m s} = { m E}_{f x} \{ f x \mid ilde x_{ m s}, ilde \sigma_{ m s,1}^2, \dots, ilde \sigma_{ m s,N}^2 \}$ // NLMMSE

The statement of GMP in Algorithm 3.1 implies that the messages are sent from one factor node $f_{y_i|x}(x)$ to one variable node x_j ("upwards") and then the messages from variable node x_j (by employing factor node $f_x(x_j)$) to factor node $f_{y_i|x}(x)$ ("downwards"); subsequently choosing the next indices to be processed. This notation is used to simplify readability; for the simulations in Chapter 6, a schedule that iterates between computing all upward messages, followed by computing all downward messages, is used. Of course also other schedules are possible.

The "upward" and "downward" messages as used in the statement of the algorithm are also visualized in Fig. 3.7. The input message for the signal factor $f_x(x_j)$ is not shown, because it depends for which outgoing edge (downwards) the estimate is computed; here it may be either $\tilde{x}_{s,j}^{\setminus 1}$, together with $\tilde{\sigma}_{s,j}^{2,\setminus 1}$, or $\tilde{x}_{s,j}^{\setminus M}$ and $\tilde{\sigma}_{s,j}^{2,\setminus M}$.



Figure 3.7: Visualization of "upward" and "downward" messages in terms of means and variances in the factor graph for CS.

3.3.5 Related Algorithms

Despite simplifying to two scalar moments instead of computing arbitrary marginals or multidimensional estimates, the resulting algorithm still inherits high computational complexity. The reason lies in the fine-granular view on the factor graph that causes edge-dependent estimates (and messages) that have to be computed, stored and propagated and the fact that the variable nodes are connected to all channel-constrained factors, i.e., the sensing matrix A is not sparse. The edge-dependent processing therefore impedes scalability, i.e., the algorithm can only be used for tasks with small dimensions.

One way to reduce the computation overhead is to utilize post-processed extrinsic, instead of pre-processed extrinsic as in Algorithm 3.1. In Eqs. (3.51) and (3.52), it has been shown that for the channel-constrained estimation both processings are mathematically the same. The post-processed extrinsic has in combination with the schedule that computes all "upward" messages together the advantage that the expectations get node-dependent, i.e., have to be performed M times, instead of $M \cdot N$ times, when they are edge-dependent as it is the case when the extrinsic is pre-processed. Of course, afterwards the computation of the extrinsic is again edge-dependent, but this involves only scalings and one subtraction (per edge and parameter), as opposed to the multiple summations of (almost) all parameters in the computation of the expectations. This form of the algorithm is not stated in detail, but referred to as GMPpostExtLin for GMP with post-processed extrinsic at the linear estimation.

If it is assumed that the non-linear estimation in the signal-constrained factor behaves similarly for pre- and post-processed extrinsic, then the same strategy can be applied again to the computation of the respective estimates, i.e., one can switch from pre- to post-processed extrinsic. The resulting algorithm is shown in Appendix E.2 as Algorithm E.2—it is called GMPpostExt—and benefits from the fact that the input parameters for the signal-constrained estimation do not have to be computed edge-dependently. Therefore, the signal-constrained estimation can be computed once for all variables x_1, \ldots, x_N and the "downward" message is obtained by subtracting only the message from the incoming edge.

GMP and its variants can be composed from the block diagrams shown in Fig. 3.8. The GMP algorithm itself utilizes pre-processed extrinsic, i.e., only blocks from the left side (and perhaps



(a) Pre-processing of extrinsic at channel-constrained node; messages $m_{{\rm s},j}^{\backslash i}$ and $\sigma_{{\rm s},j}^{2,\backslash i}$ are left out.



(c) Pre-processed extrinsic at signal-constrained node, i.e., the input values of the desired edge are left out.



(b) Post-processing of extrinsic at channel-constrained node; all entries are summed up, $m_{{\rm s},j}^{\backslash i}$ and $\sigma_{{\rm s},j}^{2,\backslash i}$ are scaled and subtracted afterwards.



(d) Extrinsic at signal-constrained node is postprocessed, i.e., the respective input is subtracted in the end.

Figure 3.8: Computations at the channel-constrained (upper row) and the signal-constrained nodes (lower row). The left side shows pre-processed, the right side shows post-processed extrinsic.

a suitable transformation from output to input parameters between the blocks). In contrast, GMPpostExt only uses the right-side blocks; and GMPpostExtLin mixes the approaches by taking the upper right and lower left blocks for its processing. Note that, for simplicity, the block representing the signal-constrained estimation accepts here scaled mean and inverse variance (similar to the implementation in (6.1)-(6.3)) as inputs and returns conditional mean and variance.

In the literature, one can also find various approaches to mitigate the scalability problem by simplifying the strategy. In [KMS⁺12] the complexity is reduced by restricting to signalconstrained and channel-constrained estimates for each variable x_j while also keeping two estimates (mean and variance) per observation y_i , the resulting algorithm is called *generalized approximate message passing* (GAMP) [Ran11]. Further simplifications, e.g., the step from individual variances to average variances and the insertion of the observation dependent parts into the calculation of the residuum, lead to the regular form of *approximate message passing* (AMP) [DMM09, Mal11]. Noteworthy, the LMMSE solution in the linear estimation is simplified to a matched filter (MF), i.e., the transpose of the sensing matrix is used without the inverse that is computed in the LMMSE estimator. The AMP and the GAMP algorithm are stated in the notation used here in Appendices E.3 and E.4, respectively. AMP works well for Gaussian sensing matrices and possesses very low computational complexity, but deteriorates for more general classes of sensing matrices [CZK14].

With the focus on factor graphs and its various representations of the CS problem, we will not follow this path but instead consider the opposite approach to the fine-granular view that is considering only two factors.

4. Turbo-type Inference

The notion of Turbo-type processing was established in [BGT93] in the form of Turbo decoding. The idea is to consider codes composed of two component codes. Decoding can then be done for the component codes, which is less computationally costly, with the results being passed in a suitable way between the decoders, i.e., an iterative scheme is present; for an overview, see [HOP96]. The overall goal is, thus, the exchange of *messages* between two major parts constituting to the problem. The suitable way to process the message for sending leads again to the notion of *extrinsic* [BG96, GH11].

The Turbo principle has also been used in detection for digital transmission schemes and is known as *Turbo equalization* [DJB⁺95] there, for an overview see [TS11].

In this chapter, the processing of Turbo-type algorithms for CS is derived and connections to the famous VAMP algorithm [RSF19] are revealed. To that end, the same procedure of considering a constrained optimization problem for the approximation of the CS posterior—only this time for a different form of the factorization—as in the chapter before is employed. This shall show that the underlying principles are the same. For inference in general, the optimization approach has in [HOW⁺05] been shown to be equivalent to expectation-consistent (EC) approximate inference [OW05], which coincides with the VAMP algorithm in the given case.

For the case of Turbo codes, the connection to message passing (MP) has already been shown in [MMC98], [KF98], and [Wib96]. Noteworthy, VAMP is also derived from a messagepassing-type view in [RSF19], however, the exposition here puts everything into one notation, states the derivation in detail and draws attention to the optimization point of view.

Other literature that considers optimization-based approaches, usually adds the optimization on top, so that a *double-loop algorithm* results [OW05, RFSK16]. Here, the focus lies on the optimization within the factor graph, yielding *single-loop algorithms* only.

The connection of Turbo-type processing and minimization of the Kullback–Leibler divergence is investigated in [MG98] for iterative decoding and in [Moh97a, Moh97b] for multi-user detection in digital transmission scenarios.



Figure 4.1: Factor graph for CS separating the two different kinds of factors.

4.1 Structure of the Factor Graph

The Turbo approach implies two opposing parts of a problem. In the compressed sensing case, this is given by considering the factorization

$$f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}) = \text{const.} \cdot f_{\mathbf{y}|\mathbf{x}}(\mathbf{x}) \cdot f_{\mathbf{x}}(\mathbf{x})$$
(4.1)

of the posterior, which combines factors of the same kind in one factor, thereby explicitly separating channel and signal part while keeping the different components of each part together. The corresponding factor graph is shown in Fig. 2.3(d) and for convenience plotted here in Fig. 4.1 again.

The approach can be motivated by the fact that the parts, which complicate the overall estimation, are separated. In order to comprehend this, recall that the computation of the conditional mean (2.9) is infeasible, since the multi-dimensional integral depends on the signal prior $f_x(x)$ and can in general not be solved analytically. Further recall, that the estimators in Sec. 2.2.1.1 and Sec. 2.2.1.3 have been derived by replacing either of the factors with a Gaussian distribution, respectively, which resulted in feasible MMSE estimators. On the one hand, the multi-dimensionality of the estimator (2.10) can be handled, since the estimator is linear, which is optimal under the assumption of a Gaussian prior $f_x(x)$. If one or more signal components x_j were supposed to be distributed differently, the linear estimator would not be optimal and the corresponding multi-dimensional non-linear estimator could not be expressed, i.e., the computational effort would increase disproportionately. On the other hand, the non-linear estimators (2.19) can be handled because they are separated into one-dimensional integrals that can be solved, i.e., the multi-dimensionality is broken down, since the dependencies induced by the observations y and the sensing matrix A are ignored.

Hence, the approach in this chapter aims to separate what needs to be separated without going further into detail, i.e., what can be considered together is considered as one factor. Thereby, the issue of scalability, which arose for MP from the huge number of edges in the factor graph, is also addressed since the number of edges is decreased drastically. However, one has to keep in mind that instead of two scalars, the messages on the edges are now N-dimensional vectors.

4.1.1 Factorization of the Substitute Distribution

This more distanced view with less focus on the details has direct consequences on the form of the approximating distribution $q(\boldsymbol{x})$. For the given factorization, the partial approximations $q_s(\boldsymbol{x})$ are used for $f_{\boldsymbol{x}}(\boldsymbol{x})$, as well as $q_c(\boldsymbol{x})$ for $f_{\boldsymbol{y}|\boldsymbol{x}}(\boldsymbol{x})$. Comparing to Chapter 3, the number of approximations is drastically reduced. Here, the adaptions by the variations may be written as follows

$$f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}) = \text{const.} \cdot f_{\mathbf{x}}(\mathbf{x}) \cdot f_{\mathbf{y}|\mathbf{x}}(\mathbf{x}) \cdot \frac{s(\mathbf{x}) \cdot c(\mathbf{x})}{s(\mathbf{x}) \cdot c(\mathbf{x})}$$
$$= \text{const.} \cdot \frac{\mathbf{q}_{s}(\mathbf{x}) \cdot \mathbf{q}_{c}(\mathbf{x})}{s(\mathbf{x}) \cdot c(\mathbf{x})} , \qquad (4.2)$$

i.e., the representatives are defined as

$$q_{s}(\boldsymbol{x}) = f_{\boldsymbol{x}}(\boldsymbol{x}) \cdot s(\boldsymbol{x}) , \qquad (4.3)$$

$$q_{c}(\boldsymbol{x}) = f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) \cdot c(\boldsymbol{x}) .$$
(4.4)

Since, only two deviations are present, the manipulated variable node obeys

$$\frac{\text{const.}}{s(\boldsymbol{x}) \cdot c(\boldsymbol{x})} = \text{const.} \int \frac{\delta(\boldsymbol{x} - \tilde{\boldsymbol{x}})}{s(\boldsymbol{x}) \cdot c(\tilde{\boldsymbol{x}})} \,\mathrm{d}\tilde{\boldsymbol{x}} , \qquad (4.5)$$

so that the representative distribution for the variable node is defined as

$$\mathbf{q}_{\mathbf{v}}(\boldsymbol{x}) = \text{const.} \cdot s(\boldsymbol{x}) \cdot c(\boldsymbol{x}) . \tag{4.6}$$

All in all, the substitute distribution is given by

$$\mathbf{q}(\boldsymbol{x}) = \frac{\mathbf{q}_{s}(\boldsymbol{x}) \cdot \mathbf{q}_{c}(\boldsymbol{x})}{\mathbf{q}_{v}(\boldsymbol{x})} .$$
(4.7)

4.1.2 The Moment-Matching Constraint

The connection between the approximated factors is routed via the (combined) variable node. Since the variable node has an edge to both factors, the consistency is now required on these edges. Our interest in moments, especially in the conditional mean (2.9), justifies the use of a moment constraint in this case as well. Therefore, exponential families are used again, i.e., a certain statistic is specified on which to rely the estimation and which should be matched in order to achieve consistency. That is, for an arbitrary sufficient statistic g(x) of the exponential family,

$$\mathbf{E}_{\mathbf{x} \sim \mathbf{q}_{v}} \{ \boldsymbol{g}(\mathbf{x}) \} \stackrel{!}{=} \mathbf{E}_{\mathbf{x} \sim \mathbf{q}_{s}} \{ \boldsymbol{g}(\mathbf{x}) \} , \qquad (4.8)$$

$$\mathbf{E}_{\mathbf{x} \sim q_{v}} \{ \boldsymbol{g}(\mathbf{x}) \} \stackrel{!}{=} \mathbf{E}_{\mathbf{x} \sim q_{c}} \{ \boldsymbol{g}(\mathbf{x}) \}$$

$$(4.9)$$

is required. Note that depending on the sufficient statistic that specifies the form of consistency, the requirements and thereby the resulting messages are based on multi-dimensional vectors. For example, in Chapter 3, the *individual variances* case was employed, which leads to 2N-dimensional messages per edge. Another possibility introduced in Appendix D is the *average variance* case yielding N + 1-dimensional messages. For now, this may stay open; later it will be specified for the respective algorithms.

4.2 Constrained Optimization

As in Chapter 3, the minimization of the Kullback–Leibler divergence is considered for the given structure to derive stationary points that are aspired to be achieved by recovery algorithms. For the structure at hand, there are vector-wise consistency constraints for the matching of moments and two approximating factors given. The constrained optimization problem reads

minimize
$$D_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) || \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}))$$

subject to $E_{\mathbf{x} \sim q_{\mathrm{v}}}\{\boldsymbol{g}(\mathbf{x})\} - E_{\mathbf{x} \sim q_{\mathrm{c}}}\{\boldsymbol{g}(\mathbf{x})\} = \mathbf{0}$,
 $E_{\mathbf{x} \sim q_{\mathrm{v}}}\{\boldsymbol{g}(\mathbf{x})\} - E_{\mathbf{x} \sim q_{\mathrm{s}}}\{\boldsymbol{g}(\mathbf{x})\} = \mathbf{0}$,
 $\int \mathbf{q}_{\mathrm{v}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 = 0$, $\int \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 = 0$, $\int \mathbf{q}_{\mathrm{c}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 = 0$.
(4.10)

4.2.1 Cost Function for Compressed Sensing

For the given case, the Kullback-Leibler divergence to be minimized splits directly into

$$D_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) \mid\mid f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})) = \int \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \log \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} + \int \mathbf{q}_{\mathrm{c}}(\boldsymbol{x}) \log \mathbf{q}_{\mathrm{c}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - \int \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \log f_{\mathbf{x}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - \int \mathbf{q}_{\mathrm{c}}(\boldsymbol{x}) \log f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - \int \mathbf{q}_{\mathrm{v}}(\boldsymbol{x}) \log \mathbf{q}_{\mathrm{v}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} .$$
(4.11)

For this form of the Kullback–Leibler divergence, the stationary points of the Lagrangian incorporating additionally the moment-matching constraints are derived. The stationary points are of interest, because they minimize the cost function under the given constraints, and will later be used to derive recovery algorithms.

4.2.2 Stationary Points of the Lagrangian

As before, the Lagrangian, cf. Appendix B, is used to insert the moment-matching constraint into the optimization problem. By introducing the Lagrangian multipliers ν_c and ν_s for the moment-matching constraints and ν_v , ν_c , and ν_s for the normalization constraints, the Lagrangian reads (again without multipliers in the argument)

$$\mathcal{L}(\mathbf{q}(\boldsymbol{x})) = D_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) || \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})) + \boldsymbol{\nu}_{\mathrm{s}}^{\top} (\mathrm{E}_{\mathbf{x},\mathsf{q}_{\mathrm{v}}} \{\boldsymbol{g}(\mathbf{x})\} - \mathrm{E}_{\mathbf{x},\mathsf{q}_{\mathrm{s}}} \{\boldsymbol{g}(\mathbf{x})\}) + \boldsymbol{\nu}_{\mathrm{c}}^{\top} (\mathrm{E}_{\mathbf{x},\mathsf{q}_{\mathrm{v}}} \{\boldsymbol{g}(\mathbf{x})\} - \mathrm{E}_{\mathbf{x},\mathsf{q}_{\mathrm{c}}} \{\boldsymbol{g}(\mathbf{x})\}) + \nu_{\mathrm{v}} \left(\int \mathbf{q}_{\mathrm{v}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 \right) + \nu_{\mathrm{s}} \left(\int \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 \right) + \nu_{\mathrm{c}} \left(\int \mathbf{q}_{\mathrm{c}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 \right) .$$

$$(4.12)$$

The optimal form for the distributions is given by the stationary points of the Lagrangian.

Hence, the Lagrangian is differentiated w.r.t. the distributions, which yields

$$\frac{\partial}{\partial q_{v}} \mathcal{L}(q(\boldsymbol{x})) = -\log q_{v}(\boldsymbol{x}) + (\boldsymbol{\nu}_{s} + \boldsymbol{\nu}_{c})^{\top} \boldsymbol{g}(\boldsymbol{x}) + \text{const.}, \qquad (4.13)$$

$$\frac{\partial}{\partial q_{s}} \mathcal{L}(q(\boldsymbol{x})) = \log q_{s}(\boldsymbol{x}) - \log f_{\boldsymbol{x}}(\boldsymbol{x}) - \boldsymbol{\nu}_{s}^{\top} \boldsymbol{g}(\boldsymbol{x}) + \text{const.}, \qquad (4.14)$$

$$\frac{\partial}{\partial q_{c}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = \log q_{c}(\boldsymbol{x}) - \log f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) - \boldsymbol{\nu}_{c}^{\top} \boldsymbol{g}(\boldsymbol{x}) + \text{const.}$$
(4.15)

By setting the above terms to zero, the respective stationary points, thus, calculate to

$$q_{v}(\boldsymbol{x}) = \frac{1}{Z_{v}(\boldsymbol{\nu}_{s}, \boldsymbol{\nu}_{c})} \exp\left((\boldsymbol{\nu}_{s} + \boldsymbol{\nu}_{c})^{\top} \boldsymbol{g}(\boldsymbol{x})\right) , \qquad (4.16)$$

$$\mathbf{q}_{s}(\boldsymbol{x}) = \frac{1}{Z_{s}(\boldsymbol{\nu}_{s})} \mathbf{f}_{\boldsymbol{x}}(\boldsymbol{x}) \exp\left(\boldsymbol{\nu}_{s}^{\top} \boldsymbol{g}(\boldsymbol{x})\right) , \qquad (4.17)$$

$$q_{c}(\boldsymbol{x}) = \frac{1}{Z_{c}(\boldsymbol{\nu}_{c})} f_{\boldsymbol{y}|\boldsymbol{x}}(\boldsymbol{x}) \exp\left(\boldsymbol{\nu}_{c}^{\top} \boldsymbol{g}(\boldsymbol{x})\right) , \qquad (4.18)$$

with normalizing partition functions¹

$$Z_{\rm v}(\boldsymbol{\nu}_{\rm s},\,\boldsymbol{\nu}_{\rm c}) = \int \exp\left((\boldsymbol{\nu}_{\rm s}+\boldsymbol{\nu}_{\rm c})^{\top}\boldsymbol{g}(\boldsymbol{x})\right)\,\mathrm{d}\boldsymbol{x}\;,\tag{4.19}$$

$$Z_{\rm s}(\boldsymbol{\nu}_{\rm s}) = \int f_{\mathbf{x}}(\boldsymbol{x}) \exp\left(\boldsymbol{\nu}_{\rm s}^{\top} \boldsymbol{g}(\boldsymbol{x})\right) \, \mathrm{d}\boldsymbol{x} \;, \tag{4.20}$$

$$Z_{\rm c}(\boldsymbol{\nu}_{\rm c}) = \int f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) \exp\left(\boldsymbol{\nu}_{\rm c}^{\top}\boldsymbol{g}(\boldsymbol{x})\right) \,\mathrm{d}\boldsymbol{x} \;. \tag{4.21}$$

These forms show directly the connection between variable and factor nodes, as the natural parameters of the variable node are given by

$$\boldsymbol{\theta} \stackrel{\text{\tiny def}}{=} \boldsymbol{\nu}_{\mathrm{c}} + \boldsymbol{\nu}_{\mathrm{s}} \; .$$
 (4.22)

Considering that the two summands are the input parameters for the estimations $E_{\mathbf{x}\sim q_c}{\mathbf{g}(\mathbf{x})}$, and $E_{\mathbf{x}\sim q_s}{\mathbf{g}(\mathbf{x})}$, respectively, it is already visible how the extrinsic calculations are performed here. Which is to say, the extrinsic is obtained via *post-processing* of the parameter θ , i.e., the input for the other estimation is achieved by subtracting the input of the currently considered estimation from the result of the estimation. This is usual for Turbo-type approaches, e.g., in coding related publications [HOP96] a-priori log-likelihood ratios (LLRs) are subtracted from the decoding results of component codes to obtain the extrinsic values.

4.2.3 The Lagrange Dual Function

Before continuing with the derivation of algorithms for this approach, we take a step back and consider the *Lagrange dual function* of the problem. It is defined as infimum of the Lagrangian w.r.t. the optimization parameter, which is q(x) in the given case, i.e.,

$$\mathcal{L}_{\mathrm{D}}(\boldsymbol{\nu}_{\mathrm{s}},\,\boldsymbol{\nu}_{\mathrm{c}}) = \inf_{\boldsymbol{\mathfrak{q}}(\boldsymbol{x})} \mathcal{L}(\boldsymbol{\mathfrak{q}}(\boldsymbol{x})) \;. \tag{4.23}$$

¹The term *partition function* stems from exponential families, for details see Appendix D.

Therefore, it is a function of the Lagrangian multipliers only. It offers a different way in approaching the optimal value, namely from below. The so-called *dual problem* is the task of maximizing the Lagrange dual function. Since *strong duality* holds, the maximum value of $\mathcal{L}_{\rm D}(\boldsymbol{\nu}_{\rm s}, \boldsymbol{\nu}_{\rm c})$ is indeed the desired optimum. For more information, see Appendix B. Noteworthy, the switch from primal to dual problem, also shifts the focus to the Lagrangian multipliers as parameters to be optimized.

Since the stationary points which minimize the Lagrangian are already given, the dual function is obtained by inserting the solution for the distributions $q_v(\boldsymbol{x})$, $q_s(\boldsymbol{x})$, and $q_c(\boldsymbol{x})$ into the Lagrangian. Note that at the stationary points the constraints must be fulfilled, i.e., they vanish. The insertion leads to

$$\begin{aligned} \mathcal{L}_{\mathrm{D}}(\boldsymbol{\nu}_{\mathrm{s}},\,\boldsymbol{\nu}_{\mathrm{c}}) &= \int \mathsf{q}_{\mathrm{s}}(\boldsymbol{x}) \left(\log f_{\mathbf{x}}(\boldsymbol{x}) + \boldsymbol{\nu}_{\mathrm{s}}^{\top} \boldsymbol{g}(\boldsymbol{x}) - \log Z_{\mathrm{s}}(\boldsymbol{\nu}_{\mathrm{s}}) \right) \, \mathrm{d}\boldsymbol{x} \\ &+ \int \mathsf{q}_{\mathrm{c}}(\boldsymbol{x}) \left(\log f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) + \boldsymbol{\nu}_{\mathrm{c}}^{\top} \boldsymbol{g}(\boldsymbol{x}) - \log Z_{\mathrm{c}}(\boldsymbol{\nu}_{\mathrm{c}}) \right) \, \mathrm{d}\boldsymbol{x} \\ &- \int \mathsf{q}_{\mathrm{s}}(\boldsymbol{x}) \log f_{\mathbf{x}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - \int \mathsf{q}_{\mathrm{c}}(\boldsymbol{x}) \log f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} \\ &- \int \mathsf{q}_{\mathrm{v}}(\boldsymbol{x}) \left((\boldsymbol{\nu}_{\mathrm{s}} + \boldsymbol{\nu}_{\mathrm{c}})^{\top} \boldsymbol{g}(\boldsymbol{x}) - \log Z_{\mathrm{v}}(\boldsymbol{\nu}_{\mathrm{s}},\,\boldsymbol{\nu}_{\mathrm{c}}) \right) \, \mathrm{d}\boldsymbol{x} \\ &= -\log Z_{\mathrm{s}}(\boldsymbol{\nu}_{\mathrm{s}}) \int \mathsf{q}_{\mathrm{s}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - \log Z_{\mathrm{c}}(\boldsymbol{\nu}_{\mathrm{c}}) \int \mathsf{q}_{\mathrm{c}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} \\ &+ \log Z_{\mathrm{v}}(\boldsymbol{\nu}_{\mathrm{s}},\,\boldsymbol{\nu}_{\mathrm{c}}) \int \mathsf{q}_{\mathrm{v}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} \\ &+ \log Z_{\mathrm{v}}(\boldsymbol{\nu}_{\mathrm{s}},\,\boldsymbol{\nu}_{\mathrm{c}}) \int \mathsf{q}_{\mathrm{v}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} \\ &+ \nu_{\mathrm{s}}^{\top} \left(\mathrm{E}_{\mathbf{x},\mathsf{q}_{\mathrm{s}}} \{ \boldsymbol{g}(\mathbf{x}) \} - \mathrm{E}_{\mathbf{x},\mathsf{q}_{\mathrm{v}}} \{ \boldsymbol{g}(\mathbf{x}) \} \right) \\ &+ \boldsymbol{\nu}_{\mathrm{c}}^{\top} \left(\mathrm{E}_{\mathbf{x},\mathsf{q}_{\mathrm{c}}} \{ \boldsymbol{g}(\mathbf{x}) \} - \mathrm{E}_{\mathbf{x},\mathsf{q}_{\mathrm{v}}} \{ \boldsymbol{g}(\mathbf{x}) \} \right) . \end{aligned}$$

Since the equality constraints are fulfilled [HOW⁺05]

$$\mathcal{L}_{\rm D}(\boldsymbol{\nu}_{\rm s},\,\boldsymbol{\nu}_{\rm c}) = -\log Z_{\rm s}(\boldsymbol{\nu}_{\rm s}) - \log Z_{\rm c}(\boldsymbol{\nu}_{\rm c}) + \log Z_{\rm v}(\boldsymbol{\nu}_{\rm s},\,\boldsymbol{\nu}_{\rm c})$$
(4.25)

results. This means that alternatively to minimizing the Kullback–Leibler divergence between densities, an approach to solve the CS problem is to adjust the parameters ν_s and ν_c such that the cost function (4.25) is maximized. Noteworthy, this cost function already incorporates the moment-matching and normalization constraints in itself due to the definition of the partition functions $Z_v(\nu_s, \nu_c), Z_s(\nu_s)$, and $Z_c(\nu_c)$.

Equation (4.25) is the cost function that the expectation-consistent (EC) approximate inference framework [OW05] builds upon. There, so-called *single-loop algorithms* are derived, which are in the given case equivalent to VAMP or variants of it. Furthermore, it is suggested to optimize the Lagrange dual (4.25) partially, which leads to so-called *double-loop algorithms*.

The connection between EC, as well as EP (expectation propagation) [Min01] and convex optimization has already been shown in [HOW⁺05] for general inference problems.

4.3 Turbo-Type Algorithms—Optimization-Based Derivation

Having the stationary points of the Lagrangian, we can now turn to the derivation of algorithms by iteratively following the restrictions imposed by the structure of the densities.



Figure 4.2: Factor graph for CS of the Turbo kind with respective messages.

The optimization-based view on the problem gives us already strict connections between the (natural) parameters, cf. (4.22). In Chapter 3, the use of this computation was called *postprocessed* extrinsic, because an outgoing message on an edge could be obtained by subtracting the influence of the incoming message on the same edge from the result of the estimation that was obtained from all incoming messages. In the given case, each factor node is connected with only one edge to the variable node. Hence, for the given structure, the extrinsic computed at one factor node is directly the input to the other factor as indicated in (4.22). The respective forwarding of the messages is depicted in Fig. 4.2.

The mapping of the estimates $E_{\mathbf{x}\sim q_s}\{g(\mathbf{x})\}$, respectively $E_{\mathbf{x}\sim q_c}\{g(\mathbf{x})\}$ to the natural parameters $\boldsymbol{\theta}$ depend on the form of the exponential family, i.e., on the definition of the sufficient statistics $g(\boldsymbol{x})$. In the following, two cases for the specification of $g(\boldsymbol{x})$ are considered, starting with average variances.

4.3.1 Average Variances—VAMP

The aspired estimate (2.9) is a conditional *mean* that specifies the value of the components of the signal vector x for the given observations y. The advantage of working with a Gaussian density as exponential family is that it additionally specifies the variance or second order moment, thereby introducing a measure of reliability for the random variable, thus, also for the estimate.

In Chapter 3, the individual variance case was considered, where each estimate had its own variance, i.e., reliability. Due to the detailed view on the factor graph and the resulting scalar estimations this specification was necessary. In the given case, where always the entire vector is processed, this constraint can be relaxed to considering a variance that represents the average reliability of the entire vector, which yields a complexity reduction in terms of the number of messages to be sent.

The average variance case of a Gaussian density is specified in (D.11), i.e.,

$$\boldsymbol{g}(\boldsymbol{x}) = [\boldsymbol{g}_{\boldsymbol{\lambda}}(\boldsymbol{x})^{\top}, \, g_{\boldsymbol{\Lambda}}(\boldsymbol{x})^{\top}]^{\top}$$
(4.26)

with

$$\boldsymbol{g}_{\lambda}(\boldsymbol{x}) = \boldsymbol{x} , \quad \boldsymbol{g}_{\Lambda}(\boldsymbol{x}) = -\frac{1}{2}\boldsymbol{x}^{\top}\boldsymbol{x} .$$
 (4.27)

This specification will eventually lead to the VAMP algorithm [RSF19].

Let us consider the moment computed at the factor nodes first. Obviously, the estimation at the channel factor is related to the LMMSE estimation in Sec. 2.2.1.1. To compare with those estimators, regard that an average variance at the input is equivalent to $\boldsymbol{\Phi}_{x} = \sigma_{x}^{2} \boldsymbol{I}_{N}$. With the identification $\boldsymbol{\nu}_{c} = [\boldsymbol{m}_{x}^{\top}/\sigma_{x}^{2}, 1/\sigma_{x}^{2}]^{\top}$, it is $(\operatorname{trace}(\boldsymbol{M}) = \sum_{j=1}^{N} [\boldsymbol{M}]_{jj})$

$$\mathbf{E}_{\mathbf{x} \sim \mathbf{q}_{\mathrm{c}}} \{ \boldsymbol{g}_{\boldsymbol{\lambda}}(\mathbf{x}) \} = \boldsymbol{m}_{\mathrm{c}} , \qquad (4.28)$$

$$\mathbf{E}_{\mathbf{x} \sim \mathbf{q}_{c}} \{ g_{A}(\mathbf{x}) \} = -\frac{1}{2} \operatorname{trace} \left(\mathbf{E}_{\mathbf{x}} \{ \mathbf{x} \mathbf{x}^{\top} \mid \mathbf{y} \} \right) = -\frac{1}{2} \operatorname{trace} \left(\mathbf{\Phi}_{c} + \mathbf{m}_{c} \mathbf{m}_{c}^{\top} \right).$$
(4.29)

The resulting average variance is then given by

$$\sigma_{\rm c}^2 \stackrel{\text{\tiny def}}{=} \frac{1}{N} \operatorname{trace}(\boldsymbol{\Phi}_{\rm c}) = -\frac{2}{N} \operatorname{E}_{\mathbf{x} \sim \mathbf{q}_{\rm c}} \{g_A(\mathbf{x})\} - \boldsymbol{m}_{\rm c}^{\top} \boldsymbol{m}_{\rm c} .$$
(4.30)

The mapping to the natural parameter θ of the variable node is then obtained via

$$\boldsymbol{\theta} = [\boldsymbol{m}_{c}^{\top}/\sigma_{c}^{2}, 1/\sigma_{c}^{2}]^{\top}, \qquad (4.31)$$

and the extrinsic (which serves as input for the signal factor) is obtained via (4.22), i.e.,

$$\boldsymbol{\nu}_{\rm s} = \boldsymbol{\theta} - \boldsymbol{\nu}_{\rm c} \; . \tag{4.32}$$

For the signal factor, the estimator similarly relates to (2.19) and (2.20) from Sec. 2.2.1.3 with $\boldsymbol{\nu}_{\rm s} = [\tilde{\boldsymbol{x}}_{\rm s}^{\top}/\tilde{\sigma}_{\rm s}^2, 1/\tilde{\sigma}_{\rm s}^2]^{\top}$, i.e., in the average variance case it is $\tilde{\sigma}_{{\rm s},j}^2 = \tilde{\sigma}_{\rm s}^2$ for all $j \in \{1, \ldots, N\}$, as follows

$$\mathbf{E}_{\mathbf{x}\sim q_{\mathrm{s}}}\{\boldsymbol{g}_{\boldsymbol{\lambda}}(\mathbf{x})\} = \mathbf{E}_{\mathbf{x}}\{\mathbf{x} \mid \tilde{\boldsymbol{x}}_{\mathrm{s}}\} = \boldsymbol{m}_{\mathrm{s}} , \qquad (4.33)$$

$$\mathbf{E}_{\mathbf{x}\sim \mathbf{q}_{s}}\{g_{A}(\mathbf{x})\} = -\frac{1}{2} \operatorname{trace}\left(\mathbf{E}_{\mathbf{x}}\left\{\mathbf{x}\mathbf{x}^{\top} \mid \tilde{\mathbf{x}}_{s}\right\}\right) = -\frac{1}{2} \sum_{j=1}^{N} (\sigma_{s,j}^{2} + m_{s,j}^{2}) , \qquad (4.34)$$

and a respective average variance is obtained via

$$\sigma_{\mathrm{s}}^{2} \stackrel{\text{\tiny def}}{=} \frac{1}{N} \sum_{j=1}^{N} \sigma_{\mathrm{s},j}^{2} = -\frac{2}{N} \mathrm{E}_{\mathbf{x} \sim \boldsymbol{q}_{\mathrm{s}}} \{ g_{A}(\mathbf{x}) \} - \boldsymbol{m}_{\mathrm{s}}^{\top} \boldsymbol{m}_{\mathrm{s}} .$$

$$(4.35)$$

The mapping is analogously $\boldsymbol{\theta} = [\boldsymbol{m}_{s}^{\top}/\sigma_{s}^{2}, 1/\sigma_{s}^{2}]^{\top}$ and the extrinsic calculation $\boldsymbol{\nu}_{c} = \boldsymbol{\theta} - \boldsymbol{\nu}_{s}$, respectively.

Given a suitable initialization, the iterative processing of the steps above yields a recovery algorithm for the CS problem, which is termed VAMP in the literature. The entire algorithm is stated in pseudo code in Appendix E.5. Note that the use of natural parameters is omitted in the description by directly computing mean and variance via the respective transformations.

Noteworthy, an algorithm called *orthogonal approximate message passing* (OAMP) [MP17] was proposed simultaneously to VAMP, essentially identical.

4.3.2 Individual Variances

The VAMP algorithm can simply be varied by specifying the sufficient statistics differently. Two cases have been considered in [FSRS16]: so-called *uniform diagonalization*, which equals the specification we termed *average variance*, and *vector-valued diagonalization*, leading to *individual variances*. This second case may lead to better insights into the estimates, as the reliability is evaluated on an individual level. However, the complexity increases in terms of messages—from N + 1-dimensional to 2N-dimensional vectors—and also the computation of the LMMSE estimate gets more difficult. In [FSRS16], generalizations of the EC framework [OW05], i.e., variants of VAMP, are considered for general inference problems, but not the compressed sensing problem itself.

By changing the specification, the basic processing steps stay the same as in VAMP; only the moment computation and the mapping from and to natural parameters differ.

Mathematically, the sufficient statistic need to be defined by $g(x) = [g_{\lambda}(x)^{\top}, g_{\Lambda}(x)^{\top}]^{\top}$ with

$$\boldsymbol{g}_{\boldsymbol{\lambda}}(\boldsymbol{x}) = \boldsymbol{x} \; , \tag{4.36}$$

$$\boldsymbol{g}_{\boldsymbol{\Lambda}}(\boldsymbol{x}) = -\frac{1}{2} \operatorname{diag}\left(\boldsymbol{x}\boldsymbol{x}^{\top}\right).$$
 (4.37)

The individual variances case for the estimates in (2.10) and (2.11) is obtained by specifying $\Phi_{x} = \operatorname{diag}(\sigma_{x,j}^{2})$. To compare with the estimators from Sections 2.2.1.1 and 2.2.1.3, the input variables are identified by

$$\boldsymbol{\nu}_{\rm c} = [m_{{\rm x},1}/\sigma_{{\rm x},1}^2, \ldots, m_{{\rm x},N}/\sigma_{{\rm x},N}^2, 1/\sigma_{{\rm x},1}^2, \ldots, 1/\sigma_{{\rm x},N}^2]^{\top}, \qquad (4.38)$$

respectively,

$$\boldsymbol{\nu}_{\rm s} = [\tilde{x}_{{\rm s},1}/\tilde{\sigma}_{{\rm s},1}^2, \dots, \tilde{x}_{{\rm s},N}/\tilde{\sigma}_{{\rm s},N}^2, 1/\tilde{\sigma}_{{\rm s},1}^2, \dots, 1/\tilde{\sigma}_{{\rm s},N}^2]^{\top} .$$
(4.39)

The expectations w.r.t. $g_{\lambda}(x)$ are the same as for the average variances case. For the second order moment, the channel-constrained estimation reads

$$E_{\mathbf{x} \sim q_{c}} \{ \boldsymbol{g}_{\boldsymbol{\Lambda}}(\mathbf{x}) \} = -\frac{1}{2} \operatorname{diag} \left(\boldsymbol{\varPhi}_{c} + \boldsymbol{m}_{c} \boldsymbol{m}_{c}^{\top} \right), \qquad (4.40)$$

respectively

$$E_{\mathbf{x}\sim q_{s}}\{\boldsymbol{g}_{\boldsymbol{\Lambda}}(\mathbf{x})\} = -\frac{1}{2}\operatorname{diag}\left(\operatorname{diag}\left(\sigma_{s,j}^{2} + m_{s,j}^{2}\right)\right), \qquad (4.41)$$

for the signal-constrained estimation. The mapping to the natural parameters is then given by

$$\boldsymbol{\theta} = [m_{\bullet,1}/\sigma_{\bullet,1}^2, \dots, m_{\bullet,N}/\sigma_{\bullet,N}^2, 1/\sigma_{\bullet,1}^2, \dots, 1/\sigma_{\bullet,N}^2]^\top$$
(4.42)

with $\bullet \in \{c, s\}$, respectively. The extrinsic calculation is the same as before, i.e., subtraction of the natural parameters. The entire algorithm is stated in Algorithm E.6; we call it VAMPind, for VAMP with <u>ind</u>ividual variances. Again, the representation with means and variances is used in the statement of the algorithm.

4.4 Estimation-Theoretic Bias Compensation

The simple additive connection (4.22) between the natural parameters of variable and factor nodes turns out to be crucial for the processing in the algorithms. So far, it was called the computation of *extrinsic*, since previously known values are subtracted [BG96, GH11].

In [SF17, SF18] a connection to *bias compensation* [For03] from an estimation theoretical point of view was drawn. This already shows that different perspectives can yield the same results, but also open the view for more possibilities, as will be visible soon.

Bias compensation describes the task of removing the *bias*, which is a systematic offset in the estimate, from the estimate. The procedure usually degrades the reliability of the estimate [For03, Fis02]. MMSE estimates are always biased, since they are by definition orthogonal to the error. This causes a part of the signal (to be estimated) to be accounted for the error [For03]. The compensation is then usually realized by scaling the estimate adequately [Fis02].

In the following, forms of post-processing for the extrinsic based on the estimation-theoretical viewpoint of bias compensation are derived and compared to the extrinsic calculations obtained so far. Similar derivations can be found in [SF18, Spa19, FS22, SF22].

There are two kinds of estimations to be considered; the linear estimation which is used for the channel part of CS and the non-linear estimation as in the signal part; starting with the derivations for the channel-constrained estimations, i.e., m_c as defined in (2.10).

4.4.1 Channel-Constrained Estimation

The bias in the channel-constrained estimate

$$\boldsymbol{m}_{c} = \boldsymbol{m}_{x} + \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} + \sigma_{n}^{2} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \boldsymbol{m}_{x})$$
(4.43)

can be seen in the fact that the end-to-end cascade K between signal x and estimate $m_{\rm c}$, which is given by [Spa19, FS22]

$$\begin{aligned} \boldsymbol{K} &= \boldsymbol{\Phi}_{\mathsf{x}} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{\Phi}_{\mathsf{x}} \boldsymbol{A}^{\top} + \sigma_{\mathsf{n}}^{2} \boldsymbol{I}_{M})^{-1} \boldsymbol{A} \\ &= (\boldsymbol{A}^{\top} \boldsymbol{A} + \sigma_{\mathsf{n}}^{2} \boldsymbol{\Phi}_{\mathsf{x}}^{-1})^{-1} \boldsymbol{A}^{\top} \boldsymbol{A} \\ &= \boldsymbol{I}_{N} - \boldsymbol{\Phi}_{\mathsf{c}} \boldsymbol{\Phi}_{\mathsf{x}}^{-1} , \end{aligned}$$
(4.44)

is not an identity matrix, which can be easily verified reformulating $\boldsymbol{\Phi}_{c}$ from (2.11) via the Sherman-Morrison-Woodbury identity [GV96]. A visualization of the end-to-end cascade by block diagrams is given in Fig. 4.3. The compensation for an unbiased estimate $\boldsymbol{m}_{c,u}$ can be done by inserting a respective scaling matrix \boldsymbol{W} in the block diagram (shown in Fig. 4.4) under various restrictions.

4.4.1.1 Average Unbiasing

First, a rather loose requirement is considered; the effect of the end-to-end cascade should be compensated *on average* [Fis16]. This requirement is suitable as long as the energies (squared



Figure 4.3: Visualization of the end-to-end cascade K between signal x and estimate m_c .



Figure 4.4: Visualization of the bias compensation procedure with scaling matrix W.

 ℓ_2 -norms) of the columns of the sensing matrix as well as the variances of signal components of \boldsymbol{x} do not vary much. For this procedure, it is therefore assumed that $\boldsymbol{\Phi}_{\mathsf{x}} = \sigma_{\mathsf{x}}^2 \boldsymbol{I}_N$. Mathematically, the requirement means, a scaling matrix \boldsymbol{W} must fulfill [Fis16, FS22]

$$\frac{1}{N} \operatorname{trace}(\boldsymbol{W} \cdot \boldsymbol{K}) \stackrel{!}{=} 1 .$$
(4.45)

The solution is

$$\boldsymbol{W} = \frac{1}{\frac{1}{N} \operatorname{trace}(\boldsymbol{K})} \boldsymbol{I}_N , \qquad (4.46)$$

i.e., there is a single unbiasing factor $w = N/\text{trace}(\mathbf{K})$, which is used to (re-)scale all elements of \mathbf{m}_{c} . The resulting <u>unbiased</u> estimate $\mathbf{m}_{c,u}$ is then given by

$$\boldsymbol{m}_{\mathrm{c,u}} = \boldsymbol{m}_{\mathrm{x}} + \boldsymbol{W}\boldsymbol{\Phi}_{\mathrm{x}}\boldsymbol{A}^{\mathrm{T}} (\boldsymbol{A}\boldsymbol{\Phi}_{\mathrm{x}}\boldsymbol{A}^{\mathrm{T}} + \sigma_{\mathrm{n}}^{2}\boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A}\boldsymbol{m}_{\mathrm{x}}) .$$
(4.47)

Noteworthy, with $\boldsymbol{\Phi}_{\mathsf{x}} = \sigma_{\mathsf{x}}^2 \boldsymbol{I}_N$, one obtains from (4.44)

$$w = \frac{N}{\operatorname{trace}(\boldsymbol{K})} = \frac{N}{\operatorname{trace}(\boldsymbol{I}_N - \boldsymbol{\Phi}_c \boldsymbol{\Phi}_x^{-1})}$$
$$= \frac{1}{1 - \frac{1}{\sigma_x^2} \frac{1}{N} \operatorname{trace}(\boldsymbol{\Phi}_c)}$$
$$= \frac{1}{1 - \frac{\sigma_c^2}{\sigma_x^2}} = \frac{\sigma_x^2}{\sigma_x^2 - \sigma_c^2} .$$
(4.48)

The covariance matrix corresponding to estimate $m_{
m c,u}$ reads [Spa19]

$$\boldsymbol{\Phi}_{c,u} = E_{\mathbf{x}} \Big\{ (\mathbf{x} - \boldsymbol{m}_{c,u}) (\mathbf{x} - \boldsymbol{m}_{c,u})^\top \mid \boldsymbol{y} \Big\} = \boldsymbol{\Phi}_{\mathbf{x}} - \boldsymbol{\Phi}_{\mathbf{x}} \boldsymbol{K}^\top \boldsymbol{W}^\top - \boldsymbol{W} \boldsymbol{K} \boldsymbol{\Phi}_{\mathbf{x}} + \boldsymbol{W} \boldsymbol{\Phi}_{\mathbf{x}} \boldsymbol{K}^\top \boldsymbol{W}^\top .$$
(4.49)

The average unbiased variance therefore reads (trace ($\boldsymbol{W}\cdot\boldsymbol{K})=\text{trace}(\boldsymbol{K}^{\top}\cdot\boldsymbol{W}^{\top})=N$)

$$\sigma_{\rm c,u}^2 = \frac{1}{N} \operatorname{trace}(\boldsymbol{\Phi}_{\rm c,u}) = \sigma_{\rm x}^2 - \sigma_{\rm x}^2 - \sigma_{\rm x}^2 + w\sigma_{\rm x}^2 = \sigma_{\rm x}^2(w-1)$$
$$= \sigma_{\rm x}^2 \left(\frac{\sigma_{\rm x}^2}{\sigma_{\rm x}^2 - \sigma_{\rm c}^2} - 1\right) = \sigma_{\rm x}^2 \frac{\sigma_{\rm x}^2 - \sigma_{\rm x}^2 + \sigma_{\rm c}^2}{\sigma_{\rm x}^2 - \sigma_{\rm c}^2} = \frac{\sigma_{\rm x}^2 \sigma_{\rm c}^2}{\sigma_{\rm x}^2 - \sigma_{\rm c}^2} = \left(\frac{1}{\sigma_{\rm c}^2} - \frac{1}{\sigma_{\rm x}^2}\right)^{-1} .$$
(4.50)

So, the unbiased estimate can be written as

$$\boldsymbol{m}_{c,u} = \boldsymbol{m}_{x} + \frac{\sigma_{x}^{2}}{\sigma_{x}^{2} - \sigma_{c}^{2}} \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} + \sigma_{n}^{2} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \boldsymbol{m}_{x})$$

$$= \left(1 - \frac{\sigma_{x}^{2}}{\sigma_{x}^{2} - \sigma_{c}^{2}} + \frac{\sigma_{x}^{2}}{\sigma_{x}^{2} - \sigma_{c}^{2}}\right) \boldsymbol{m}_{x} + \frac{\sigma_{x}^{2}}{\sigma_{x}^{2} - \sigma_{c}^{2}} \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} + \sigma_{n}^{2} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \boldsymbol{m}_{x})$$

$$= \frac{\sigma_{c}^{2}}{\sigma_{c}^{2} - \sigma_{x}^{2}} \boldsymbol{m}_{x} + \frac{\sigma_{x}^{2}}{\sigma_{x}^{2} - \sigma_{c}^{2}} \boldsymbol{m}_{c} = \sigma_{c,u}^{2} (\boldsymbol{m}_{c} / \sigma_{c}^{2} - \boldsymbol{m}_{x} / \sigma_{x}^{2}) .$$
(4.51)

By slightly rearranging these equations to

$$\boldsymbol{m}_{c,u}/\sigma_{c,u}^2 = \boldsymbol{m}_c/\sigma_c^2 - \boldsymbol{m}_x/\sigma_x^2$$
, $1/\sigma_{c,u}^2 = 1/\sigma_c^2 - 1/\sigma_x^2$, (4.52)

one can make the connection $\tilde{\boldsymbol{x}}_{c} = \boldsymbol{m}_{x}$, respectively $\tilde{\sigma}_{c}^{2} = \sigma_{x}^{2}$, as well as $\boldsymbol{m}_{c,u} = \tilde{\boldsymbol{x}}_{s}$ and $\sigma_{c,u}^{2} = \tilde{\sigma}_{s}^{2}$, to see that the bias compensation procedure equals the exchange after the linear estimation in the VAMP algorithm. This means that VAMP performs (at this point) average bias compensation with an average variance.

4.4.1.2 Individual Unbiasing

A stricter requirement is the compensation of the individual scaling effects that K causes on the signal components. In this case, the scaling matrix W must fulfill [FS07, FS22]

$$\operatorname{diag}(\boldsymbol{W}\cdot\boldsymbol{K}) \stackrel{!}{=} \boldsymbol{I}_N . \tag{4.53}$$

This is obtained for scaling matrix

$$\boldsymbol{W} = \operatorname{diag}(1/[\boldsymbol{K}]_{jj}) \ . \tag{4.54}$$

Being more careful in the compensation, allows us to handle variations in the energies (squared ℓ_2 -norms) of both, columns of the sensing matrix, as well as the signal components itself. This procedure is, thus, suitable for the use of individual variances, i.e., the more general case $\boldsymbol{\Phi}_{x} = \operatorname{diag}(\sigma_{x,j}^2)$ is considered. The diagonal entries of the compensation matrix can be calculated via (4.44) [Spa19]

$$[\mathbf{W}]_{jj} = \frac{1}{[\mathbf{K}]_{jj}} = \frac{1}{1 - \sigma_{c,j}^2 / \sigma_{x,j}^2} = \frac{\sigma_{x,j}^2}{\sigma_{x,j}^2 - \sigma_{c,j}^2}, \qquad (4.55)$$

which resembles the unbiasing factor before, only with the individual variances. The (individual) unbiased variances compute to $([\mathbf{W}]_{jj} \cdot [\mathbf{K}]_{jj} = 1)$

$$\sigma_{c,u,j}^{2} \stackrel{\text{def}}{=} [\boldsymbol{\Phi}_{c,u}]_{jj} = \sigma_{x,j}^{2} - \sigma_{x,j}^{2} - \sigma_{x,j}^{2} + [\boldsymbol{W}]_{jj}\sigma_{x,j}^{2} = \sigma_{x,j}^{2}([\boldsymbol{W}]_{jj} - 1)$$
$$= \sigma_{x,j}^{2} \left(\frac{\sigma_{x,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}} - 1\right) = \frac{\sigma_{x,j}^{2}\sigma_{c,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}} = \left(\frac{1}{\sigma_{c,j}^{2}} - \frac{1}{\sigma_{x,j}^{2}}\right)^{-1}, \quad (4.56)$$

leading to an unbiased estimate of the form

$$\boldsymbol{m}_{c,u} = \boldsymbol{m}_{x} + \operatorname{diag}\left(\frac{\sigma_{x,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}}\right) \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} + \sigma_{n}^{2} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \boldsymbol{m}_{x})$$

$$= \left(\boldsymbol{I}_{N} - \operatorname{diag}\left(\frac{\sigma_{x,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}}\right) + \operatorname{diag}\left(\frac{\sigma_{x,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}}\right)\right) \boldsymbol{m}_{x}$$

$$+ \operatorname{diag}\left(\frac{\sigma_{x,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}}\right) \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{\Phi}_{x} \boldsymbol{A}^{\top} + \sigma_{n}^{2} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \boldsymbol{m}_{x})$$

$$= \operatorname{diag}\left(\frac{\sigma_{c,j}^{2}}{\sigma_{c,j}^{2} - \sigma_{x,j}^{2}}\right) \boldsymbol{m}_{x} + \operatorname{diag}\left(\frac{\sigma_{x,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}}\right) \boldsymbol{m}_{c}, \qquad (4.57)$$

i.e., the components of the estimate compute to $(j \in \{1, ..., N\})$

$$m_{c,u,j} = \frac{\sigma_{x,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}} m_{c,j} - \frac{\sigma_{c,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}} m_{x,j} = \frac{\sigma_{x,j}^{2} \sigma_{c,j}^{2}}{\sigma_{x,j}^{2} - \sigma_{c,j}^{2}} \left(\frac{m_{c,j}}{\sigma_{c,j}^{2}} - \frac{m_{x,j}}{\sigma_{x,j}^{2}}\right)$$
$$= \sigma_{c,u,j}^{2} \left(\frac{m_{c,j}}{\sigma_{c,j}^{2}} - \frac{m_{x,j}}{\sigma_{x,j}^{2}}\right), \qquad (4.58)$$

or slightly rephrased

$$m_{\mathrm{c},\mathrm{u},j}/\sigma_{\mathrm{c},\mathrm{u},j}^2 = m_{\mathrm{c},j}/\sigma_{\mathrm{c},j}^2 - m_{\mathrm{x},j}/\sigma_{\mathrm{x},j}^2$$
, $1/\sigma_{\mathrm{c},\mathrm{u},j}^2 = 1/\sigma_{\mathrm{c},j}^2 - 1/\sigma_{\mathrm{x},j}^2$. (4.59)

Hence, there is a connection to the recovery algorithm with individual variances from Sec. 4.3.2, when making the identifications $m_{c,u,j} = \tilde{x}_{s,j}$, $\sigma_{c,u,j}^2 = \tilde{\sigma}_{s,j}^2$ and $m_{x,j} = \tilde{x}_{c,j}$, respectively, $\sigma_{x,j}^2 = \tilde{\sigma}_{c,j}^2$.

The resulting average unbiased variance of individual unbiasing computes with the definition of the *arithmetic mean* as $M_A(x_j) = \frac{1}{N} \sum_{j=1}^N x_j$ to [Fis16]

$$\sigma_{\mathrm{c,u,i}}^2 = \frac{1}{N} \sum_{j=1}^N \sigma_{\mathrm{c,u,j}}^2 = \mathrm{M}_{\mathrm{A}} \left(\sigma_{\mathrm{x,j}}^2 [\boldsymbol{W}]_{jj} \right) - \mathrm{M}_{\mathrm{A}} \left(\sigma_{\mathrm{x,j}}^2 \right) \,. \tag{4.60}$$

In case of the specification by a single (average) variance, i.e., $\boldsymbol{\Phi}_{\mathsf{x}} = \sigma_{\mathsf{x}}^{2} \boldsymbol{I}_{N}$, with harmonic mean $M_{\mathrm{H}}(x_{j}) = \frac{N}{\sum_{j=1}^{N} 1/x_{j}} = \frac{1}{M_{\mathrm{A}}(1/x_{j})}$ this simplifies to

$$\sigma_{\rm c,u,i}^2 = \sigma_{\rm x}^2 \left(M_{\rm A}([\boldsymbol{W}]_{jj}) - 1 \right) = \sigma_{\rm x}^2 \left(M_{\rm A}(1/[\boldsymbol{K}]_{jj}) - 1 \right) = \sigma_{\rm x}^2 \left(\frac{1}{M_{\rm H}([\boldsymbol{K}]_{jj})} - 1 \right) .$$
(4.61)

Comparing this with the average variance (4.50) for average bias compensation

$$\sigma_{\rm c,u}^2 = \sigma_{\rm x}^2(w-1) = \sigma_{\rm x}^2 \left(\frac{N}{\rm trace}(\mathbf{K}) - 1\right) = \sigma_{\rm x}^2 \left(\frac{1}{M_{\rm A}([\mathbf{K}]_{jj})} - 1\right) , \qquad (4.62)$$

shows that $\sigma_{c,u}^2 < \sigma_{c,u,i}^2$, since $M_A([\boldsymbol{K}]_{jj}) > M_H([\boldsymbol{K}]_{jj})$ for varying $[\boldsymbol{K}]_{jj}$.

Noteworthy, in the optimization-based view, average variances and average unbiasing is inherently connected and so are individual variances and individual unbiasing. From the estimation-theoretic perspective, there is no need in restricting to these combinations. A respective consideration of individual variances with average unbiasing is unlikely to yield good results, as the "loose" requirement of compensating only the average effect of the end-to-end cascade K cannot capture the variability in the variances $\sigma_{x,j}^2$. This case will therefore not be considered here. However, at first sight, there is no reason not to use individual unbiasing, when an average variance is present (case $\Phi_x = \sigma_x^2 I_N$). Hence, algorithms with this combination will be composed below. Before that, the non-linear estimation and the respective unbiasing procedure are considered from an estimation-theoretic perspective.

4.4.2 Signal-Constrained Estimation

In this section, we recapitulate suitable bias compensation techniques for the scalar non-linear MMSE estimation² (2.19) similar to [Fis16, SF18, Spa19, FSG20, FS22], where the derivation is based on a geometric point of view as used in [Kay93, Fis02, For03]. Random variables allow a geometrical interpretation, since joint estimation can be interpreted as inner product of a vector space [Kay93, Sec. 12.4]. To understand this, recall that the observation model for this case is given by (cf. also (2.18))

$$\tilde{\mathsf{x}}_{\mathrm{s}} = \mathsf{x} + \mathsf{n}_{\mathrm{s}}, \quad \text{with} \quad \mathsf{x} \sim \mathsf{f}_{\mathsf{x}}(x), \quad \mathsf{n}_{\mathrm{s}} \sim \mathcal{N}\left(0, \sigma_{\mathsf{n}_{\mathrm{s}}}^{2}\right).$$
 (4.63)

Note that the noise variance is denoted by $\sigma_{n_s}^2$ here, instead of $\tilde{\sigma}_s^2$, as in (2.18). The denomination shall stress the connection to the observation noise n_s . For the algorithms, $\tilde{\sigma}_s^2$ will be used again to point out the connection to \tilde{x}_s .

In CS, the sparse prior **x** usually is assumed to have zero-mean, i.e., $m_x = E_x\{x\} = 0$. Furthermore, assuming that the noise / error n_s is *uncorrelated* to the signal, it is [Kay93]

$$E_{xn_s}\{xn_s\} = E_x\{x\}E_{n_s}\{n_s\} = 0, \qquad (4.64)$$

where the first estimation is w.r.t. the joint pdf $f_{xn_s}(x, n_s)$. Interpreting the joint expectation $E_{xn_s}{xn_s}$ as inner product, implies that the random variables x and n_s are *orthogonal*, when considered as vectors in a vector space, cf. [PP02, p. 211]; for a summary of the notion of vector spaces, the reader may be referred to Appendix A. The norm of the vector space is accordingly given by "joint" expectation of the random variable with itself, i.e., when considering a Euclidean vector space, the squared "length", i.e., the squared ℓ_2 -norm of the vector representing observation noise n_s is given by

$$\mathbf{E}_{\mathsf{n}_{\mathrm{s}}}\left\{\mathsf{n}_{\mathrm{s}}^{2}\right\} = \sigma_{\mathsf{n}_{\mathrm{s}}}^{2} \,. \tag{4.65}$$

Respectively, the signal x has a squared length of $\sigma_x^2 = E_x \{x^2\}$.

Let us get back to the MMSE criterion. There, the estimation error $e_s = x - m_s$ between signal x and estimate m_s is considered. The criterion aims to minimize the squared version of this error given the observation \tilde{x}_s , i.e.,

$$m_{\rm s} = \underset{m}{\operatorname{argmin}} \operatorname{E}_{\mathsf{x}} \left\{ (\mathsf{x} - m)^2 \mid \tilde{x}_{\rm s} \right\} \,. \tag{4.66}$$

²For simplicity, the index j is omitted in this section.

The resulting conditional mean is justified by the fact that [Kay93, PP02]

$$\frac{\partial}{\partial m} \mathbf{E}_{\mathsf{x}} \left\{ (\mathsf{x} - m)^2 \mid \tilde{x}_{\mathsf{s}} \right\} = -2(\mathbf{E}_{\mathsf{x}} \{ \mathsf{x} \mid \tilde{x}_{\mathsf{s}} \} - m) \stackrel{!}{=} 0 , \qquad (4.67)$$

requires $m = m_s = E_x \{x \mid \tilde{x}_s\}$. The corresponding reliability of the estimate for the given observation is defined as [PP02, GWSV11]

$$\sigma_{\rm s}^2 = {\rm E}_{\rm x} \{ ({\rm x} - m_{\rm s})^2 \mid \tilde{x}_{\rm s} \} = {\rm E}_{\rm x} \{ {\rm x}^2 \mid \tilde{x}_{\rm s} \} - m_{\rm s}^2 , \qquad (4.68)$$

cf. also (2.20). As both estimates are functions of \tilde{x}_s , we may in the following use the notation $m_s(\tilde{x}_s)$ and $\sigma_s^2(\tilde{x}_s)$ to emphasize this dependency. The so-called (minimum) *mean-squared error* (MSE) [Kay93] of this strategy is obtained by integration over the statistic of the observation \tilde{x}_s , i.e., it defines the observation-independent reliability of the estimation, [Kay93, GWSV11]

$$\sigma^{2} \stackrel{\text{\tiny def}}{=} \operatorname{E}_{\tilde{\mathsf{x}}_{\mathrm{s}}} \left\{ \sigma_{\mathrm{s}}^{2}(\tilde{\mathsf{x}}_{\mathrm{s}}) \right\} = \int \mathsf{f}_{\tilde{\mathsf{x}}_{\mathrm{s}}}(\tilde{x}_{\mathrm{s}}) \sigma_{\mathrm{s}}^{2}(\tilde{x}_{\mathrm{s}}) \, \mathrm{d}\tilde{x}_{\mathrm{s}} \; . \tag{4.69}$$

Calculating this integration within a recovery algorithm would lead to an unnecessary complexity overhead. Since the MSE is only a function of the variance $\sigma_{n_s}^2$ ($\tilde{\sigma}_s^2$ within the algorithms), the computation can be done beforehand and stored in a look-up table (LUT).

The MSE σ^2 is depicted over the input $\sigma_{n_s}^2$ (a.k.a. $\tilde{\sigma}_s^2$) in Fig. 4.5 for the distribution (2.4), which is called discrete ternary (DT) prior, for N = 250. Also, two bounds from [GWSV11] are given. The rougher one (dashed) states that the MSE σ^2 is always smaller than the input variance $\sigma_{n_s}^2$ and the variance of the prior σ_x^2 . The second bound (dotted) is the MSE of a Gaussian random variable, which reads

$$\sigma^{2} = (1/\sigma_{x}^{2} + 1/\sigma_{n_{s}}^{2})^{-1} = \frac{\sigma_{x}^{2}}{1 + \sigma_{x}^{2}/\sigma_{n_{s}}^{2}} = \frac{\sigma_{x}^{2}\sigma_{n_{s}}^{2}}{\sigma_{n_{s}}^{2} + \sigma_{x}^{2}}.$$
(4.70)

According to [GWSV11], the MSE of a Gaussian random variable is always larger (or equal) than the MSE of any other random variable with the same variance σ_x^2 . Noteworthy, this feature is connected to the maximum-entropy property of Gaussian random variables. On the left side, the MSE of the DT prior is compared to the two bounds and the Bernoulli-Gaussian (BG) prior (2.5) with a sparsity of s = 5, i.e., $\sigma_x^2 = s/N = 5/250 = 0.02$. One can see that the MSEs of both priors converge to the bound given by σ_x^2 for large values of the input variance $\sigma_{n_s}^2$. It is also visible that the DT prior has a tremendous decrease in the MSE for small input variances, when compared to the BG prior; only around $\sigma_{n_s}^2 = 10^{-1}$ the MSE of the DT prior is larger. The reason for the faster descent is the discrete nature of the prior, which simplifies the estimation in less noise. On the right side, the DT prior is compared for different sparsities that will be used in the simulations in Chapter 6. According to the bounds, the value of convergence towards large input variances changes, but the descents towards small input variances is similar. For more properties of the MSE, cf. [GSV05, GWSV11].

It can be shown that joint estimation of estimation error e_s and observation \tilde{x}_s , in fact any function of \tilde{x}_s , is zero, i.e., [Fis16, PP02]

$$\begin{split} \mathbf{E}_{\mathbf{e}_{s}\tilde{\mathbf{x}}_{s}}\{\mathbf{e}_{s}f(\tilde{\mathbf{x}}_{s})\} &= \mathbf{E}_{\mathbf{x}\tilde{\mathbf{x}}_{s}}\{(\mathbf{x} - m_{s}(\tilde{\mathbf{x}}_{s}))f(\tilde{\mathbf{x}}_{s})\} = \mathbf{E}_{\tilde{\mathbf{x}}_{s}}\{f(\tilde{\mathbf{x}}_{s})\mathbf{E}_{\mathbf{x}}\{\mathbf{x} - m_{s}(\tilde{\mathbf{x}}_{s}) \mid \tilde{\mathbf{x}}_{s}\}\}\\ &= \mathbf{E}_{\tilde{\mathbf{x}}_{s}}\{f(\tilde{\mathbf{x}}_{s})(\mathbf{E}_{\mathbf{x}}\{\mathbf{x} \mid \tilde{\mathbf{x}}_{s}\} - \mathbf{E}_{\mathbf{x}}\{m_{s}(\tilde{\mathbf{x}}_{s}) \mid \tilde{\mathbf{x}}_{s}\})\}\\ &= \mathbf{E}_{\tilde{\mathbf{x}}_{s}}\{f(\tilde{\mathbf{x}}_{s})(m_{s}(\tilde{\mathbf{x}}_{s}) - m_{s}(\tilde{\mathbf{x}}_{s}))\} = 0 , \end{split}$$
(4.71)



Figure 4.5: Comparison of the MSE for two different priors for N = 250. On the left, the priors are compared; on the right, the MSE of the DT prior (2.4) is compared for different sparsities. Note that $\sigma_{n_s}^2$ is called $\tilde{\sigma}_s^2$, when referring to the algorithmic procedures.

with $f(\tilde{x}_s)$ being an arbitrary function of \tilde{x}_s . Since m_s is a function of \tilde{x}_s , this means that the MMSE estimate is orthogonal to the resulting estimation error e_s . Naturally, also the observation \tilde{x}_s when considered a random variable, is orthogonal to e_s .

Since there are right angles between different vectors, if the estimation is not linear, this requires a three-dimensional space to be displayed, because the estimate (considered as random variable) needs to lie on a Thales circle between x and $\sigma_x^2/(\sigma_x^2 + \sigma_{n_s}^2)\tilde{x}_s$, which is the corresponding linear estimator and optimal if the signal x is Gaussian distributed, for a visualization see [For03]. The random variables are depicted as vectors in a vector space in Fig. 4.6, similarly to a figure in [FSG20]. The vector space is spanned by the direction of x, n_s, and a third axis, which we call v. The right angles are depicted as squares in the respective perspective. The vector for the estimate m_s is here constructed by the intersection of two tilted Thales circles, that span vector x and n_s, respectively, thereby ensuring the above derived right angles between the vectors. The displayed measures are squared distances, i.e., the variances of the corresponding random variables.

Since the estimation error e_s is not orthogonal to the signal x itself, i.e., a projection towards x does not yield zero, which means a part of the signal is accounted to the error. This may be interpreted as *bias*, for which two ways of compensation are introduced in the following sections.

Noteworthy, by use of total expectation, the expected value of the MMSE estimate is obtained as [Fis18]

$$E_{m_s}\{m_s\} = E_{\tilde{x}_s}\{E_x\{x \mid \tilde{x}_s\}\} = E_x\{x\} = 0, \qquad (4.72)$$

which means that on average an MMSE estimate is expectation-consistent and thereby unbiased, however, this notion of bias is not considered here.



Figure 4.6: Visualization of the random variables involved in the estimation as vectors (cf. [Fis16]).

4.4.2.1 Individual Unbiasing

The precise handling of the bias compensation procedure is obtained by scaling realizations of the random variables, which are involved in the estimation, adequately. In this work, only what is called *noise-oriented* unbiasing in [FSG20] and [Spa19] is considered, since it has better properties for the recovery in CS, for more details, see [Spa19]. This means, that instead of scaling directly the estimate m_s , the error between observation \tilde{x}_s and estimate m_s is regarded, which needs to be decomposed into a part of the actual observation noise n_s and an uncorrelated distortion d_n . For this, one needs to set up the model [Fis16, Ban13]

$$\mathbf{e}_{\mathbf{n}} = \tilde{\mathbf{x}}_{\mathbf{s}} - \mathbf{m}_{\mathbf{s}} = k_{\mathbf{n}}\mathbf{n}_{\mathbf{s}} + \mathbf{d}_{\mathbf{n}} , \qquad (4.73)$$

where k_n is the factor that is necessary to scale the observation noise such that it is uncorrelated to the distortion, i.e., $E_{n_sd_n}\{n_sd_n\} = 0$. The error e_n is already depicted in Fig. 4.6 as vector. The scaling factor can be obtained by considering similar triangles in the shaded area of Fig. 4.6. The shaded triangle of Fig. 4.6 is drawn again in Fig. 4.7, such that the triangle lies in the drawing plane. The corner points of the triangle are the endpoints of x, m_s , and \tilde{x}_s , the edges are given by n_s , e_s , and e_n . The estimation error e_s has by definition a squared length equal to the MSE, i.e., $E_{e_s}\{e_s^2\} = \sigma^2$. With the Pythagorean theorem, one can infer the respective squared length of e_n . The relations between sides of the two right triangles yields

$$\frac{\mathrm{E}_{\mathsf{n}_{\mathrm{s}}}\{(k_{\mathsf{n}}\mathsf{n}_{\mathrm{s}})^{2}\}}{\mathrm{E}_{\mathsf{e}_{\mathrm{n}}}\{\mathsf{e}_{\mathrm{n}}^{2}\}} = \frac{k_{\mathsf{n}}^{2}\sigma_{\mathsf{n}_{\mathrm{s}}}^{2}}{\sigma_{\mathsf{n}_{\mathrm{s}}}^{2} - \sigma^{2}} \stackrel{!}{=} \frac{\mathrm{E}_{\mathsf{e}_{\mathrm{n}}}\{\mathsf{e}_{\mathrm{n}}^{2}\}}{\mathrm{E}_{\mathsf{n}_{\mathrm{s}}}\{\mathsf{n}_{\mathrm{s}}^{2}\}} = \frac{\sigma_{\mathsf{n}_{\mathrm{s}}}^{2} - \sigma^{2}}{\sigma_{\mathsf{n}_{\mathrm{s}}}^{2}} , \qquad (4.74)$$

and finally

$$k_{\mathsf{n}} = \frac{\sigma_{\mathsf{n}_{\mathrm{s}}}^2 - \sigma^2}{\sigma_{\mathsf{n}_{\mathrm{s}}}^2} \,. \tag{4.75}$$



Figure 4.7: Visualization of the random variables involved in the noise-oriented bias compensation.

The scaling factor tells us, how much the observation noise n_s is scaled in order to obtain the difference e_n . Returning to the estimates, an unbiased version $m_{s,u}$ of estimate m_s is obtained by inversely applying the scaling to the difference between observation and estimate and subtracting it from the observation, i.e.,

$$m_{\rm s,u} = \tilde{x}_{\rm s} + \frac{1}{k_{\rm n}} (\tilde{x}_{\rm s} - m_{\rm s}) = \tilde{x}_{\rm s} + \frac{\sigma_{\rm n_s}^2}{\sigma_{\rm n_s}^2 - \sigma^2} (m_{\rm s} - \tilde{x}_{\rm s}) .$$
(4.76)

The corresponding unbiased variance is obtained as

$$\sigma_{\rm s,u}^2 = \sigma_{\rm s}^2 + \left(\frac{\sigma^2}{\sigma^2 - \sigma_{\rm n_s}^2}\right)^2 (m_{\rm s} - \tilde{x}_{\rm s})^2 .$$
(4.77)

This is called the *individual* unbiasing case, as it is exactly tailored for one observation \tilde{x}_s with respective noise variance $\sigma_{n_s}^2$. The case, where the reliability of an entire vector is specified in one variance is discussed in the next section.

4.4.2.2 Average Unbiasing

In the following, the connection between the previously explained *noise-oriented* unbiasing and the extrinsic calculations performed in VAMP is shown, which has similarly been stated in [SF18].

First, consider that since VAMP uses one variance $\sigma_{n_s}^2$ to adjust the non-linear estimators for all components of the signal vector \boldsymbol{x} , the corresponding MSE σ^2 is equal for all components because it is a function of the input variance only. Furthermore, for a high-dimensional signal, the empirical average over the conditional variances $\sigma_{s,j}^2$ suitably approximates the MSE, since the averaging may be interpreted as Monte-Carlo integration with many samples / observations $\tilde{x}_{s,j}$, i.e., it is

$$\sigma^{2} = \lim_{N \to \infty} \frac{1}{N} \sum_{j=1}^{N} \sigma_{\mathrm{s},j}^{2}(\tilde{x}_{\mathrm{s},j}) .$$
(4.78)
With this in mind, (4.76) can be reformulated to obtain

$$m_{\rm s,u} = \tilde{x}_{\rm s} + \frac{\sigma_{\rm n_s}^2}{\sigma_{\rm n_s}^2 - \sigma^2} (m_{\rm s} - \tilde{x}_{\rm s}) = \tilde{x}_{\rm s} + \frac{\sigma_{\rm n_s}^2 \sigma^2}{\sigma_{\rm n_s}^2 - \sigma^2} (m_{\rm s}/\sigma^2 - \tilde{x}_{\rm s}/\sigma^2)$$

$$= \left(\frac{1}{\sigma^2} - \frac{1}{\sigma_{\rm n_s}^2}\right)^{-1} \left(m_{\rm s}/\sigma^2 + \tilde{x}_{\rm s} \left(\frac{\sigma_{\rm n_s}^2 - \sigma^2}{\sigma_{\rm n_s}^2 \sigma^2} - \frac{1}{\sigma^2}\right)\right)$$

$$= (1/\sigma^2 - 1/\sigma_{\rm n_s}^2)^{-1} (m_{\rm s}/\sigma^2 - \tilde{x}_{\rm s}/\sigma_{\rm n_s}^2) , \qquad (4.79)$$

which already reminds us of the update in VAMP without indices. Note that in VAMP $\sigma_{n_s}^2$ is denoted $\tilde{\sigma}_s^2$ and as already discussed, σ_s^2 takes on the role of σ^2 .

Last but not least, the unbiased MSE is stated, which is obtained as the observation-independent version of $\sigma_{s,u}^2$, i.e., with (4.68), $E_{\tilde{x}_s}{\{\tilde{x}_s^2\}} = \sigma_x^2 + \sigma_{n_s}^2$,

$$\mathbf{E}_{\tilde{\mathbf{x}}_{\mathrm{s}}}\{m_{\mathrm{s}}(\tilde{\mathbf{x}}_{\mathrm{s}})\tilde{\mathbf{x}}_{\mathrm{s}}\} = \mathbf{E}_{\tilde{\mathbf{x}}_{\mathrm{s}}}\{(\mathbf{x} - \mathbf{e}_{\mathrm{s}})\tilde{\mathbf{x}}_{\mathrm{s}}\} = \mathbf{E}_{\tilde{\mathbf{x}}_{\mathrm{s}}}\{\mathbf{x}\tilde{\mathbf{x}}_{\mathrm{s}}\} = \mathbf{E}_{\tilde{\mathbf{x}}_{\mathrm{s}}}\{\mathbf{x}(\mathbf{x} + \mathbf{n}_{\mathrm{s}})\} = \sigma_{\mathbf{x}}^{2}, \qquad (4.80)$$

$$\mathbf{E}_{\tilde{\mathbf{x}}_{s}}\left\{m_{s}(\tilde{\mathbf{x}}_{s})^{2}\right\} = \mathbf{E}_{\tilde{\mathbf{x}}_{s}}\left\{\mathbf{E}_{x}\left\{\mathbf{x}^{2} \mid \tilde{\mathbf{x}}_{s}\right\}\right\} - \mathbf{E}_{\tilde{\mathbf{x}}_{s}}\left\{\sigma_{s}^{2}(\tilde{\mathbf{x}}_{s})\right\} = \mathbf{E}_{x}\left\{\mathbf{x}^{2}\right\} - \sigma^{2} = \sigma_{x}^{2} - \sigma^{2} , \qquad (4.81)$$

it is

$$\begin{aligned} \sigma_{\rm u}^2 &= {\rm E}_{\tilde{\mathsf{x}}_{\rm s}} \Big\{ \sigma_{\rm s,u}^2(\tilde{\mathsf{x}}_{\rm s}) \Big\} = {\rm E}_{\tilde{\mathsf{x}}_{\rm s}} \Big\{ \sigma_{\rm s}^2(\tilde{\mathsf{x}}_{\rm s}) \Big\} + \left(\frac{\sigma^2}{\sigma^2 - \sigma_{\rm n_s}^2} \right)^2 {\rm E}_{\tilde{\mathsf{x}}_{\rm s}} \Big\{ (m_{\rm s}(\tilde{\mathsf{x}}_{\rm s}) - \tilde{\mathsf{x}}_{\rm s})^2 \Big\} \\ &= \sigma^2 + \left(\frac{\sigma^2}{\sigma^2 - \sigma_{\rm n_s}^2} \right)^2 \left(\sigma_{\rm x}^2 - \sigma^2 - 2\sigma_{\rm x}^2 + \sigma_{\rm x}^2 + \sigma_{\rm n_s}^2 \right) = \sigma^2 + \frac{(\sigma^2)^2 (\sigma_{\rm n_s}^2 - \sigma^2)}{(\sigma^2 - \sigma_{\rm n_s}^2)^2} \\ &= \frac{\sigma^2 (\sigma_{\rm n_s}^2 - \sigma^2) + (\sigma^2)^2}{\sigma_{\rm n_s}^2 - \sigma^2} = \frac{\sigma^2 \sigma_{\rm n_s}^2}{\sigma_{\rm n_s}^2 - \sigma^2} = \left(\frac{1}{\sigma^2} - \frac{1}{\sigma_{\rm n_s}^2} \right)^{-1}, \end{aligned}$$
(4.82)

which justifies the use of the respective update for the variance as well. This can also be derived graphically, the resulting squared length is shown in Fig. 4.7. Noteworthy, the averaging of variances performed in VAMP approximates $MSEs^3$, which are the parameters to use in bias compensation when deriving it from an estimation-theoretic perspective. Therefore, the procedure in VAMP is valid especially for high-dimensional signals x.

All in all, this shows that the extrinsic computation derived for VAMP is also justified from an estimation-theoretic perspective.

4.5 Turbo-Type Algorithms—Estimation-Theoretic Adaption

The previous sections have opened different possibilities of treating the MMSE estimates to prepare them as input for the other estimator. This insight will be used in the following to compose algorithms by replacing the strategies at respective steps. At first, an algorithm with individual variances is considered.

³In case of linear estimation conditional variance and MSE are the same, so $\sigma_{n_s}^2$ a.k.a. $\tilde{\sigma}_s^2$ is indeed an MSE.



Figure 4.8: Inverse of the unbiased variance $\tilde{\sigma}_{c,j}^2$ for the two strategies: optimization-based approach (4.83) and estimation-theoretic bias compensation (4.87), for N = 250, s = 5, $10 \log_{10}(1/\tilde{\sigma}_{s,j}^2) = 16 \text{ dB}$.

4.5.1 Individual Variances and Estimation-Theoretic Bias Compensation

When comparing the individual unbiasing procedure in Sec. 4.4.2.1 with the extrinsic calculation in the algorithm of Sec. 4.3.2, one can see certain differences. For example, although the update for the mean estimates shares the same form, the initial algorithm does not make use of the MSE itself, but instead works with conditional variances. In [FSG20] this treatment of variances was shown to be disadvantageous as negative variances may result frequently in the algorithm, when computing

$$\frac{1}{\tilde{\sigma}_{\mathrm{s},j}^2} = \frac{1}{\sigma_{\mathrm{s},j}^2} - \frac{1}{\tilde{\sigma}_{\mathrm{s},j}^2} , \quad j \in \{1, \dots, N\} .$$
(4.83)

Since negative variances have no useful meaning, this has to be prevented by, e.g., clipping the result to a positive interval [RSF19]. However, this affects the processing of the algorithm and ultimately deteriorates the performance.

A way to solve this problem is naturally given by the estimation-theoretic bias compensation view in Sec. 4.4.2.1. The corresponding algorithm is obtained by replacing the updates after the signal-constrained estimation with the derived Eqs. (4.76) and (4.77), which is stated below in Eqs. (4.86) and (4.87) in the notation of the algorithms again. The two computations of the unbiased variance (4.83) and (4.87) are compared in Fig. 4.8 for parameters that are used in the simulations in Chapter 6. The procedure with the subtraction leads to negative values in the gray area, while the estimation-theoretic unbiasing stays positive for all possible values of $\tilde{x}_{s,j}$ and $\tilde{\sigma}_{s,j}^2$.

The resulting algorithm was introduced in [FSG20] first, a similar version is called euIMS in [Spa19], here it is called *VAMPire*, for VAMP with individual reliabilities enhanced.

The main steps in the algorithm can be summarized as follows (in estimation-theoretic notation with mean and variances). Starting with a suitable initialization for the input variables of the channel-constrained estimation, e.g., $\tilde{\boldsymbol{x}}_{c} = \boldsymbol{0}$ and $\tilde{\boldsymbol{\Phi}}_{c} = s/N \cdot \boldsymbol{I}_{N}$, the unbiased

conditional mean and variance can directly be computed via

$$\tilde{\boldsymbol{x}}_{s} = \tilde{\boldsymbol{x}}_{c} + \boldsymbol{W} \tilde{\boldsymbol{\Phi}}_{c} \boldsymbol{A}^{\top} (\boldsymbol{A} \tilde{\boldsymbol{\Phi}}_{c} \boldsymbol{A}^{\top} + \sigma_{n}^{2} \boldsymbol{I}_{N})^{-1} (\boldsymbol{y} - \boldsymbol{A} \tilde{\boldsymbol{x}}_{c}) , \qquad (4.84)$$

$$\tilde{\sigma}_{s,j}^2 = \tilde{\sigma}_{c,j}^2 ([\boldsymbol{W}]_{jj} - 1) , \qquad (4.85)$$

where the diagonal entries of the compensation matrix read $[\mathbf{W}]_{jj} = 1/[\tilde{\boldsymbol{\Phi}}_{c} \mathbf{A}^{\top} (\mathbf{A} \tilde{\boldsymbol{\Phi}}_{c} \mathbf{A}^{\top} + \sigma_{n}^{2} \mathbf{I}_{N})^{-1} \mathbf{A}]_{jj}$. The results are used to compute signal-constrained conditional mean $m_{s,j}$ from (2.19) and conditional variance $\sigma_{s,j}^{2}$ from (2.20) for $j \in \{1, \ldots, N\}$.

For the subsequent bias compensation, the procedure derived in Sec. 4.4.2.1 is used, i.e., one needs to compute

$$\tilde{x}_{c,j} = \tilde{x}_{s,j} + \frac{\tilde{\sigma}_{s,j}^2}{\tilde{\sigma}_{s,j}^2 - \sigma_j^2} (m_{s,j} - \tilde{x}_{s,j}) , \qquad (4.86)$$

$$\tilde{\sigma}_{c,j}^2 = \sigma_{s,j}^2 + \left(\frac{\sigma_j^2}{\sigma_j^2 - \tilde{\sigma}_{s,j}^2}\right)^2 (m_{s,j} - \tilde{x}_{s,j})^2 .$$
(4.87)

The MSEs σ_j^2 are a direct function of $\tilde{\sigma}_{c,j}^2$ and can therefore be tabulated to reduce the complexity [FSG20]. In the next iteration, $\tilde{\boldsymbol{x}}_s = [\tilde{x}_{s,1}, \ldots, \tilde{x}_{s,N}]^\top$ and $\tilde{\boldsymbol{\Phi}}_c = \operatorname{diag}(\tilde{\sigma}_{c,j}^2)$ are used as input variables for the channel-constrained estimation and the computations are repeated until a suitable stopping criterion is met. The Algorithm is also stated in Appendix E.7.

It has to be noted that the computational complexity of this algorithm ranges in⁴ $\mathcal{O}(M \cdot N^2)$, whereas VAMP is in $\mathcal{O}(MN)$, for the derivations see [Spa19]. The reason for the increased complexity of VAMPire is that in the channel-constrained estimation, a $M \times M$ matrix has to be inverted. For the VAMP case, the inverse can be simplified to the inversion of a diagonal matrix, i.e., scalar inversions are performed, by utilizing the singular value decomposition (SVD) of the sensing matrix [RSF19, Spa19]. Switching from scaled identity matrix $\tilde{\boldsymbol{\Phi}}_{c} = \tilde{\sigma}_{c}^{2} \boldsymbol{I}_{N}$ to a diagonal version $\tilde{\boldsymbol{\Phi}}_{c} = \operatorname{diag}(\tilde{\sigma}_{c,j}^{2})$, i.e., individual variances, the decomposition cannot yield a diagonal inverse, so the complexity reduction is not possible. This is of course a huge drawback and means that the gained insight into the individual reliabilities of the linear estimate is relative costly.

4.5.2 VAMP with Individually Unbiased LMMSE Estimate

It has been shown in Sec. 4.4.1.1 that the usual version of the VAMP algorithm inherently performs average bias compensation after the linear estimation. The individual bias compensation strategy has so far only been used, when individual variances were present. With individual variances in the linear estimator, this leads to a huge increase in complexity, since the computation of the inverse can not be facilitates by the SVD trick, as suggested in [RSF19]. However, there is no reason not to use individual unbiasing in a setting with average variances. As will be shown in numerical simulations in Chapter 6, individual unbiasing is particularly useful if the columns of the sensing matrix carry different energies, i.e., have different (squared) ℓ_2 norms. In the following, we therefore compose algorithms that combine the use of an average

⁴By $f(x) = \mathcal{O}(g(x))$ the Landau notation is denoted, which bounds the growth of f(x) for $x \to \infty$ by the function g(x) in the sense that $|f(x)| \leq \text{const.} \cdot g(x)$ for $x \geq x_0$ from some point $x_0 \in \mathbb{R}$ on [Lan09].

variance in the linear estimation with individual bias compensation after the linear estimation. The respective scheme has been published in [SF22].

The channel-constrained estimation is for this case given as

$$\tilde{\boldsymbol{x}}_{s} = \tilde{\boldsymbol{x}}_{c} + \boldsymbol{W}\boldsymbol{A}^{\top}(\boldsymbol{A}\boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}}\boldsymbol{I}_{M})^{-1}(\boldsymbol{y} - \boldsymbol{A}\tilde{\boldsymbol{x}}_{c}) , \qquad (4.88)$$

where $[\mathbf{W}]_{jj} = 1/[\mathbf{A}^{\top}(\mathbf{A}\mathbf{A}^{\top} + \frac{\sigma_n^2}{\tilde{\sigma}_c^2}\mathbf{I}_M)^{-1}\mathbf{A}]_{jj}$. As reliability for the signal-constrained estimation in (2.19) and (2.20), one can either keep the individual unbiased variances

$$\tilde{\sigma}_{s,j}^2 = \sigma_x^2 \left([\boldsymbol{W}]_{jj} - 1 \right) ,$$
 (4.89)

or average them to obtain

$$\tilde{\sigma}_{s}^{2} = \sigma_{x}^{2} \left(M_{A}([\boldsymbol{W}]_{jj}) - 1 \right) .$$
 (4.90)

This leads to two different variants of the VAMP algorithm, which is denoted by *VAMPii* for individual unbiasing and individual (unbiased) variances and *VAMPia* for individual unbiasing and average (unbiased) variances.

For the unbiasing after the signal-constrained estimation, the procedure of average unbiasing is utilized, since the resulting input variance for the channel-constrained estimation needs to be an average one. Therefore, it is in both cases

$$\tilde{\boldsymbol{x}}_{c} = \tilde{\sigma}_{c}^{2} (\boldsymbol{m}_{s} / \sigma_{s}^{2} - \tilde{\boldsymbol{x}}_{s} / \tilde{\sigma}_{s}^{2}) , \quad \tilde{\sigma}_{c}^{2} = (1 / \sigma_{s}^{2} - 1 / \tilde{\sigma}_{s}^{2})^{-1} ,$$
(4.91)

where $\sigma_{\rm s}^2 = M_{\rm A}(\sigma_{{\rm s},j}^2)$ has to be averaged for both algorithms and $\tilde{\sigma}_{\rm s}^2 = M_{\rm A}(\tilde{\sigma}_{{\rm s},j}^2)$ additionally for VAMPii.

The complete algorithms are stated in Appendix E.8.

On the next double page, the Turbo algorithms are summarized in the form of block diagrams.

4.6 Block Diagrams of the Turbo Algorithms

In the following, the Turbo-type algorithms are summarized in the form of block diagrams. Here, the block for the signal-constrained estimation accepts the mean \tilde{x}_s and corresponding variance(s) as input parameters and returns the conditional mean m_s . The depiction shows exemplary the signal-constrained estimation for the DT prior over input values $\tilde{x}_{s,j}$ for a given variance.

4.6.1 VAMP



4.6.2 VAMPind







4.6.4 VAMPia





65

5. Sequential Inference

In the last two chapters, we have seen the two extreme cases of working with the factor graph for CS. Chapter 3 factorized the conditional pdf (2.7) as much as possible to yield the most simplest estimators, which led to a relative high complexity because of the huge number ($\mathcal{O}(M \cdot N)$ due to the number of connecting edges) of messages that needed to be broadcasted.

In Chapter 4, the opposite approach was considered by factorizing to a minimum number of components, i.e., separating only what really needed to be apart. The result were two major components, between which the recovery algorithms iterated to achieve a suitable solution.

The chapter at hand, describes cases that lie in between these two extreme approaches. The idea is to factorize the posterior $f_{x|y}(x)$ more than in the Turbo case (Chapter 4) but not as much as in the MP case (Chapter 3). In contrast to the Turbo procedure, this leads to the possibility of (different) schedules, as in the MP case, but without the corresponding high computational complexity in terms of messages to be propagated. The processing will result in a *sequential* version of the (parallel) Turbo case.

Subsequently two different approaches will be considered that follow the remaining factorization possibilities of Fig. 2.3, starting with separated variables nodes, but combined channel factors, i.e., Fig. 2.3(c) and finishing off with the final part Fig. 2.3(b), which depicts separated channel factor nodes and a combined variable node.

5.1 Sequential Processing of Variables

When considering the processing in the Turbo case, one may realize that the computation of the signal-constraint estimates (2.19) and (2.20), although being processed in parallel, not necessarily need to be computed together, since the variables are considered i.i.d., cf. (2.18). A natural idea for sequentializing the Turbo process is, therefore, to consider only one component x_j of the signal vector in the signal-constrained estimation to find the impact of the single variable on the channel-constrained estimation. In the next step, another component is considered, eventually leading to a sequential processing of the variables x_1, \ldots, x_N .

Sequential processing of variables can for the CS literature be found, e.g., in [MKTZ15],

where the AMP algorithm is sequentialized. In contrast, in this chapter sequential versions of VAMPind, i.e., VAMP with individual variances, and VAMPire are considered. Other examples for sequential (variable) processing in CS are [SN08], where the focus lies on the design of the sensing matrix A—in contrast, we consider A fixed and focus on the algorithmic part— and [BMPP20], which considers a Bernoulli-Gaussian prior whose sparsity parameter is learned throughout the algorithm. A similar algorithm to the one in [BMPP20] is also described in [OW05, Appendix D] for a general inference problem with binary prior. The respective algorithm for CS with arbitrary prior is explained in Sec. 5.1.4 below. The algorithm was improved in [SF21] by an estimation-theoretic consideration of the bias compensation; here, this version of the algorithm is discussed in Sec. 5.1.5.

Noteworthy, the mentioned algorithms in [OW05, BMPP20, SN08] are all based on the *expectation propagation* (EP) [Min01] framework, which considers (in its initial version) sequential updates in a general inference problem but is often interpreted in a parallel schedule, last but not least also for VAMP [RSF19].

Recent work on sequential processing is also given in [ÇLO22], which analyzes a sequential version of VAMP with random schedule for a Gaussian model with latent parameter.

Another study on sequential CS, which is further away than the previously mentioned references, is [LCJ15], where dynamically changing signals with spatial and temporal correlation are processed sequentially in a wireless sensor network (WSN).

The derivation given here, follows the optimization-based approach that was also considered in the previous chapters.

5.1.1 Structure of the Factor Graph

The idea for the sequential processing of the variables x_1, \ldots, x_N that is pursued in this part of the chapter is motivated by the fact that the components of the signal vector **x** are i.i.d.. Consequently, the signal factors $f_x(x_j)$ are considered individually in the factorization, which yields individual processing of the signal-constrained estimates. The channel-constrained part, on the other hand, is kept together, thereby enabling the possibility to propagate effects of a single signal-constrained estimation to the other signal components in one step. The factorization that is considered here is therefore given by

$$f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}) = \text{const.} \cdot f_{\mathbf{y}|\mathbf{x}}(\mathbf{x}) \cdot \prod_{j=1}^{N} f_{\mathbf{x}}(x_j) .$$
(5.1)

The corresponding factor graph is the one shown in Fig. 2.3(c), which is plotted in Fig. 5.1 for convenience again.

5.1.1.1 Factorization of the Substitute Distribution

Following the factorization in (5.1), $q_c(x)$ approximates $f_{\mathbf{y}|\mathbf{x}}(x)$ as in (4.4) and $q_{s,j}(x_j)$ approximates $f_{\mathbf{x}}(x_j)$ ($j \in \{1, \ldots, N\}$) as in (3.3), where it is assumed that the deviation in the



Figure 5.1: Factor graph for sequential processing of variable nodes.

approximation for the channel-constrained factor splits according to

$$c(\boldsymbol{x}) = \prod_{j=1}^{N} c_j(x_j) .$$
(5.2)

The individual marginals $q_{v,j}(x_j)$ are then given by

$$\mathbf{q}_{\mathbf{v},j}(x_j) = \text{const.} \cdot s_j(x_j) \cdot c_j(x_j) .$$
(5.3)

The individual representation is used (instead of an entire representative $q_v(x)$) because the variable nodes are considered separately in the factor graph. The resulting factorization of the substitute distribution reads

...

$$\mathbf{q}(\boldsymbol{x}) = \frac{\mathbf{q}_{c}(\boldsymbol{x}) \cdot \prod_{j=1}^{N} \mathbf{q}_{s,j}(x_{j})}{\prod_{j=1}^{N} \mathbf{q}_{v,j}(x_{j})} .$$
(5.4)

5.1.1.2 The Moment-Matching Constraint

The consistency constraints, which need to be defined per edge in the factor graph, are in the given case set up for all variables x_1, \ldots, x_N . Since each variable node is connected to a signal factor and the channel factor, two constraints need to be established per variable node. As before, exponential families are used, especially in the form of Gaussian distributions, because eventually the (first-order) moment of the MMSE estimate (2.9) is desired. However, since the signal factors are considered individually per variable, it is necessary to specify the reliability of the estimation individually. Therefore, the sufficient statistics of the exponential family is determined by

$$\boldsymbol{g}(\boldsymbol{x}) = [\boldsymbol{g}_{\boldsymbol{\lambda}}(\boldsymbol{x})^{\top}, \, \boldsymbol{g}_{\boldsymbol{\Lambda}}(\boldsymbol{x})^{\top}]^{\top}$$
(5.5)

with

$$\boldsymbol{g}_{\boldsymbol{\lambda}}(\boldsymbol{x}) = [x_1, \ldots, x_N]^{\top}, \quad \boldsymbol{g}_{\boldsymbol{\Lambda}}(\boldsymbol{x}) = -\frac{1}{2}[x_1^2, \ldots, x_N^2]^{\top},$$
 (5.6)

cf. Ex. D.2 in Appendix D. For convenience, the notation $g(x_j) = [g_{\lambda_j}(x_j), g_{\Lambda_j}(x_j)]^\top = [x_j, -x_j^2/2]^\top$ may be used. With these specifications, the constraints on the moments read

$$\mathbf{E}_{\mathsf{x}_{j} \sim \mathsf{q}_{\mathsf{x},j}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} \stackrel{!}{=} \mathbf{E}_{\mathsf{x}_{j} \sim \mathsf{q}_{\mathsf{s},j}} \{ \boldsymbol{g}(\mathsf{x}_{j}) \} , \qquad j \in \{1, \ldots, N\} , \qquad (5.7)$$

$$\mathbf{E}_{\mathsf{x}_{j}\sim\mathsf{q}_{\mathsf{v},j}}\{\boldsymbol{g}(\mathsf{x}_{j})\}\stackrel{!}{=} \mathbf{E}_{\mathsf{x}_{j}\sim\mathsf{q}_{c}}\{\boldsymbol{g}(\mathsf{x}_{j})\}, \qquad j\in\{1,\ldots,N\}.$$
(5.8)

5.1.2 Constrained Optimization

The optimization-based approach that is used to derive the parameter connections for the algorithm takes on the form of a constrained optimization problem, where the task of minimizing the Kullback–Leibler divergence is combined with the moment constraints from above and normalization constraints for the representatives $q_c(x)$, $q_{s,j}(x_j)$, and $q_{v,j}(x_j)$ ($j \in \{1, ..., N\}$). The entire problem reads

$$\begin{array}{ll} \text{minimize} & \mathrm{D}_{\mathrm{KL}} \Big(\mathbf{q}(\boldsymbol{x}) \mid| \, \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) \Big) \\ \text{subject to} & \mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathbf{v},j}} \{ \boldsymbol{g}(\mathbf{x}_{j}) \} - \mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathbf{c}}} \{ \boldsymbol{g}(\mathbf{x}_{j}) \} = \boldsymbol{0} \,, \qquad j \in \{1, \ldots, N\} \,, \\ & \mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathbf{v},j}} \{ \boldsymbol{g}(\mathbf{x}_{j}) \} - \mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathbf{s},j}} \{ \boldsymbol{g}(\mathbf{x}_{j}) \} = \boldsymbol{0} \,, \qquad j \in \{1, \ldots, N\} \,, \\ & \int \mathbf{q}_{\mathbf{v},j}(x_{j}) \, \mathrm{d}x_{j} = 1 \,, \quad \int \mathbf{q}_{\mathbf{s},j}(x_{j}) \, \mathrm{d}x_{j} = 1 \,, \quad j \in \{1, \ldots, N\} \,, \\ & \int \mathbf{q}_{\mathbf{c}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = 1 \,. \end{array} \right.$$

5.1.2.1 Cost Function for Compressed Sensing

For the sake of completeness, the form of the Kullback–Leibler divergence for the given case is stated, which is somewhat a mixture of the cases in the chapters before

$$D_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) \mid\mid \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})) = \sum_{j=1}^{N} \int \mathbf{q}_{\mathrm{s},j}(x_{j}) \log \mathbf{q}_{\mathrm{s},j}(x_{j}) \,\mathrm{d}x_{j} + \int \mathbf{q}_{\mathrm{c}}(\boldsymbol{x}) \log \mathbf{q}_{\mathrm{c}}(\boldsymbol{x}) \,\mathrm{d}\boldsymbol{x}$$
$$- \sum_{j=1}^{N} \int \mathbf{q}_{\mathrm{s},j}(x_{j}) \log \mathbf{f}_{\mathbf{x}}(x_{j}) \,\mathrm{d}x_{j} - \int \mathbf{q}_{\mathrm{c}}(\boldsymbol{x}) \log \mathbf{f}_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) \,\mathrm{d}\boldsymbol{x}$$
$$- \sum_{j=1}^{N} \int \mathbf{q}_{\mathrm{v},j}(x_{j}) \log \mathbf{q}_{\mathrm{v},j}(x_{j}) \,\mathrm{d}x_{j} \,. \tag{5.10}$$

5.1.2.2 Stationary Points of the Lagrangian

To see the parameter connections that are necessary for the recovery algorithm, a Lagrangian is set up and the stationary points are derived, because the stationary points indicate the form of the solution to the problem. For the Lagrangian, multipliers $\boldsymbol{\nu}_{c,j}$ and $\boldsymbol{\nu}_{s,j}$ $(j \in \{1, ..., N\})$ are defined for the moment-matching constraint, as well as ν_c , ν_j , and $\nu_{s,j}$ $(j \in \{1, ..., N\})$ for the normalization constraint. The resulting Lagrangian reads

$$\mathcal{L}(\mathbf{q}(\boldsymbol{x})) = D_{\mathrm{KL}}\left(\mathbf{q}(\boldsymbol{x}) \mid\mid \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})\right) + \sum_{j=1}^{N} \boldsymbol{\nu}_{\mathrm{c},j}^{\top} \left(\mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathrm{v},j}} \{\boldsymbol{g}(\mathbf{x}_{j})\} - \mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathrm{c}}} \{\boldsymbol{g}(\mathbf{x}_{j})\}\right) + \sum_{j=1}^{N} \boldsymbol{\nu}_{\mathrm{s},j}^{\top} \left(\mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathrm{v},j}} \{\boldsymbol{g}(\mathbf{x}_{j})\} - \mathrm{E}_{\mathbf{x}_{j} \sim \mathbf{q}_{\mathrm{s},j}} \{\boldsymbol{g}(\mathbf{x}_{j})\}\right) + \sum_{j=1}^{N} \left(\nu_{j} \left(\int \mathbf{q}_{\mathrm{v},j}(x_{j}) \, \mathrm{d}x_{j} - 1 \right) + \nu_{\mathrm{s},j} \left(\int \mathbf{q}_{\mathrm{s},j}(x_{j}) \, \mathrm{d}x_{j} - 1 \right) \right) + \nu_{\mathrm{c}} \left(\int \mathbf{q}_{\mathrm{c}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 \right) .$$
(5.11)

Similar to the previous chapters, the stationary points are obtained by differentiating w.r.t. the representatives $q_c(x)$, $q_{v,j}(x_j)$, and $q_{s,j}(x_j)$ ($j \in \{1, ..., N\}$), respectively, setting the results to zero and solve for the respective representative, i.e.,

$$\frac{\partial}{\partial \mathbf{q}_{\mathbf{v},j}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = -\log \mathbf{q}_{\mathbf{v},j}(x_j) + (\boldsymbol{\nu}_{\mathbf{c},j} + \boldsymbol{\nu}_{\mathbf{s},j})^{\top} \boldsymbol{g}(\mathbf{x}_j) + \text{const.} \stackrel{!}{=} 0, \qquad (5.12)$$

$$\frac{\partial}{\partial \mathbf{q}_{\mathrm{s},j}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = \log \mathbf{q}_{\mathrm{s},j}(x_j) - \log \mathbf{f}_{\mathrm{x}}(x_j) - \boldsymbol{\nu}_{\mathrm{s},j}^{\top} \boldsymbol{g}(\mathbf{x}_j) + \text{const.} \stackrel{!}{=} 0, \qquad (5.13)$$

$$\frac{\partial}{\partial \mathbf{q}_{c}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = \log \mathbf{q}_{c}(\boldsymbol{x}) - \log f_{\mathbf{y}|\mathbf{x}}(\boldsymbol{x}) - \sum_{j=1}^{N} \boldsymbol{\nu}_{c,j}^{\top} \boldsymbol{g}(\mathbf{x}_{j}) + \text{const.} \stackrel{!}{=} 0.$$
(5.14)

Finally, one obtains

$$\mathbf{q}_{\mathbf{v},j}(x_j) = \frac{1}{Z_{\mathbf{v},j}} \exp\left((\boldsymbol{\nu}_{\mathrm{s},j} + \boldsymbol{\nu}_{\mathrm{c},j})^\top \boldsymbol{g}(x_j)\right) , \qquad (5.15)$$

$$q_{c}(\boldsymbol{x}) = \frac{1}{Z_{c}} f_{\boldsymbol{y}|\boldsymbol{x}}(\boldsymbol{x}) \exp\left(\sum_{j=1}^{N} \boldsymbol{\nu}_{c,j}^{\top} \boldsymbol{g}(x_{j})\right) , \qquad (5.16)$$

$$\mathbf{q}_{\mathbf{s},j}(x_j) = \frac{1}{Z_{\mathbf{s},j}} \mathbf{f}_{\mathbf{x}}(x_j) \exp\left(\boldsymbol{\nu}_{\mathbf{s},j}^{\top} \boldsymbol{g}(x_j)\right) , \qquad (5.17)$$

with normalizing partition functions $Z_{v,j}$, Z_c , and $Z_{s,j}$.

The basic result of the derivations is that again the connection of the natural parameters in the variable node is given by

$$\boldsymbol{\theta}_{j} \stackrel{\text{\tiny def}}{=} \boldsymbol{\nu}_{\mathrm{c},j} + \boldsymbol{\nu}_{\mathrm{s},j} , \qquad (5.18)$$

which means that in the algorithm, *post-processing* is done as extrinsic calculation, i.e., the inputs to the estimations are obtained by subtractions of previously computed estimates. The procedure will be discussed in detail below.

5.1.2.3 Connection to Expectation Propagation

The basic idea behind expectation propagation (EP) [Min01] is to refine an approximation q(x) by replacing a representing factor of the approximation by the corresponding exact term

and projecting the refinement back into the space of distributions that is considered. This is iteratively conducted for all the parts / factors, that are available in the structure of the considered problem.

In terms of exponential families, it is directly visible how this connects to the derivation here. Starting with a representative $q(\boldsymbol{x}) = \frac{1}{Z(\theta)} \exp(\theta^{\top} \boldsymbol{g}(\boldsymbol{x}))$, which encodes the structure of the problem by a respective summation of its natural parameters, e.g., $\theta_j = \boldsymbol{\nu}_{c,j} + \boldsymbol{\nu}_{s,j}$ in the case above (then the entire natural parameter might be $\boldsymbol{\theta} = [\boldsymbol{\theta}_1^{\top}, \ldots, \boldsymbol{\theta}_N^{\top}]^{\top}$). The representing factor is removed by subtracting its corresponding natural parameter from the natural parameter representing the entire distribution. The resulting parameter is called the *cavity parameter*, because it excludes the influence of the respective factor. Inserting the exact factor—in the given case any of $f_x(x_j)$ ($j \in \{1, \ldots, N\}$)—yields something of the form of the stationary points $q_{s,j}(x_j) = \frac{1}{Z_{s,j}} f_x(x_j) \exp(\boldsymbol{\nu}_{s,j}^{\top} \boldsymbol{g}(x_j))$ of the Lagrangian, i.e., the cavity parameter is here $\boldsymbol{\nu}_{s,j} = \boldsymbol{\theta}_j - \boldsymbol{\nu}_{c,j}$. The projection is then performed as minimization of the Kullback–Leibler divergence such that the result is a member of an exponential family, i.e., of the form, e.g., [Min05]

$$\mathbf{q}_{\mathbf{v},j}(x_j) = \min_{\mathbf{q}(x_j)} \mathcal{D}_{\mathrm{KL}}(\mathbf{q}_{\mathrm{s},j}(x_j) \mid\mid \mathbf{q}(x_j)) .$$
(5.19)

As indicated, the result of the projection is of the form of $q_{v,j}(x_j)$ from (5.15). The restriction to members of an exponential family causes a matching of moments, as shown in the insert below. Noteworthy, the minimization of the Kullback–Leibler divergence considered here has the position of substitute distribution and given distribution swapped compared to the instances used so far and it requires arguments of a different form.

The matching of moments can be interpreted as in Chapter 3 and be realized by computing moments and the resulting natural parameters of a corresponding member of an exponential family, which yields an updated θ_j . The remaining task is to update the so-called *site parameters*, i.e., the complementary part of the cavity parameters, here by $\nu_{c,j} = \theta_j - \nu_{s,j}$.

Another interesting aspect is that the computation of extrinsic is directly realized in this procedure and can intuitively be explained by the fact that one has to make room in the approximation (by subtracting the corresponding natural parameter representation) for the exact part to be inserted.

Insert .

We show that, when restricting the result of the minimization (5.19) to an exponential family, a matching of moments is obtained. Assume that $q(x_j) = \frac{1}{Z(\theta)} \exp(\theta^{\top} g(x_j))$. Under this assumption, the Kullback–Leibler divergence can be written as

$$D_{\mathrm{KL}}(\mathbf{q}_{\mathrm{s},j}(x_j) \mid\mid \mathbf{q}(x_j)) = \int \mathbf{q}_{\mathrm{s},j}(x_j) \log \mathbf{q}_{\mathrm{s},j}(x_j) \,\mathrm{d}x_j - \boldsymbol{\theta}^\top \int \mathbf{q}_{\mathrm{s},j}(x_j) \boldsymbol{g}(x_j) \,\mathrm{d}x_j + \log Z(\boldsymbol{\theta}) \,.$$
(5.20)

Then differentiation of $D_{KL}(\mathbf{q}_{s,j}(x_j)||\mathbf{q}(x_j))$ w.r.t. $\boldsymbol{\theta}$ yields

$$\frac{\partial}{\partial \boldsymbol{\theta}} \mathrm{D}_{\mathrm{KL}}(\mathbf{q}_{\mathrm{s},j}(x_j) \mid\mid \mathbf{q}(x_j)) = \frac{\partial}{\partial \boldsymbol{\theta}} \left(\log Z(\boldsymbol{\theta}) - \boldsymbol{\theta}^\top \mathrm{E}_{\mathbf{x}_j \sim \mathbf{q}_{\mathrm{s},j}} \{ \boldsymbol{g}(\mathbf{x}_j) \} \right) \\ = \mathrm{E}_{\mathbf{x}_j \sim \mathbf{q}} \{ \boldsymbol{g}(\mathbf{x}_j) \} - \mathrm{E}_{\mathbf{x}_j \sim \mathbf{q}_{\mathrm{s},j}} \{ \boldsymbol{g}(\mathbf{x}_j) \} ,$$
(5.21)

where (D.18) was used for the differentiation of $\log Z(\boldsymbol{\theta})$. Subsequently setting this to zero, yields the matching of moments $E_{x_j \sim q} \{ \boldsymbol{g}(x_j) \} = E_{x_j \sim q_{s,j}} \{ \boldsymbol{g}(x_j) \}$, which means that any extremum must fulfill this relation.

_ End Insert

5.1.3 Expectations and Rank-One Update

The approach considered here has been motivated by the idea of refining the LMMSE (channelconstrained) estimate sequentially based on the impact of single estimations on the signalconstrained side of the problem. For the signal-constrained estimations, the individual calculations are needed for which it is known from Chapter 3 that with $\boldsymbol{\nu}_{s,j} = [\tilde{x}_{s,j}/\tilde{\sigma}_{s,j}^2, 1/\tilde{\sigma}_{s,j}^2]^{\top}$, it is

$$\mathbf{E}_{\mathsf{x}_{j} \sim \mathsf{q}_{\mathsf{s},j}} \left\{ g_{\lambda_{j}}(\mathsf{x}_{j}) \right\} = \mathbf{E}_{\mathsf{x}_{j}} \left\{ \mathsf{x}_{j} \mid \tilde{x}_{\mathsf{s},j} \right\} = m_{\mathsf{s},j} , \qquad (5.22)$$

$$E_{x_{j} \sim q_{s,j}} \left\{ g_{A_{j}}(x_{j}) \right\} = -\frac{1}{2} E_{x_{j}} \left\{ x_{j}^{2} \mid \tilde{x}_{s,j} \right\} = -\frac{1}{2} \left(\sigma_{s,j}^{2} + m_{s,j}^{2} \right) .$$
(5.23)

On the side of the channel-constrained estimation it is useful to stack the input parameters such that

$$\sum_{j=1}^{N} \boldsymbol{\nu}_{\mathrm{c},j}^{\top} \boldsymbol{g}(x_j) = \boldsymbol{\nu}_{\mathrm{c}}^{\top} \boldsymbol{g}(\boldsymbol{x}) .$$
(5.24)

With the identification

$$\boldsymbol{\nu}_{\rm c} = [m_{{\rm x},1}/\sigma_{{\rm x},1}^2, \, \dots, \, m_{{\rm x},N}/\sigma_{{\rm x},N}^2, \, 1/\sigma_{{\rm x},1}^2, \, \dots, \, 1/\sigma_{{\rm x},N}^2]^\top \,, \tag{5.25}$$

it can be referred to Chapter 4, where for individual variances it has been found that

$$\mathbf{E}_{\mathbf{x}\sim q_{c}}\{\boldsymbol{g}_{\boldsymbol{\lambda}}(\mathbf{x})\} = \boldsymbol{m}_{c} , \qquad (5.26)$$

respectively,

$$E_{\mathbf{x} \sim \mathbf{q}_{c}} \{ \boldsymbol{g}_{\boldsymbol{\Lambda}}(\mathbf{x}) \} = -\frac{1}{2} \operatorname{diag} \left(\boldsymbol{\varPhi}_{c} + \boldsymbol{m}_{c} \boldsymbol{m}_{c}^{\top} \right).$$
(5.27)

For the moment-matching, the respective jth components are required, i.e.,

$$\mathbf{E}_{\mathbf{x}_{j}\sim\mathbf{q}_{c}}\left\{g_{\lambda_{j}}(\mathbf{x}_{j})\right\} = m_{c,j} , \qquad (5.28)$$

and

$$E_{x_{j} \sim q_{c}} \left\{ g_{A_{j}}(x_{j}) \right\} = -\frac{1}{2} (\sigma_{c,j}^{2} + m_{c,j}^{2}) .$$
(5.29)

In the following, the effect of the changes in the channel-constrained estimates induced by a single signal-constrained estimation at position j is considered. To that end, assume that $u_{\rm c}$ and the corresponding estimates $m_{\rm c}$ and $\Phi_{\rm c}$ are fixed. In order to simplify the notation, the following definitions are set up

$$\boldsymbol{\lambda}_{\mathrm{c}} \stackrel{\text{\tiny def}}{=} [\lambda_{\mathrm{c},1}, \, \dots, \, \lambda_{\mathrm{c},N}]^{\top} = [m_{\mathrm{x},1}/\sigma_{\mathrm{x},1}^2, \, \dots, \, m_{\mathrm{x},N}/\sigma_{\mathrm{x},N}^2]^{\top} , \qquad (5.30)$$

$$\mathbf{\Lambda}_{\mathrm{c}} \stackrel{\text{\tiny def}}{=} [\Lambda_{\mathrm{c},1}, \ldots, \Lambda_{\mathrm{c},N}]^{\top} = [1/\sigma_{\mathrm{x},1}^2, \ldots, 1/\sigma_{\mathrm{x},N}^2]^{\top}, \qquad (5.31)$$

then the combined input parameter reads $\boldsymbol{\nu}_{c} = [\boldsymbol{\lambda}_{c}^{\top}, \boldsymbol{\Lambda}_{c}^{\top}]^{\top}$. With these input parameters, the channel-constrained estimation can be written as follows by reformulating (2.10) and making the identifications $\boldsymbol{\varPhi}_{x}^{-1} = \operatorname{diag}(\Lambda_{c,j})$ and $\boldsymbol{\lambda}_{c} = \boldsymbol{\varPhi}_{x}^{-1}\boldsymbol{m}_{x} = \operatorname{diag}(\Lambda_{c,j})\boldsymbol{m}_{x}$

$$\begin{split} \boldsymbol{\Phi}_{c} &= \sigma_{n}^{2} (\boldsymbol{A}^{\top} \boldsymbol{A} + \sigma_{n}^{2} \operatorname{diag}(\boldsymbol{\Lambda}_{c,j}))^{-1} \\ &= (\frac{1}{\sigma_{n}^{2}} \boldsymbol{A}^{\top} \boldsymbol{A} + \operatorname{diag}(\boldsymbol{\Lambda}_{c,j}))^{-1} , \end{split}$$
(5.32)
$$\boldsymbol{m}_{c} &= \boldsymbol{m}_{x} + (\boldsymbol{A}^{\top} \boldsymbol{A} + \sigma_{n}^{2} \operatorname{diag}(\boldsymbol{\Lambda}_{c,j}))^{-1} \boldsymbol{A}^{\top} (\boldsymbol{y} - \boldsymbol{A} \boldsymbol{m}_{x}) \\ &= (\boldsymbol{A}^{\top} \boldsymbol{A} + \sigma_{n}^{2} \operatorname{diag}(\boldsymbol{\Lambda}_{c,j}))^{-1} ((\boldsymbol{A}^{\top} \boldsymbol{A} + \sigma_{n}^{2} \operatorname{diag}(\boldsymbol{\Lambda}_{c,j})) \boldsymbol{m}_{x} + \boldsymbol{A}^{\top} \boldsymbol{y} - \boldsymbol{A}^{\top} \boldsymbol{A} \boldsymbol{m}_{x}) \\ &= (\boldsymbol{A}^{\top} \boldsymbol{A} + \sigma_{n}^{2} \operatorname{diag}(\boldsymbol{\Lambda}_{c,j}))^{-1} (\boldsymbol{A}^{\top} \boldsymbol{y} + \sigma_{n}^{2} \operatorname{diag}(\boldsymbol{\Lambda}_{c,j}) \boldsymbol{m}_{x}) \\ &= (\frac{1}{\sigma_{n}^{2}} \boldsymbol{A}^{\top} \boldsymbol{A} + \operatorname{diag}(\boldsymbol{\Lambda}_{c,j}))^{-1} (\frac{1}{\sigma_{n}^{2}} \boldsymbol{A}^{\top} \boldsymbol{y} + \boldsymbol{\lambda}_{c}) \\ &= \boldsymbol{\Phi}_{c} (\frac{1}{\sigma_{n}^{2}} \boldsymbol{A}^{\top} \boldsymbol{y} + \boldsymbol{\lambda}_{c}) . \end{aligned}$$
(5.33)

The change in position j affects only the input parameter $\boldsymbol{\nu}_{c,j} = [\lambda_{c,j}, \Lambda_{c,j}]^{\top}$. The changed (new) parameters at position j are denoted by $\lambda_{c,j}^{\circ}$ and $\Lambda_{c,j}^{\circ}$, respectively. With jth unit vector \boldsymbol{e}_{j} , which contains a 1 at the jth position and zeros elsewhere, and $\Delta \lambda_{c,j} = \lambda_{c,j}^{\circ} - \lambda_{c,j}$, as well as $\Delta \Lambda_{c,j} = \Lambda_{c,j}^{\circ} - \Lambda_{c,j}$ one can write the resulting input parameter vectors as (cf. [OW05])

$$\boldsymbol{\lambda}_{\rm c}^{\circ} = [\lambda_{{\rm c},1}, \ldots, \lambda_{{\rm c},j-1}, \lambda_{{\rm c},j}^{\circ}, \lambda_{{\rm c},j+1}, \ldots, \lambda_{{\rm c},N}]^{\top} = \boldsymbol{\lambda}_{\rm c} + \Delta \lambda_{{\rm c},j} \boldsymbol{e}_{j} , \qquad (5.34)$$

$$\boldsymbol{\Lambda}_{\mathrm{c}}^{\circ} = [\Lambda_{\mathrm{c},1}, \ldots, \Lambda_{\mathrm{c},j-1}, \Lambda_{\mathrm{c},j}^{\circ}, \Lambda_{\mathrm{c},j+1}, \ldots, \Lambda_{\mathrm{c},N}]^{\top} = \boldsymbol{\Lambda}_{\mathrm{c}} + \Delta \Lambda_{\mathrm{c},j} \boldsymbol{e}_{j} .$$
(5.35)

The last line can be written in matrix form as

$$\operatorname{diag}(\boldsymbol{\Lambda}_{\mathrm{c}}^{\circ}) = \operatorname{diag}(\boldsymbol{\Lambda}_{\mathrm{c}}) + \Delta \boldsymbol{\Lambda}_{\mathrm{c},j} \boldsymbol{e}_{j} \boldsymbol{e}_{j}^{\top} .$$
(5.36)

The additive form of the change in the vector helps us to formulate the change in the resulting channel-constrained as follows. Using the Sherman-Morrison-Woodbury identity [GV96], the new covariance matrix $\boldsymbol{\Phi}_{c}^{\circ}$ can be written as a rank-one update of the previous one $\boldsymbol{\Phi}_{c}$ via

$$\boldsymbol{\Phi}_{c}^{\circ} = \left(\frac{1}{\sigma_{n}^{2}}\boldsymbol{A}^{\top}\boldsymbol{A} + \operatorname{diag}(\boldsymbol{\Lambda}_{c}^{\circ})\right)^{-1} = \left(\boldsymbol{\Phi}_{c}^{-1} + \Delta\boldsymbol{\Lambda}_{c,j}\boldsymbol{e}_{j}\boldsymbol{e}_{j}^{\top}\right)^{-1}$$
$$= \boldsymbol{\Phi}_{c} - \frac{\Delta\boldsymbol{\Lambda}_{c,j}}{1 + \Delta\boldsymbol{\Lambda}_{c,j}\boldsymbol{e}_{j}^{\top}\boldsymbol{\Phi}_{c}\boldsymbol{e}_{j}}\boldsymbol{\Phi}_{c}\boldsymbol{e}_{j}\boldsymbol{e}_{j}^{\top}\boldsymbol{\Phi}_{c} .$$
(5.37)

Since the unit vector picks the *j*th column, i.e., it is $\boldsymbol{\Phi}_{c}\boldsymbol{e}_{j} = \boldsymbol{\phi}_{c,j}$, respectively $\boldsymbol{e}_{j}^{\top}\boldsymbol{\Phi}_{c} = \boldsymbol{\phi}_{c,j}^{\top}$ for $\boldsymbol{\Phi}_{c} = [\boldsymbol{\phi}_{c,1}, \ldots, \boldsymbol{\phi}_{c,N}]^{\top}$ and $\boldsymbol{e}_{j}^{\top}\boldsymbol{\Phi}_{c}\boldsymbol{e}_{j} = \sigma_{c,j}^{2}$, this can be simplified to

$$\boldsymbol{\Phi}_{\rm c}^{\circ} = \boldsymbol{\Phi}_{\rm c} - \frac{\Delta \Lambda_{{\rm c},j}}{1 + \Delta \Lambda_{{\rm c},j} \sigma_{{\rm c},j}^2} \boldsymbol{\phi}_{{\rm c},j} \boldsymbol{\phi}_{{\rm c},j}^{\top} .$$
(5.38)



Figure 5.2: Procedure for sequential processing of variable nodes [SF21].

The corresponding new conditional mean m_c° can be calculated from the "old" one m_c by inserting the result into (5.33) and using (5.34)

$$\boldsymbol{m}_{c}^{\circ} = \boldsymbol{\varPhi}_{c}^{\circ} \left(\frac{1}{\sigma_{n}^{2}} \boldsymbol{A}^{\top} \boldsymbol{y} + \boldsymbol{\lambda}_{c}^{\circ}\right)$$

$$= \left(\boldsymbol{\varPhi}_{c} - \frac{\Delta \Lambda_{c,j}}{1 + \Delta \Lambda_{c,j} \sigma_{c,j}^{2}} \boldsymbol{\varPhi}_{c,j} \boldsymbol{\varPhi}_{c,j}^{\top}\right) \left(\frac{1}{\sigma_{n}^{2}} \boldsymbol{A}^{\top} \boldsymbol{y} + \boldsymbol{\lambda}_{c} + \Delta \lambda_{c,j} \boldsymbol{e}_{j}\right)$$

$$= \boldsymbol{\varPhi}_{c} \left(\frac{1}{\sigma_{n}^{2}} \boldsymbol{A}^{\top} \boldsymbol{y} + \boldsymbol{\lambda}_{c}\right) + \boldsymbol{\varPhi}_{c} \Delta \lambda_{c,j} \boldsymbol{e}_{j} - \frac{\Delta \Lambda_{c,j}}{1 + \Delta \Lambda_{c,j} \sigma_{c,j}^{2}} \boldsymbol{\varPhi}_{c,j} \boldsymbol{\varPhi}_{c,j}^{\top} \left(\frac{1}{\sigma_{n}^{2}} \boldsymbol{A}^{\top} \boldsymbol{y} + \boldsymbol{\lambda}_{c}\right)$$

$$- \frac{\Delta \Lambda_{c,j} \Delta \lambda_{c,j}}{1 + \Delta \Lambda_{c,j} \sigma_{c,j}^{2}} \boldsymbol{\varPhi}_{c,j} \boldsymbol{\varPhi}_{c,j}^{\top} \boldsymbol{e}_{j}.$$
(5.39)

Noticing that $m_{c,j} = \boldsymbol{\phi}_{c,j}^{\top} (\frac{1}{\sigma_n^2} \boldsymbol{A}^{\top} \boldsymbol{y} + \boldsymbol{\lambda}_c)$ and $\boldsymbol{\phi}_{c,j} \boldsymbol{e}_j = \sigma_{c,j}^2$, this may be simplified to

$$\boldsymbol{m}_{c}^{\circ} = \boldsymbol{m}_{c} + \boldsymbol{\phi}_{c,j} \Delta \lambda_{c,j} \left(1 - \frac{\Delta \Lambda_{c,j} \sigma_{c,j}^{2}}{1 + \Delta \Lambda_{c,j} \sigma_{c,j}^{2}} \right) - \frac{\Delta \Lambda_{c,j} m_{c,j}}{1 + \Delta \Lambda_{c,j} \sigma_{c,j}^{2}} \boldsymbol{\phi}_{c,j}$$
$$= \boldsymbol{m}_{c} + \frac{\Delta \lambda_{c,j} - \Delta \Lambda_{c,j} m_{c,j}}{1 + \Delta \Lambda_{c,j} \sigma_{c,j}^{2}} \boldsymbol{\phi}_{c,j} .$$
(5.40)

This way, a change induced by signal-constrained estimation of one variable, can simply be propagated to the other variables. The procedure is depicted in Fig. 5.2; starting with a signal-constrained estimation, the change in the extrinsic parameters is fed into the channel-constrained estimation and via the rank-one update, the impact on other variables can be computed. Then, another variable can be chosen.

The procedure of the resulting algorithm is summarized in the section below.

5.1.4 Sequential Version of VAMPind

The basis for the algorithm is the linear estimate m_c , which is sequentially refined by single signal-constrained estimations. Therefore, the algorithm is initialized by computing the LMMSE estimate (5.33) and (5.32) once with suitable values for λ_c and Λ_c .

For the refinement, a component x_j of the signal vector x is chosen. As it is known from (5.18), the extrinsic is calculated via a subtraction. Hence, the input to the signal-constrained

estimator is computed via

$$\tilde{x}_{s,j}/\tilde{\sigma}_{s,j}^2 = m_{c,j}/\sigma_{c,j}^2 - \lambda_{c,j} , \qquad 1/\tilde{\sigma}_{s,j}^2 = 1/\sigma_{c,j}^2 - \Lambda_{c,j} ,$$
 (5.41)

so that the signal-constrained estimate can be computed as usual by

$$m_{\mathbf{s},j} = \mathbf{E}_{\mathbf{x}_j} \left\{ \mathbf{x}_j \mid \tilde{x}_{\mathbf{s},j}, \, \tilde{\sigma}_{\mathbf{s},j}^2 \right\}, \tag{5.42}$$

$$\sigma_{s,j}^{2} = E_{x_{j}} \left\{ (x_{j} - m_{s,j})^{2} \mid \tilde{x}_{s,j}, \, \tilde{\sigma}_{s,j}^{2} \right\} \,.$$
(5.43)

The connection (5.18) is used afterwards for the extrinsic calculation again, so that

$$\lambda_{c,j}^{\circ} = m_{s,j} / \sigma_{s,j}^2 - \tilde{x}_{s,j} / \tilde{\sigma}_{s,j}^2 , \quad \Lambda_{c,j}^{\circ} = 1 / \sigma_{s,j}^2 - 1 / \tilde{\sigma}_{s,j}^2 , \qquad (5.44)$$

The additive parameters $\Delta \lambda_{c,j}$ and $\Delta \Lambda_{c,j}$ in (5.34) and (5.35), respectively, thus, read

$$\begin{aligned} \Delta\lambda_{\mathrm{c},j} &= \lambda_{\mathrm{c},j}^{\circ} - \lambda_{\mathrm{c},j} = m_{\mathrm{s},j}/\sigma_{\mathrm{s},j}^{2} - \tilde{x}_{\mathrm{s},j}/\tilde{\sigma}_{\mathrm{s},j}^{2} - \lambda_{\mathrm{c},j} \\ &= m_{\mathrm{s},j}/\sigma_{\mathrm{s},j}^{2} - (m_{\mathrm{c},j}/\sigma_{\mathrm{c},j}^{2} - \lambda_{\mathrm{c},j}) - \lambda_{\mathrm{c},j} = m_{\mathrm{s},j}/\sigma_{\mathrm{s},j}^{2} - m_{\mathrm{c},j}/\sigma_{\mathrm{c},j}^{2} , \end{aligned}$$

$$\begin{aligned} \Delta\Lambda_{\mathrm{c},j} &= \Lambda_{\mathrm{c},j}^{\circ} - \Lambda_{\mathrm{c},j} \end{aligned}$$
(5.45)

$$= 1/\sigma_{s,j}^{2} - 1/\tilde{\sigma}_{s,j}^{2} - \lambda_{c,j}$$

= $1/\sigma_{s,j}^{2} - (1/\sigma_{c,j}^{2} - \Lambda_{c,j}) - \Lambda_{c,j}$
= $1/\sigma_{s,j}^{2} - 1/\sigma_{c,j}^{2}$, (5.46)

With this, the rank-one update (5.38) simplifies to

$$\begin{split} \boldsymbol{\varPhi}_{c}^{o} &= \boldsymbol{\varPhi}_{c} - \frac{1/\sigma_{s,j}^{2} - 1/\sigma_{c,j}^{2}}{1 + (1/\sigma_{s,j}^{2} - 1/\sigma_{c,j}^{2})\sigma_{c,j}^{2}} \boldsymbol{\varPhi}_{c,j} \boldsymbol{\varPhi}_{c,j}^{\top} \\ &= \boldsymbol{\varPhi}_{c} - \frac{\sigma_{s,j}^{2}}{\sigma_{c,j}^{2}} \left(\frac{1}{\sigma_{s,j}^{2}} - \frac{1}{\sigma_{c,j}^{2}} \right) \boldsymbol{\varPhi}_{c,j} \boldsymbol{\varPhi}_{c,j}^{\top} \\ &= \boldsymbol{\varPhi}_{c} + \frac{1}{(\sigma_{c,j}^{2})^{2}} (\sigma_{s,j}^{2} - \sigma_{c,j}^{2}) \boldsymbol{\varPhi}_{c,j} \boldsymbol{\varPhi}_{c,j}^{\top} , \end{split}$$
(5.47)

which yields the interesting interpretation that the covariance matrix is adjusted by the difference of the variances of the two estimators. Basically that means, at position j it is

$$[\boldsymbol{\Phi}_{\mathrm{c}}^{\circ}]_{jj} = \sigma_{\mathrm{s},j}^2 , \qquad (5.48)$$

i.e., the previous variance $\sigma_{c,j}^2$ is replaced by the variance $\sigma_{s,j}^2$ obtained in the signal-constrained estimation and at any other position $j' \neq j$

$$[\boldsymbol{\Phi}_{\rm c}^{\circ}]_{j'j'} = \sigma_{{\rm c},j'}^2 + \frac{(\sigma_{{\rm c},j'}^2)^2}{(\sigma_{{\rm c},j}^2)^2} (\sigma_{{\rm s},j}^2 - \sigma_{{\rm c},j}^2) .$$
(5.49)

For the update of the conditional mean, one obtains ($m{m}^\circ_{
m c} = [m^\circ_{
m c,1},\,\ldots,\,m^\circ_{
m c,N}]^ op$)

$$\boldsymbol{m}_{c}^{\circ} = \boldsymbol{m}_{c} + \frac{\sigma_{s,j}^{2}}{\sigma_{c,j}^{2}} \left(\Delta \lambda_{c,j} - \Delta \Lambda_{c,j} \boldsymbol{m}_{c,j} \right) \boldsymbol{\phi}_{c,j}$$

$$= \boldsymbol{m}_{c} + \frac{\sigma_{s,j}^{2}}{\sigma_{c,j}^{2}} \left(\frac{\boldsymbol{m}_{s,j}}{\sigma_{s,j}^{2}} - \frac{\boldsymbol{m}_{c,j}}{\sigma_{c,j}^{2}} - \left(\frac{1}{\sigma_{s,j}^{2}} - \frac{1}{\sigma_{c,j}^{2}} \right) \boldsymbol{m}_{c,j} \right) \boldsymbol{\phi}_{c,j}$$

$$= \boldsymbol{m}_{c} + \frac{1}{\sigma_{c,j}^{2}} (\boldsymbol{m}_{s,j} - \boldsymbol{m}_{c,j}) \boldsymbol{\phi}_{c,j} , \qquad (5.50)$$

i.e., the update is based on the difference between the conditional mean estimates. Here, again $m_{c,j}$ is replaced by $m_{s,j}$ and at other positions $j' \neq j$

$$m_{c,j'}^{\circ} = m_{c,j'} + \frac{\sigma_{c,j'}^2}{\sigma_{c,j}^2} (m_{s,j} - m_{c,j}) , \qquad (5.51)$$

which means the estimate is updated by the scaled difference between the conditional means of the two opposing estimators, where the scaling is defined by the relation between the variance at the considered position j' and the variance of the position that has been updated.

At this point, the next position needs to be chosen to calculate another signal-constrained estimate and further refine the channel-constrained estimate. The complete algorithm is stated in Appendix E.9 as seqVAMPind. The denomination is chosen because it is the sequential version of VAMPind. The algorithm can be found in [OW05, Appendix D] for a binary prior and in [BMPP20] for a Bernoulli-Gaussian prior.

5.1.5 Sequential Version with Estimation-Theoretic Bias Compensation

In Chapter 4 we have seen that an estimation-theoretic view on the extrinsic calculation can lead to different procedures than derived from the optimization-based view. These updates have been used to compose the VAMPire algorithm, which has been proposed in [FSG20]. Similarly, the respective updates will be used here to obtain a variant of the previously explained algorithm seqVAMPind. The resulting algorithm is called seqVAMPire, as it is a <u>seq</u>uential version of the VAMPire algorithm. It was first published in [SF21].

The start of the algorithm is identically to before, i.e., the channel-constrained estimates (5.33) and (5.32) are computed with suitably chosen input parameters λ_c and Λ_c . The ensuing extrinsic computation is equal for estimation-theoretic and optimization-based view, i.e., position j is chosen and (5.41), as well as the signal-constrained estimates as in (5.42) and (5.43) are computed. The estimation-theoretic bias compensation after the signal-constrained estimation differs from the optimization-based extrinsic, since the estimator is non-linear. Instead of the subtraction of natural parameters, the updates (4.86) and (4.86) are used, which are stated here for convenience again

$$\tilde{x}_{\mathrm{c},j} = \tilde{x}_{\mathrm{s},j} + \frac{\tilde{\sigma}_{\mathrm{s},j}^2}{\tilde{\sigma}_{\mathrm{s},j}^2 - \sigma_j^2} (m_{\mathrm{s},j} - \tilde{x}_{\mathrm{s},j}) , \qquad (5.52)$$

$$\tilde{\sigma}_{c,j}^2 = \sigma_{s,j}^2 + \left(\frac{\sigma_j^2}{\sigma_j^2 - \tilde{\sigma}_{s,j}^2}\right)^2 (m_{s,j} - \tilde{x}_{s,j})^2 , \qquad (5.53)$$

where $\sigma_j^2 = E_{\tilde{x}_{s,j}} \{ \sigma_{s,j}^2(\tilde{x}_{s,j}) \}$ is the MSE of the respective prior. Noteworthy, with

$$1/\tilde{\sigma}_{j}^{2} = 1/\sigma_{j}^{2} - 1/\tilde{\sigma}_{\mathrm{s},j}^{2}$$
(5.54)

the first line can be rewritten to

$$\tilde{x}_{\mathrm{c},j}/\tilde{\sigma}_j^2 = m_{\mathrm{s},j}/\sigma_j^2 - \tilde{x}_{\mathrm{s},j}/\tilde{\sigma}_{\mathrm{s},j}^2 \,. \tag{5.55}$$

One has to be aware, that this differs from the update in (5.44) since the MSE σ_j^2 and its unbiased version $\tilde{\sigma}_j^2$ are used instead of $\sigma_{\mathrm{s},j}^2$ and $\tilde{\sigma}_{\mathrm{c},j}^2$, respectively. Also, $\tilde{\sigma}_{\mathrm{c},j}^2$ is computed differently here. All in all, the updated parameters are obtained as

$$\lambda_{\mathrm{c},j}^{\circ} = \tilde{x}_{\mathrm{c},j} / \tilde{\sigma}_{\mathrm{c},j}^2 = \frac{\tilde{\sigma}_j^2}{\tilde{\sigma}_{\mathrm{c},j}^2} (m_{\mathrm{s},j} / \sigma_j^2 - \tilde{x}_{\mathrm{s},j} / \tilde{\sigma}_{\mathrm{s},j}^2) , \qquad (5.56)$$

$$\Lambda^{\circ}_{\mathrm{c},j} = 1/\tilde{\sigma}^2_{\mathrm{s},j} , \qquad (5.57)$$

and the additive components, thus, as

$$\Delta\lambda_{\mathrm{c},j} = \frac{\tilde{\sigma}_{j}^{2}}{\tilde{\sigma}_{\mathrm{c},j}^{2}} (m_{\mathrm{s},j}/\sigma_{j}^{2} - \tilde{x}_{\mathrm{s},j}/\tilde{\sigma}_{\mathrm{s},j}^{2}) - \lambda_{\mathrm{c},j}$$

$$= \frac{\tilde{\sigma}_{j}^{2}}{\tilde{\sigma}_{\mathrm{c},j}^{2}} (m_{\mathrm{s},j}/\sigma_{j}^{2} - m_{\mathrm{c},j}/\sigma_{\mathrm{c},j}^{2}) + \lambda_{\mathrm{c},j} \left(\frac{\tilde{\sigma}_{j}^{2}}{\tilde{\sigma}_{\mathrm{c},j}^{2}} - 1\right) , \qquad (5.58)$$

$$\Delta \Lambda_{\mathrm{c},j} = 1/\tilde{\sigma}_{\mathrm{c},j}^2 - \Lambda_{\mathrm{c},j} .$$
(5.59)

These equations have the effect that the rank-one updates (5.38) and (5.40) cannot be simplified to easier scaling operations as in the previous section. Nevertheless, the changes induced by the single signal-constrained estimation can be computed and another position for refinement may be chosen to proceed with the recovery.

The entire procedure is stated in Appendix E.10 as seqVAMPire(sequential version of VAM-Pire).

So far, we have focused on the steps of updating the changes induced by a single signalconstrained estimation. For the overall procedure of the algorithm, one has to think of the order of processing the variables one after the other, as the sequential approach opens the possibilities for various schedules. The discussion of these opportunities is postponed to Chapter 6.

5.2 Sequential Processing of Observations

In this section, a CS model is considered that slightly differs from the one that is introduced in Chapter 2. So far, it has been assumed that the noise vector **n** is i.i.d., which means each element of the vector is distributed in the same way, and especially the noise variance σ_n^2 is the same for all elements. In the following, the case of noise with individual variances is considered, i.e., **n** ~ $\mathcal{N}(\mathbf{0}, \operatorname{diag}(\sigma_{n,i}^2))$, respectively $\mathbf{n}_i \sim \mathcal{N}(0, \sigma_{n,i}^2)$ for $i \in \{1, \ldots, M\}$, where it is assumed that the individual variances are known. The entries of the noise vector are thereby still independent, but not identically distributed anymore. We call this case the non-i.i.d. noise case. Such a model is justified, if some form of channel equalization, which colors the noise, is performed on the measurements before reconstructing the signal, which happens, e.g., in several PAR reduction schemes for OFDM [KPN⁺14, LHD20], or when sparse noise is assumed [SSF13].

Since the observations y_i carry different reliabilities in this model, one may benefit from a sequential processing of the observations. This will be examined in the following based on the optimization-based approach that was considered in the previous chapters as well.



Figure 5.3: Factor graph for sequential processing of the observations.

Sequential processing of observations has been considered in [HTP12, MSW10], where one observation is processed after the other and the reconstruction is stopped if enough observations have been processed. In contrast, here, a particular observation vector \boldsymbol{y} is measured once and the algorithm works on the individual observations of the vector, so that the particular schedule is in the focus. Other work that considered variable noise variances in the CS scenario is given in [BG18], which introduces an algorithm based on AMP. Also the GAMP algorithm [Ran11], which is also stated in Appendix E can cope with the variable noise, if the noise variance σ_n^2 is replaced by the respective individual noise variances. As before, the work at hand is based on VAMP instead of AMP.

5.2.1 Structure of the Factor Graph

For the considerations that are pursued here, the channel factors are considered $f_{y_i|x}(x)$ independently, since each one should get its own efficacy. In order not to complicate the problem too much—separating the signal components as well, would lead to the MP case of Chapter 3 anyway—the signal-constrained part is kept together. Hence, the factorization of the posterior reads

$$f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) = \text{const.} \cdot f_{\mathbf{x}}(\boldsymbol{x}) \cdot \prod_{i=1}^{M} f_{\mathbf{y}_{i}|\mathbf{x}}(\boldsymbol{x}) , \qquad (5.60)$$

and the resulting factor graph is the remaining one in Fig. 2.3, i.e., Fig. 2.3(b), which is plotted in Fig. 5.3 for convenience again.

5.2.1.1 Factorization of the Substitute Distribution

The factorization above tells us that $q_s(\boldsymbol{x})$ is needed as approximation for $f_{\boldsymbol{x}}(\boldsymbol{x})$, as well as the individual $q_{c,i}(\boldsymbol{x})$ for approximating $f_{y_i|\boldsymbol{x}}(\boldsymbol{x})$ ($i \in \{1, \ldots, M\}$), i.e.,

$$q_s(\boldsymbol{x}) = f_{\boldsymbol{x}}(\boldsymbol{x})s(\boldsymbol{x}) , \qquad (5.61)$$

$$\mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) = \mathbf{f}_{\mathbf{y}_i|\mathbf{x}}(\boldsymbol{x})c_i(\boldsymbol{x}) .$$
(5.62)

The (combined) variable node is therefore constituted by M + 1 instances, which leads to an M-fold integration as in Chapter 3, i.e., the internal behavior of the variable node is given by

$$\frac{\text{const.}}{s(\boldsymbol{x})\prod_{i=1}^{M}c_{i}(\boldsymbol{x})} = \text{const.} \cdot \int \cdots \int \left(s(\boldsymbol{x})\prod_{i=1}^{M}c_{i}(\boldsymbol{x}_{i})\right)^{-1}\prod_{i=1}^{M}\delta(\boldsymbol{x}-\boldsymbol{x}_{i})\,\mathrm{d}\boldsymbol{x}_{1}\ldots\boldsymbol{x}_{M}$$
(5.63)

with slack variables x_1, \ldots, x_M (the argument of s(x) is chosen to be the original argument for convenience). So, the connection between variable node and deviations is obtained as

$$\mathbf{q}_{\mathbf{v}}(\boldsymbol{x}) = \text{const.} \cdot \sqrt[M]{s(\boldsymbol{x}) \prod_{i=1}^{M} c_i(\boldsymbol{x})} .$$
 (5.64)

The resulting factorization reads

$$\mathbf{q}(\boldsymbol{x}) = \frac{\mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \cdot \prod_{i=1}^{M} \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x})}{(\mathbf{q}_{\mathrm{v}}(\boldsymbol{x}))^{M}} .$$
(5.65)

5.2.1.2 The Moment-Matching Constraint

As can be seen in the factorization (5.60), respectively the factor graph in Fig. 5.3 above, all factors are connected to the (combined) variable node, which represents the entire vector \boldsymbol{x} , in the center. The matching of moments for the sufficient statistics $\boldsymbol{g}(\boldsymbol{x})$ is therefore set up between the representative of the variable node $q_v(\boldsymbol{x})$ and the respective representatives of the factors, yielding

$$\mathbf{E}_{\mathbf{x}\sim q_{\mathbf{v}}}\{\boldsymbol{g}(\mathbf{x})\} \stackrel{!}{=} \mathbf{E}_{\mathbf{x}\sim q_{\mathbf{s}}}\{\boldsymbol{g}(\mathbf{x})\}, \qquad (5.66)$$

$$\mathbf{E}_{\mathbf{x}\sim q_{\mathbf{v}}}\{\boldsymbol{g}(\mathbf{x})\} \stackrel{!}{=} \mathbf{E}_{\mathbf{x}\sim q_{\mathbf{c},i}}\{\boldsymbol{g}(\mathbf{x})\}, \qquad i \in \{1, \ldots, M\}.$$
(5.67)

Here, the expectation constraint is vector-wise, i.e., the focus is not on individual signal components x_j . Therefore it is possible to restrict to the average variance case, i.e., $\boldsymbol{g}(\boldsymbol{x}) = [\boldsymbol{x}^{\top}, -\frac{1}{2}\boldsymbol{x}^{\top}\boldsymbol{x}]^{\top}$ from Ex. D.1 and (D.12), which means the reliability of the entire vector is represented with a single variance.

5.2.2 Constrained Optimization

The constrained optimization problem that is obtained, when combining minimization of the Kullback–Leibler divergence with the moment-matching constraints above (and additional normalization constraints for valid pdfs) reads

$$\begin{array}{ll} \text{minimize} & \mathrm{D}_{\mathrm{KL}} \Big(\mathbf{q}(\boldsymbol{x}) \mid\mid \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x}) \Big) \\ \text{subject to} & \mathrm{E}_{\mathbf{x} \sim q_{\mathrm{v}}} \{ \boldsymbol{g}(\mathbf{x}) \} - \mathrm{E}_{\mathbf{x} \sim q_{\mathrm{c},i}} \{ \boldsymbol{g}(\mathbf{x}) \} = \mathbf{0} , \qquad i \in \{1, \ldots, M\} , \\ & \mathrm{E}_{\mathbf{x} \sim q_{\mathrm{v}}} \{ \boldsymbol{g}(\mathbf{x}) \} - \mathrm{E}_{\mathbf{x} \sim q_{\mathrm{s}}} \{ \boldsymbol{g}(\mathbf{x}) \} = \mathbf{0} , \qquad (5.68) \\ & \int \mathbf{q}_{\mathrm{v}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = 1 , \quad \int \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = 1 , \\ & \int \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = 1 , \qquad i \in \{1, \ldots, M\} . \end{array}$$

5.2.2.1 Cost Function for Compressed Sensing

The corresponding Kullback–Leibler divergence between substitute distribution q(x) and posterior $f_{x|y}(x)$, which needs to be minimized, is here given by

$$D_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) \mid\mid f_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})) = \int \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \log \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \,\mathrm{d}\boldsymbol{x} + \sum_{i=1}^{M} \int \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \log \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \,\mathrm{d}\boldsymbol{x}$$
$$- \int \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \log f_{\mathbf{x}}(\boldsymbol{x}) \,\mathrm{d}\boldsymbol{x} - \sum_{i=1}^{M} \int \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \log f_{\mathbf{y}_{i}|\mathbf{x}}(\boldsymbol{x}) \,\mathrm{d}\boldsymbol{x}$$
$$- M \int \mathbf{q}_{\mathrm{v}}(\boldsymbol{x}) \log \mathbf{q}_{\mathrm{v}}(\boldsymbol{x}) \,\mathrm{d}\boldsymbol{x} .$$
(5.69)

5.2.2.2 Stationary Points of the Lagrangian

In order to find the basic relations that a solution to the CS problem with the given structural conditions need to fulfill, the stationary points of the respective Lagrangian are computed. To that end, the Lagrangian multipliers $\nu_{c,i}$ and ν_s are defined for the moment-matching constraints, as well as ν_v , ν_s , and $\nu_{c,i}$ for the normalization constraints, respectively. Then, the Lagrangian reads

$$\mathcal{L}(\mathbf{q}(\boldsymbol{x})) = D_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) || \mathbf{f}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{x})) + \sum_{i=1}^{M} \boldsymbol{\nu}_{\mathrm{c},i}^{\top} \left(\mathrm{E}_{\mathbf{x} \sim \mathbf{q}_{\mathrm{v}}} \{ \boldsymbol{g}(\mathbf{x}) \} - \mathrm{E}_{\mathbf{x} \sim \mathbf{q}_{\mathrm{c},i}} \{ \boldsymbol{g}(\mathbf{x}) \} \right) + \boldsymbol{\nu}_{\mathrm{s}}^{\top} \left(\mathrm{E}_{\mathbf{x} \sim \mathbf{q}_{\mathrm{v}}} \{ \boldsymbol{g}(\mathbf{x}) \} - \mathrm{E}_{\mathbf{x} \sim \mathbf{q}_{\mathrm{s}}} \{ \boldsymbol{g}(\mathbf{x}) \} \right) + \sum_{j=1}^{N} \left(\nu_{\mathrm{v}}(\int \mathbf{q}_{\mathrm{v}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1) + \nu_{\mathrm{s}}(\int \mathbf{q}_{\mathrm{s}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1) \right) + \sum_{i=1}^{M} \nu_{\mathrm{c},i}(\int \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1) .$$
(5.70)

As usual, the stationary points are obtained by differentiating w.r.t. the representatives $q_v(\boldsymbol{x})$, $q_{c,i}(\boldsymbol{x})$, and $q_s(\boldsymbol{x})$, respectively, and setting the result to zero, which yields

$$\frac{\partial}{\partial \mathbf{q}_{\mathrm{v}}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = -M \log \mathbf{q}_{\mathrm{v}}(\boldsymbol{x}) + \left(\boldsymbol{\nu}_{\mathrm{s}} + \sum_{i=1}^{M} \boldsymbol{\nu}_{\mathrm{c},i}\right)^{\top} \boldsymbol{g}(\boldsymbol{x}) + \text{const.} \stackrel{!}{=} 0, \quad (5.71)$$

$$\frac{\partial}{\partial \mathbf{q}_{\mathrm{c},i}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = \log \mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) - \log \mathbf{f}_{\mathbf{y}_i|\mathbf{x}}(\boldsymbol{x}) - \boldsymbol{\nu}_{\mathrm{c},i}^{\top} \boldsymbol{g}(\boldsymbol{x}) + \text{const.} \stackrel{!}{=} 0 , \qquad (5.72)$$

$$\frac{\partial}{\partial q_{s}} \mathcal{L}(\mathbf{q}(\boldsymbol{x})) = \log q_{s}(\boldsymbol{x}) - \log f_{\mathbf{x}}(\boldsymbol{x}) - \boldsymbol{\nu}_{s}^{\top} \boldsymbol{g}(\boldsymbol{x}) + \text{const.} \stackrel{!}{=} 0.$$
 (5.73)

With $Z_{\rm v}$, $Z_{{\rm c},i}$, and $Z_{\rm s}$ functioning as normalizing factors, the stationary-point forms for the

representatives read

$$\mathbf{q}_{\mathbf{v}}(\boldsymbol{x}) = \frac{1}{Z_{\mathbf{v}}} \exp\left(\frac{1}{M} (\boldsymbol{\nu}_{\mathbf{s}} + \sum_{i=1}^{M} \boldsymbol{\nu}_{\mathbf{c},i})^{\mathsf{T}} \boldsymbol{g}(\boldsymbol{x})\right) , \qquad (5.74)$$

$$\mathbf{q}_{\mathrm{c},i}(\boldsymbol{x}) = \frac{1}{Z_{\mathrm{c},i}} \mathbf{f}_{\mathbf{y}_i|\mathbf{x}}(\boldsymbol{x}) \exp\left(\boldsymbol{\nu}_{\mathrm{c},i}^{\top} \boldsymbol{g}(\boldsymbol{x})\right) , \qquad (5.75)$$

$$q_{s}(\boldsymbol{x}) = \frac{1}{Z_{s}} f_{\boldsymbol{x}}(\boldsymbol{x}) \exp\left(\boldsymbol{\nu}_{s}^{\top} \boldsymbol{g}(\boldsymbol{x})\right) .$$
(5.76)

Again, as expected, exponential families are involved.

The respective natural parameter of the (combined) variable node representative $q_v(x)$ is given by

$$\boldsymbol{\theta} \stackrel{\text{\tiny def}}{=} \frac{1}{M} \left(\boldsymbol{\nu}_{s} + \sum_{i=1}^{M} \boldsymbol{\nu}_{c,i} \right) .$$
(5.77)

This reminds us of the connection (3.25) of Chapter 3 and will be used in the following to derive the parameter update for the recovery algorithm tailored to the structure given by the factor graph considered here. The correspondences between edges and the respective edge-dependent natural parameters are depicted in Fig. 5.4(a).

5.2.3 Algorithm

Based on the stationary points of the given structure, an algorithm for the recovery of signal x is derived in the following.

5.2.3.1 Expectations

First, the computations at the factor nodes are regarded, which are expectations with input parameters $\boldsymbol{\nu}_{c,i}$ ($i \in \{1, ..., M\}$) and $\boldsymbol{\nu}_s$, respectively. Since only the average variance case is considered and the estimates of all signal vector elements are computed at once, one can (as in Sec. 4.3.1) connect the signal-constrained estimation at factor $f_x(\boldsymbol{x})$ to (2.19) and (2.20) by making the identification $\boldsymbol{\nu}_s = [\tilde{\boldsymbol{x}}_s^\top / \tilde{\sigma}_s^2, 1/ \tilde{\sigma}_s^2]^\top$, i.e.,

$$\mathbf{E}_{\mathbf{x}\sim q_{\mathrm{s}}}\{\boldsymbol{g}_{\boldsymbol{\lambda}}(\mathbf{x})\} = \boldsymbol{m}_{\mathrm{s}} , \qquad (5.78)$$

$$E_{\mathbf{x} \sim \mathbf{q}_{s}}\{g_{\Lambda}(\mathbf{x})\} = -\frac{1}{2} \sum_{j=1}^{N} (\sigma_{s,j}^{2} + m_{s,j}^{2}) .$$
(5.79)

The corresponding average variance σ_s^2 is obtained as in (4.35), i.e.,

$$\sigma_{\rm s}^2 = \frac{1}{N} \sum_{j=1}^N \sigma_{{\rm s},j}^2 = -\frac{2}{N} \operatorname{E}_{\mathbf{x} \sim \mathfrak{q}_{\rm s}} \{g_A(\mathbf{x})\} - \boldsymbol{m}_{\rm s}^\top \boldsymbol{m}_{\rm s} .$$
(5.80)

The expectations at channel factors $f_{y_i|\mathbf{x}}(\mathbf{x})$ can be related to the estimates in Sec. 2.2.1.2 by considering the average variance case, i.e., $\boldsymbol{\Phi}_{\mathbf{x}} = \sigma_{\mathbf{x}}^2 \boldsymbol{I}_N$, and making the identification $\boldsymbol{\nu}_{\mathrm{c},i} = [\boldsymbol{m}_{\mathrm{x}}^{\top}/\sigma_{\mathrm{x}}^2,\,1/\sigma_{\mathrm{x}}^2]^{\top}.$ The estimation reads then

$$E_{\mathbf{x} \sim q_{c,i}} \{ \boldsymbol{g}_{\boldsymbol{\lambda}}(\mathbf{x}) \} = \boldsymbol{m}_{c,i} = \boldsymbol{m}_{\mathbf{x}} + \sigma_{\mathbf{x}}^{2} \boldsymbol{a}_{i} (\sigma_{\mathbf{x}}^{2} \boldsymbol{a}_{i}^{\top} \boldsymbol{a}_{i} + \sigma_{\mathbf{n},i}^{2})^{-1} (y_{i} - \boldsymbol{a}_{i}^{\top} \boldsymbol{m}_{\mathbf{x}})$$
$$= \boldsymbol{m}_{\mathbf{x}} + \frac{y_{i} - \sum_{j=1}^{N} a_{ij} m_{\mathbf{x},j}}{\sigma_{\mathbf{n},i}^{2} / \sigma_{\mathbf{x}}^{2} + \sum_{j=1}^{N} a_{ij}^{2}} \boldsymbol{a}_{i} , \qquad (5.81)$$

$$\mathbf{E}_{\mathbf{x}\sim \mathbf{q}_{\mathrm{c},i}}\{g_{A}(\mathbf{x})\} = -\frac{1}{2} \operatorname{trace}\left(\boldsymbol{\varPhi}_{\mathrm{c},i} + \boldsymbol{m}_{\mathrm{c},i}\boldsymbol{m}_{\mathrm{c},i}^{\top}\right) \\
= -\frac{1}{2} \operatorname{trace}\left(\sigma_{\mathsf{n},i}^{2}(\boldsymbol{a}_{i}\boldsymbol{a}_{i}^{\top} + \frac{\sigma_{\mathsf{n},i}^{2}}{\sigma_{\mathsf{x}}^{2}}\boldsymbol{I}_{N})^{-1} + \boldsymbol{m}_{\mathrm{c},i}\boldsymbol{m}_{\mathrm{c},i}^{\top}\right).$$
(5.82)

Furthermore, the resulting average variance of the estimate $m_{ ext{c},i}$ reads

$$\sigma_{\mathrm{c},i}^{2} = \frac{1}{N} \operatorname{trace}(\boldsymbol{\Phi}_{\mathrm{c},i}) = \frac{\sigma_{\mathrm{n},i}^{2}}{N} \operatorname{trace}\left((\boldsymbol{a}_{i}\boldsymbol{a}_{i}^{\top} + \frac{\sigma_{\mathrm{n},i}^{2}}{\sigma_{\mathrm{x}}^{2}}\boldsymbol{I}_{N})^{-1}\right).$$
(5.83)

5.2.3.2 Projection to Exponential Family

The moment-matching between the variable node and the respective factor node, which shall ensure consistency, means that the computed moments of the section above have to represent the variable node distribution. Therefore, the natural parameter of $q_v(\boldsymbol{x})$ can directly be obtained from these moments via

$$\boldsymbol{\theta} = [\boldsymbol{m}_{\rm s}^{\top} / \sigma_{\rm s}^2, \, 1 / \sigma_{\rm s}^2]^{\top} \,, \tag{5.84}$$

respectively,

$$\boldsymbol{\theta} = [\boldsymbol{m}_{c,i}^{\top} / \sigma_{c,i}^2, 1 / \sigma_{c,i}^2]^{\top}, \quad i \in \{1, \dots, M\}.$$
(5.85)

For the further processing, one needs to consider the restrictions on the natural parameters induced by the structure of the problem, which finds its expression in the natural parameter (5.77) of the representative $q_v(x)$ of the variable node.

5.2.3.3 Extrinsic

In order to make use of the connection (5.77) for the processing of (natural) parameters in the algorithm, recall that the basic idea for the iterative schemes aiming towards stationary points is the notion of extrinsic, i.e., only knowledge obtained via the processing is submitted for further processing (in order to prevent feedback loops). Compared to in Chapter 3, the derivation is reversed, i.e., it starts with the processing that sums up all incoming edges except the one that the message is going to be sent on, to see that in the end, the same connections hold. To apply the notion here, the Lagrangian multipliers are interpreted as messages on the corresponding edges and a distinction between the messages in different directions is made. As the multipliers itself specify the representative of the variable node $q_v(x)$ entirely, they are interpreted as *outgoing* messages of the variable node and denoted by $\theta^{\setminus s} = \nu_s$ and $\theta^{\setminus i} = \nu_{c,i}$ ($i \in \{1, \ldots, M\}$), respectively. The messages on the backward direction, from the



(a) Depiction of edge-dependent multipliers.

(b) Depiction of edge- and direction-dependent messages.

Figure 5.4: Factor graph for the sequential processing of observations and edge-dependent parameters.

factor nodes towards the variable node will be denoted by $\theta^{(s)}$ and $\theta^{(i)}$ ($i \in \{1, ..., M\}$), respectively. The direction-dependent messages are shown on the respective edges in the factor graph in Fig. 5.4(b). The notation is chosen in order to show the connection to the expectation propagation (EP) [Min01] framework, where the idea is to replace parts of the approximations by the respective exact factor, in order to refine the estimation by the knowledge obtained by the factor. In the EP literature, the variables $\theta^{\setminus \bullet}$ are called the *cavity* parameters and $\theta^{(\bullet)}$ are the so-called *site* parameters.

The notion of extrinsic can now be established by ensuring that an outgoing message from variable node to factor node is only specified by incoming messages from factor nodes other than the one to send to, i.e.,

$$\boldsymbol{\theta}^{\mathsf{s}} = \sum_{i=1}^{M} \boldsymbol{\theta}^{(i)} , \qquad (5.86)$$

$$\boldsymbol{\theta}^{\setminus i} = \boldsymbol{\theta}^{(\mathrm{s})} + \sum_{\substack{i'=1\\i'\neq i}}^{M} \boldsymbol{\theta}^{(i')} , \quad i \in \{1, \dots, M\} .$$
(5.87)

By inserting this into the connection (5.77), one finds

$$\boldsymbol{\theta} = \frac{1}{M} \left(\sum_{i=1}^{M} \boldsymbol{\theta}^{(i)} + \sum_{i=1}^{M} \left(\boldsymbol{\theta}^{(s)} + \sum_{\substack{i'=1\\i'\neq i}}^{M} \boldsymbol{\theta}^{(i')} \right) \right)$$
$$= \frac{1}{M} \sum_{i=1}^{M} \boldsymbol{\theta}^{(i)} + \boldsymbol{\theta}^{(s)} + \frac{1}{M} (M-1) \sum_{i=1}^{M} \boldsymbol{\theta}^{(i)}$$
$$= \boldsymbol{\theta}^{(s)} + \sum_{i=1}^{M} \boldsymbol{\theta}^{(i)} .$$
(5.88)

A comparison to the definitions for the outgoing messages of the variable node, which are also the incoming messages for the factor nodes, shows that the inputs for the estimations at the factors can be obtained from the natural parameter θ of the variable node by

$$\boldsymbol{\theta}^{\mathrm{s}} = \boldsymbol{\theta} - \boldsymbol{\theta}^{\mathrm{s}}, \qquad (5.89)$$

$$\boldsymbol{\theta}^{i} = \boldsymbol{\theta} - \boldsymbol{\theta}^{(i)}, \quad i \in \{1, \dots, M\},$$
(5.90)

respectively, i.e., *post-processed* extrinsic calculations are at hand and the sum of messages in different directions yields, as in Chapter 3, the knowledge about the signal. This means that in order to be able to get an input to the expectations, it is necessary to store and keep track of the site parameters $\theta^{(\bullet)}$, which can be updated after the estimation via $\theta^{(\bullet)} = \theta - \theta^{\setminus \bullet}$ ($\bullet \in \{s, 1, \ldots, M\}$). With this knowledge, it is possible to state the entire procedure.

5.2.3.4 Procedure

When starting the algorithm with expectations, first suitable input parameters $\boldsymbol{\theta}^{\setminus \bullet}$ are required for all factors. For the simplification of notation, we introduce $\boldsymbol{\theta}^{\setminus \bullet} = [(\tilde{\boldsymbol{x}}^{\setminus \bullet})^\top / \sigma_{\setminus \bullet}^2, 1/\sigma_{\setminus \bullet}^2]^\top$, and $\boldsymbol{\theta}^{(\bullet)} = [(\tilde{\boldsymbol{x}}^{(\bullet)})^\top / \sigma_{(\bullet)}^2, 1/\sigma_{(\bullet)}^2]^\top$ ($\bullet \in \{s, 1, ..., M\}$) and use $\boldsymbol{\theta} = [\boldsymbol{m}^\top / \sigma^2, 1/\sigma^2]^\top$ to represent the moments that are mapped to the natural parameters of the variable node. The computation of the signal-constrained estimation is performed as usual,

$$\boldsymbol{m} = \mathbf{E}_{\mathbf{x}} \left\{ \mathbf{x} \mid \tilde{\boldsymbol{x}}^{\backslash \mathrm{s}}, \, \sigma_{\backslash \mathrm{s}}^{2} \right\}, \tag{5.91}$$

$$\sigma^{2} = \frac{1}{N} \sum_{j=1}^{N} \mathrm{E}_{\mathsf{x}} \{ (\mathsf{x} - m_{j})^{2} \mid \tilde{x}_{j}^{\backslash \mathsf{s}}, \, \sigma_{\backslash \mathsf{s}}^{2} \} \,, \tag{5.92}$$

and the site parameters are obtained by

$$\tilde{\boldsymbol{x}}^{(s)}/\sigma_{(s)}^2 = \boldsymbol{m}/\sigma^2 - \tilde{\boldsymbol{x}}^{\setminus s}/\sigma_{\setminus s}^2$$
, $1/\sigma_{(s)}^2 = 1/\sigma^2 - 1/\sigma_{\setminus s}^2$. (5.93)

For the channel-constrained estimations, one can spare the computation of the $N \times N$ inverse for the covariance matrix by using the knowledge about bias compensation from Sec. 4.4.1. Since average unbiasing is performed here, one can compute the unbiased estimate directly by computing the scaling factor based on the end-to-end cascade between \boldsymbol{x} and $\boldsymbol{m}_{\mathrm{c},i}$ via

$$w = \frac{1}{N} \operatorname{trace} \left(\boldsymbol{a}_i (\boldsymbol{a}_i^{\top} \boldsymbol{a}_i + \sigma_{\mathsf{n},i}^2 / \sigma_{\backslash i}^2)^{-1} \boldsymbol{a}_i^{\top} \right) = \frac{1}{N} \frac{\sum_{j=1}^N a_{ij}^2}{\sigma_{\mathsf{n},i}^2 / \sigma_{\backslash i}^2 + \sum_{j=1}^N a_{ij}^2} , \qquad (5.94)$$

which involves as visible only a scalar inverse. The unbiased a.k.a. site estimates then read

$$\tilde{\boldsymbol{x}}^{(i)} = \tilde{\boldsymbol{x}}^{\setminus i} + \frac{\sigma_{\setminus i}^2}{w} \boldsymbol{a}_i (\sigma_{\setminus i}^2 \boldsymbol{a}_i^\top \boldsymbol{a}_i + \sigma_{\mathsf{n},i}^2)^{-1} (y_i - \boldsymbol{a}_i^\top \tilde{\boldsymbol{x}}^{\setminus i}) , \qquad (5.95)$$

$$\sigma_{(i)}^2 = \sigma_{i}^2 (1/w - 1) , \qquad (5.96)$$

and the biased estimates can be retrieved by

$$\sigma^2 = (1/\sigma_{(i)}^2 + 1/\sigma_{\backslash i}^2)^{-1} , \qquad (5.97)$$

$$\boldsymbol{m} = \sigma^2 (\tilde{\boldsymbol{x}}^{(i)} / \sigma_{(i)}^2 + \tilde{\boldsymbol{x}}^{\setminus i} / \sigma_{\setminus i}^2)$$
 (5.98)

Noteworthy, the estimates \boldsymbol{m} and σ^2 need to proceed with the algorithm at the next factor and the site parameters $\tilde{\boldsymbol{x}}^{(i)}$ and $\sigma_{(i)}^2$ must be stored for when the factor $f_{y_i|\mathbf{x}}(\boldsymbol{x})$ is addressed the next time. Whenever a factor is accessed again, one can simply compute $\boldsymbol{\theta}^{\setminus \bullet} = \boldsymbol{\theta} - \boldsymbol{\theta}^{(\bullet)}$ ($\bullet \in \{s, 1, ..., M\}$) to obtain the corresponding input parameter. The entire algorithm is called seqVAMPobsand stated in Appendix E.11.

6. Numerical Results

This chapter comprises the results of numerical simulations that compare the algorithms presented in this thesis in different scenarios. Before the simulation results are shown, an overview over the algorithms that were presented and whose results will be shown in the following is given.

6.1 Overview over Recovery Algorithms

The algorithms that are presented in this thesis are summarized and categorized according to their derivations in Fig. 6.1. The first two columns show the different approaches in terms of factorization of the factor graph, which leads to a separation into message-passing-type, Turbo-type, and sequential algorithms. The last two columns specify, whether the algorithm is directly obtained from the processing in the factor graph, or if it has been adjusted by an estimation-theoretic bias compensation strategy. The arrows indicate which algorithms are related to each other. The background color highlights algorithms utilizing individual variances.

The message-passing-type algorithms that are considered here are Gaussian message passing (GMP), GMPpostExt, which uses post-processed extrinsic, instead of pre-processed extrinsic, and GMPpostExtLin, which is a mixture of both with post-processed extrinsic only at the linear (channel-constrained) estimation.

For the Turbo algorithms, there are two starting points: VAMP, which utilizes average variances as reliabilities and VAMPind with individual variances. The algorithms VAMPia, and VAMPii are obtained from VAMP by employing individual bias compensation after the channel-constrained estimation (instead of average bias compensation). The difference between VAMPia and VAMPii is that VAMPia uses an average variance in the signal-constrained estimation, whereas VAMPii keeps the individual ones.

The algorithm VAMPire results from VAMPind by inserting the estimation-theoretic bias compensation after the signal-constrained estimation.

seqVAMPind and seqVAMPire can be seen as sequential versions (sequentially process-



Figure 6.1: Categorization of the algorithms presented in this thesis. The arrows indicate which algorithms are related to each other; the background color highlights algorithms utilizing individual variances. VAMPii uses individual variances only in the signal-constrained estimation and is therefore highlighted differently.

Algorithm	Main Part	Appendix	no.	Reference(s)
AMP	Sec. 3.3.5	Sec. E.3	Alg. E.3	[DMM09, Mal11]
GAMP	Sec. 3.3.5	Sec. E.4	Alg. E.4	[Ran11]
GMP	Sec. 3.3.4	Sec. E.1	Alg. E.1	[KMS ⁺ 12]
GMPpostExt	Sec. 3.3.5	Sec. E.2	Alg. E.2	
GMPpostExtLin	Sec. 3.3.5			
seqVAMPind	Sec. 5.1.4	Sec. E.9	Alg. E.10	[OW05, BMPP20]
seqVAMPire	Sec. 5.1.5	Sec. E.10	Alg. E.11	[SF21]
seqVAMPobs	Sec. 5.2	Sec. E.11	Alg. E.12	
VAMP	Sec. 4.3.1	Sec. E.5	Alg. E.5	[RSF19]
VAMPia	Sec. 4.5.2	Sec. E.8	Alg. E.8	[SF22]
VAMPii	Sec. 4.5.2	Sec. E.8	Alg. E.9	[SF22]
VAMPind	Sec. 4.3.2	Sec. E.6	Alg. E.6	[FSRS16]
VAMPire	Sec. 4.5.1	Sec. E.7	Alg. E.7	[FSG20]

Table 6.1: References for the algorithms used in the simulations.

ing variables) of VAMPind, respectively VAMPire, because this sequential processing requires individual variances (so it is not a direct descendant from VAMP).

However, the algorithm that processes observations sequentially (seqVAMPobs) utilizes average variances and can therefore be interpreted as a sequential version (in terms of processing observations) of VAMP.

Additionally, all the algorithms used in the simulations are listed in Table 6.1 with references to the respective sections in the main part of the thesis, as well as section and algorithm number in the appendix. If the algorithms were published previous to this thesis, the reference is also given in the last column.

6.2 Simulation Setup

For the simulations, 10^5 realizations of sensing matrices A, as well as signal vectors x obeying the pdf $f_x(x)$, are drawn. If not stated otherwise, the elements of the sensing matrices are drawn i.i.d. Gaussian and the columns are normalized to unit ℓ_2 -norm, i.e., $\|\tilde{a}_j\|_2^2 = 1$ for $j \in$ $\{1, \ldots, N\}$. As signal model, the discrete ternary (DT) prior from (2.4) with non-zero signal alphabet $\mathcal{X} \in \{-1, +1\}$ is used. The Bernoulli-Gaussian prior is not considered, because in engineering applications often a specific (discrete) signal model is given [Pro00, SF14].

The computations for the signal-constrained estimation in case of the DT prior are given as follows, cf. also [SF16, Spa19]. Defining the input variables $\lambda = \tilde{x}_{\mathrm{s},j}/\tilde{\sigma}_{\mathrm{s},j}^2$ and $\Lambda = 1/\tilde{\sigma}_{\mathrm{s},j}^2$, the estimates read

$$m_{\mathrm{s},j} = \frac{\exp(\lambda - \Lambda/2) - \exp(-\lambda - \Lambda/2)}{\exp(\lambda - \Lambda/2) + \exp(-\lambda - \Lambda/2) + \frac{2(N-s)}{s}},$$
(6.1)

$$E_{x_j} \{ x_j^2 \mid \tilde{x}_{s,j} \} = \frac{\exp(\lambda - \Lambda/2) + \exp(-\lambda - \Lambda/2)}{(\lambda - \Lambda/2) + \exp(-\lambda - \Lambda/2)},$$
(6.2)

$$\exp(\lambda - \Lambda/2) + \exp(-\lambda - \Lambda/2) + \frac{2(N-s)}{s}, \qquad (0.2)$$

$$\sigma_{s,j}^{2} = E_{x_{j}} \left\{ x_{j}^{2} \mid \tilde{x}_{s,j} \right\} - m_{s,j}^{2} .$$
(6.3)

Note that this implementation is numerically more stable than the one used in [SF16], because the assembly of the two parameters λ and Λ in the exponential functions have a neutralizing effect, which makes the terms being computable for a wide range of values. If numerical instabilities occur, this makes itself noticeable already at the computation of $m_{s,j}$. Using the implementation above this only occurs, for very small variances $\tilde{\sigma}_{s,j}^2$ when $m_{s,j} \in \{+1, -1\}$ should actually be the result. In this case, the result can be obtained by setting $m_{s,j} = \text{sgn}(\lambda)$, where $\text{sgn}(\cdot)$ is the signum-function, i.e.,

$$m_{\mathbf{s},j} = \operatorname{sgn}(\lambda) \stackrel{\text{def}}{=} \begin{cases} +1 , & \text{if } \lambda > 0 , \\ 0 , & \text{if } \lambda = 0 , , \\ -1 , & \text{if } \lambda < 0 , \end{cases}$$
(6.4)

and the variance $\sigma_{\mathbf{s},j}^2$ to an arbitrary small positive value.

For a detailed explanation of the construction of the estimators for other priors, the reader may be referred to [Bir19].

Since the signal amplitude of the DT prior only attains discrete values, the estimates resulting from the algorithms are quantized before evaluation. The model (2.4) assumes to know the sparsity s of the signal, indeed the signal vectors are drawn such that they entail sparsity s. Therefore, this knowledge is applied by setting the smallest (in magnitude) N - s values in the final estimate m_s to zero. Other methods are discussed in [Spa19]. Subsequently, the remaining s entries of m_s are quantized to \mathcal{X} . Here, this can easily be done by utilizing the signum-function, i.e., $\hat{m}_{s,j} = \text{sgn}(m_{s,j})$ for $j \in \{1, \ldots, N\}$.

The noise vector is drawn i.i.d. zero-mean Gaussian with variance σ_n^2 , which results in a signal-to-noise ratio (SNR) defined as $10 \log_{10}(1/\sigma_n^2)$ in decibel (dB).

The performance is evaluated after a fixed number of iterations. As performance measure the symbol error rate (SER) is averaged over the realizations. The SER per realization is measured as

SER
$$(\hat{\boldsymbol{m}}_{s}, \boldsymbol{x}) = \frac{1}{N} |\{\hat{m}_{s,j} \neq x_{j} \mid j \in \{1, \dots, N\}\}|,$$
 (6.5)

where $|\cdot|$ denotes the cardinality of a set. If not stated otherwise, the performance is evaluated after 20 iterations. As pointed out in [RSF19] and [FSG20], the occurrence of negative variances in any of the recovery algorithm causes a drastic deterioration in performance, since a negative variance does not have a suitable meaning. Hence, this has to be mitigated. As suggested in [RSF19], the precisions, a.k.a. inverse variances are clipped to a suitable interval. The result of the signal-constrained estimation is clipped stricter than other variances, i.e., $1/\sigma_{s,j}^2 \in [10^{-8}, 10^8]$ and $1/\sigma_{c,j}^2, 1/\tilde{\sigma}_{s,j}^2 \in [10^{-12}, 10^{12}]$, since this resulted in stable simulations for a wide range of parameters. Note that clipping the precisions (instead of the variances) results in negative variances being transferred to the upper end of the interval, which is more useful since an estimation result corresponding to a negative variance (even if only produced by numerical inaccuracy) is usually not reliable. If one uses an implementation with precision ($\Lambda = 1/\sigma^2$) and scaled mean ($\lambda = m/\sigma^2$) and the precision is clipped, one has to make sure that the parameter λ is scaled respectively, in order not to change the underlying mean value m.

6.3 I.i.d. Gaussian Sensing Matrix

The standard case of sensing matrices for recovery in CS in the literature [Don06, EK12] is to draw the elements of the sensing matrices i.i.d. from a zero-mean Gaussian distribution with variance 1/M, i.e., $a_{ij} \sim \mathcal{N}(0, 1/M)$. The choice of the variance causes the ℓ_2 -norm of the columns to be roughly one, i.e., $\|\tilde{\boldsymbol{a}}_j\|_2^2 \approx 1$ for $j \in \{1, \ldots, N\}$. For the first scenario, such Gaussian matrices are considered, but with a scaling factor that ensures unit column norms, i.e., $\|\tilde{\boldsymbol{a}}_j\|_2^2 = 1$, cf. [Spa19].

In this scenario, the ability of the different algorithms to compress the signal x is investigated by displaying steady-state performance, which denotes the performance after convergence, over the observation dimension M.

To that end, all the algorithms listed above are simulated, except for VAMPia, VAMPii, and seqVAMPobs, because separate sections (Sec. 6.4 and Sec. 6.6) are spent for the discussion of their properties and they perform exactly as VAMP in the given scenario, which will be shown in the respective sections.

The first results in Fig. 6.2 show the performance over the observation dimension M for a signal dimension N = 100 and different values of sparsity s. The upper plot (Fig. 6.2(a)) compares the Turbo algorithm VAMP, with the message-passing-type algorithms GMP and GMPpostExt, as well as the approximate versions AMP [Mal11] and our implementation (Algorithm E.4) of GAMP [Ran11]. Throughout this chapter, the message-passing-type algorithms utilize schedules that iterate between the upward and downward messages by computing all messages of one type, e.g., downward, sequentially and then switching to the other type. In this scenario, AMP and VAMP show a very similar performance in terms of steady-state error, which is not surprising. Our implementation of the GAMP algorithm shows worse performance, especially for larger observation dimensions M, because the neglection of the offdiagonal elements in the LMMSE inverse (for details see Appendix E.4) gets harder to compensate for larger M and GAMP cannot benefit from the averaging effects of usual AMP. The message-passing-type algorithms perform similar to or only slightly worse than VAMP, but GMPpostExt shows some deterioration for small to mid-level values of M. The fact that the performance of GMPpostExt is generally very close to the one of GMP shows that the "postprocessing" of the extrinsics after the signal-constrained estimation does not differ much from the "pre-processing" case; at least for larger sparsities s.

In Fig. 6.2(b) the focus lies on the comparison of the Turbo algorithms with VAMP on the one hand, which uses average variances, and on the other hand VAMPind and VAMPire using individual variances. It is visible that the algorithms with individual variances can improve over VAMP in terms of compressibility and general steady-state performance, except for very small *M*. Moreover, VAMPire usually is slightly better than VAMPind, because of the estimation-theoretic bias compensation after the signal-constrained estimation that omits negative variances, thereby guaranteeing sensible variance values and improving performance compared to the treatment in VAMPind.

This difference is even more visible in Fig. 6.2(c), where the Turbo algorithms with individual variances are compared to its versions that process the signal components sequentially



(a) Comparison of the message-passing-type algorithms (GMP and GMPpostExt), VAMP as representative of a Turbo-type algorithm, and the AMP-type algorithms.



(b) Comparison of Turbo algorithms with average (VAMP) and individual variances (VAMPind, VAMPire).



(c) Comparison between parallel and sequential algorithms (processing of variable nodes, individual variances).

Figure 6.2: Performance over the observation dimension M. N = 100, i.i.d. Gaussian sensing matrices with columns normalized to unit ℓ_2 -norm, DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 20 iterations.



Figure 6.3: Performance over FLOPs per iteration. N = 100, i.i.d.Gaussian sensing matrices with columns normalized to unit ℓ_2 -norm, DT prior with sparsity s = 4, $M \in \{20, 40, 50, 80\}$ with increasing mark size, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 50 iterations.

(seqVAMPind and seqVAMPire). Here, the performance of the purely optimization-based sequential approach seqVAMPind performs visibly worse than the other algorithms. The reason therefor is that a negative variance affects the processing even more, since the result is directly broadcasted to all other variables, without a compensating effect from other variables. The sequential version of the estimation-theoretically improved algorithm (seqVAMPire) performs very close to the Turbo-type version VAMPire.

Before turning to more results that investigate the compressibility, let us have a quick look at the complexity of the algorithms. The computational complexity of the algorithms is compared by counting the floating-point operations (FLOPs) per iteration. The term FLOPs shall here comprise additions, subtractions, multiplications, and divisions. The counting is done as explained in [Spa19, Appendix E]. Note that for the signal-constrained estimation the FLOP count in [Spa19, Appendix E] involves the cardinality of the alphabet of non-zero entries, which is for the given prior (2.4) $|\mathcal{X}| = 2$.

The counting for the GMPpostExt algorithm is based on the notation that can be found in Appendix E.2. For GMP a version that utilizes scaled means and inverse variances to save the inversions in Lines 7 and 8 is used as basis for the complexity calculations. Of course, all algorithms could further be optimized, both in terms of storage and computational complexity, but since the requirements depend a lot on the architecture on which the algorithm runs, only a rough overview over the range that the algorithms span in terms of complexity is given.

The results are shown in Fig. 6.3 for the scenario that was investigated before with sparsity s = 4 and observation dimension $M \in \{20, 40, 50, 80\}$. On the horizontal axis, the counted FLOPs per iteration are shown; the vertical axis depicts, as usual, the SER. For increasing M the performance gets, as expected, better but in most algorithms also the FLOP count increases.

The AMP and the VAMP algorithm have roughly the same complexity, because the SVD implementation is used for the inverse in the LMMSE, as suggested in [RSF19, Algorithm 2]. The effort of computing the SVD is not added, because it can be computed once before the

execution of the algorithm. So, especially if a sensing matrix is used several times, the expense can be neglected. For the algorithms with individual variances (VAMPind and VAMPire), this simplification is not possible, because the decomposition does not lead to a diagonal to be inverted, as for the case with an average variance [Spa19, Appendix C]; this leads to a drastic increase in complexity compared to AMP and VAMP. The FLOP count for GAMP is slightly higher than the one of VAMP and AMP, because it involves individual variances, but not as high as VAMPind and VAMPire, because it does not compute the entire LMMSE inverse, but instead only the diagonal elements.

The computational complexity of the GMP-based algorithms (diamonds) spreads over a wide range, i.e., the FLOP counts depends very much on the implementation. The version with pre-processed extrinsic (GMP) has the highest FLOP count, because the expectations have to be computed per edge ($M \cdot N$). Post-processing the extrinsic after the channel-constrained estimation (GMPpostExtLin) already saves so much that the algorithm costs less than the Turbotype algorithms with individual variances. Note that this implementation performs exactly as GMP, because the updates are mathematically the same. The second simplification (GMPpost-Ext), which uses post-processed extrinsic also at the signal-constrained estimation, may lead to a slight loss in performance, but can further decrease the computational complexity. Hence, it makes sense to implement in message-passing-type algorithms for CS at least the channel-constrained estimation with post-processed extrinsic.

The computational complexity of the sequential approaches, seqVAMPind and seqVAM-Pire, does not scale well. Since the rank-one update operates on an $N \times N$ matrix, one such update is in the order of $\mathcal{O}(N^2)$ and as will be visible in simulations later on, at least N inner iterations are required so that the overall complexity is of order $\mathcal{O}(N^3)$. Nevertheless, the number of FLOPs does not depend on M so that these algorithms can compete with its Turbo-type versions for values of M close to N. A slight reduction in complexity may be possible by utilizing a Cholesky factorization as suggested in [SN08]; however the factorization dictates the succession of processing the variables, so the flexibility of most of the schedules to be proposed in Sec. 6.5 is lost. Recalling that this kind of processing is related to the (sequential) EP algorithm [Min01, SN11], cf. also Sec. 5.1.2.3, one may find other approaches for simplification in the respective literature. The idea of average EP (AEP) [DB18] and stochastic EP (SEP) [LHLT15] is for example to reduce the parameter space and thereby the storage complexity by in case of AEP averaging the natural parameters suitably such that only the averaged parameter has to be kept. This also simplifies the analysis of the procedure [DB18]. In the CS case this seems to remove to much of the necessary information, so that these techniques were not considered in this thesis. Consequently, these sequential algorithms are only suitable for small systems.

The version of VAMP that processes the observations sequentially (seqVAMPobs) is shown here for completeness; its FLOP count is in the order of GAMP, because the sequentialization increases the computational effort w.r.t. VAMP.

The results in Fig. 6.4 show the scenario with a signal dimension N = 250, i.e., the dimensionality of x is increased. The results are plotted for the same relative sparsities s/N as in the figures for N = 100. For this choice of signal dimension, the computation of the






(b) Comparison of Turbo algorithms with average and individual variances.



(c) Comparison between parallel and sequential algorithms (processing of variable nodes, individual variances).

Figure 6.4: Performance over the observation dimension M. N = 250, i.i.d. Gaussian sensing matrices with columns normalized to unit ℓ_2 -norm, DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 20 iterations.



Figure 6.5: Performance over M/N for the two considered values of N and $s/N \in \{0.02, 0.04, 0.08\}$). I.i.d. Gaussian sensing matrices with columns normalized to unit ℓ_2 -norm, DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 20 iterations.

GMP algorithm is already very tedious, especially for large M. Therefore, the results for the GMP algorithm are not shown. The performances of AMP and VAMP are now even closer, the degradation in GAMP stays, cf. Fig. 6.4(a).

In the middle (Fig. 6.4(b)) the comparison of the behavior of the Turbo-type algorithms with average and individual variances is shown. Here, the performances are also closer than before, i.e., the performance of VAMP improves compared to the algorithms with individual variances (VAMPind and VAMPire). The improvement of AMP and VAMP can be explained by the fact that for larger signal dimensions the performance of recovery algorithms in CS (with average variance) generally improves for growing N because the averaging pulls the (conditional) variances towards the MSE, which is the actual desired measure. For higher sparsities the algorithms with individual variances tend to have a worse compression capability than VAMP, but the steady-state error for large values of M is the same. This shows that the algorithms with individual variances are most suited for systems of small dimensions and very sparse vectors.

The comparison of parallel and sequential algorithms in Fig. 6.4(c) shows similar features as for N = 100. The performance of seqVAMPind is worst; the other algorithms differ only slightly in compressibility and steady-state error.

In Fig. 6.5 the effect of an increased signal dimension on the recovery performance for AMP, VAMP, VAMPind, and VAMPire is assessed. The main difference is that the relative compressibility gets better, i.e., the so-called *phase-transition region* [WV12] gets steeper, which is a well known fact in the literature [DMM09]. Indeed, the algorithms with average variances, AMP and VAMP, play to their strengths the larger the system size. Therefore, the algorithms with individual variances can improve mostly for smaller system sizes.

In Fig. 6.6 the convergence of the Turbo algorithms over the iterations is depicted; on the left for N = 100, M = 50 and on the right for N = 250, M = 125 with the respective sparsity values used above. On the right side, AMP and VAMP are compared to the algorithms with individual variances, VAMPind and VAMPire, for N = 250. AMP shows the slowest



Figure 6.6: Performance over iterations for the two considered values of N and several sparsities (the relative sparsity is the same, i.e., $s/N \in \{0.02, 0.04, 0.08\}$). M = N/2, i.i.d. Gaussian sensing matrices with columns normalized to unit ℓ_2 -norm, DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$.

convergence but reaches (as VAMP) better steady-state performance than the algorithms with individual variances for s = 20. For smaller sparsities (in the examples $s \in \{5, 10\}$), VAMPire supersedes the other algorithms in terms of steady-state performance, whereas VAMPind only achieves the performance of AMP and VAMP, despite using the more complex computation with individual variances. Generally, the algorithms with individual variances show the fastest convergence. This shows that the algorithms benefit from utilizing individual variances only if the treatment is performed correctly, as derived in the estimation-theoretic approach.

On the left, additionally to the algorithms on the right, also GMP and GMPpostExt are displayed in the scenario with N = 100. Both message-passing-type algorithms show slower convergence than AMP, because the fine-granular processing requires more steps to settle. The steady-state error is close to the one of AMP and VAMP, which is in general higher than on the right, which is due to less compressibility for smaller signal dimensions, as discussed beforehand. Again, one can see that the benefit of individual variances has most effect for the smaller system size of N = 100.

6.4 Model for Wireless Sensor Networks

Wireless sensor networks (WSNs) [YMG08] are an interesting application for CS [CSD⁺17, LCJ15, SD11, XDS13]. For several applications, it is necessary to gather information from a certain area. If the respective sensors are spread within this area, the direct access of the sensors might not be given or not desired. For this use, wireless sensor networks are deployed, in which the information of the sensors is collected at a so-called *fusion center*.

In order to reduce the power consumption at the sensors they only send if it is necessary, i.e., their activity may be sporadic [SD11]. Considering a large amount of sensors of whom only few are active in one time step, results in a sparse vector of high dimensionality. This brings

compressed sensing into play, which can simplify the hardware at the fusion center, e.g., by measuring with fewer antennas. Subsequently, the sparse vector needs to be recovered from the measurements. Therefore, the suitability of the discussed algorithms in such a scenario is investigated in the following.

The main effect that has to be addressed in this scenario is the fact that, since the sensors are distributed arbitrarily around the fusion center, the distances to the fusion center vary and together with it the receive power with which the signals of the sensors reach the fusion center. The assumption of unit- ℓ_2 -norm columns as in the previous section would require power control of the senders, i.e., the sensors need to send with a signal strength such that the receive power at the fusion center is equal for all sensors. In order to control this, additional communication is required, which is usually not desired since it reduces the battery lifetime of the sensors.

Alternatively, the effect may be interpreted as a non-uniform power distribution over the signal components x_1, \ldots, x_N that is modeled into the sensing matrix A by scaling the columns accordingly, which means that the power distribution is known to the receiver. The power distribution, given by p_i $(j \in \{1, \ldots, N\})$ is applied as follows

$$\boldsymbol{A} = c_{\boldsymbol{A}} \cdot \tilde{\boldsymbol{A}} \cdot \operatorname{diag}\left(\sqrt{p_{j}}\right) = \left[\tilde{\boldsymbol{a}}_{1}, \ldots, \tilde{\boldsymbol{a}}_{N}\right], \qquad (6.6)$$

where A consists of i.i.d. Gaussian entries and the columns are normalized to unit ℓ_2 -norm. For comparability, indeed to ensure that each sensing matrix carries the same energy, the factor c_A is adjusted such that $\|A\|_{\rm F}^2 = N$ by

$$c_{\boldsymbol{A}} = \frac{\sqrt{N}}{\left\|\tilde{\boldsymbol{A}} \cdot \operatorname{diag}(\sqrt{p_{j}})\right\|_{\mathrm{F}}}, \qquad (6.7)$$

i.e., the Frobenius norm (for the definition see Appendix A) is equal to the unit- ℓ_2 -norm case. In the simulations, the power distribution is specified by parameter $v \in (0, 1]$ in the form of

$$p_{\ell} = v^{(\ell-1)/(N-1)}, \qquad \ell \in \{1, \dots, N\},$$
(6.8)

and the scaling factor is assigned to the *j*th column by a random permutation $j = \pi(\ell)$. For v = 1, the unit- ℓ_2 -norm case is recovered; with decreasing v the variations in the scalings vary more and more. Furthermore, also the condition number (cf. Appendix A) of the sensing matrix increases. Noteworthy, the largest column power (before Frobenius normalization) is 1, the smallest is v.

An example for such a power distribution is plotted for v = 0.5 and N = 250 in Fig. 6.7. It shows also the effect of the Frobenius normalization for a sensing matrix of size 125×250 .

As pointed out in Sec. 4.4, VAMP is challenged by this scenario because the effects of the sensing matrix are treated on average, not individually. As alternative, the algorithms that utilize individual unbiasing instead are considered. These are VAMPind and VAMPire, which use individual variances throughout the entire iteration, as well as VAMPia and VAMPii, the algorithms with average variance in the channel-constrained estimation. This comparison



Figure 6.7: Power distribution (6.8) for N = 250, v = 0.5, i.i.d. Gaussian matrices before and after Frobenius normalization.



Figure 6.8: Performance over the iterations. N = 250, M = 125, i.i.d. Gaussian sensing matrices with power profile (6.8), DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$.

has first been considered in [SF22], the scenario with power distribution has slightly different earlier been used in [SF21].

Assuming simple binary transmission of the sensors in the WSN, translates into the DT prior (2.4) being a suitable choice for the signal model. In this section, we stick to signal dimension N = 250, observation dimension M = 125 and SNR $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$ if not stated otherwise.

First, the convergence over the iterations is depicted in Fig. 6.8. The left side compares VAMP and the variants for the usual case where the columns are normalized to unit ℓ_2 -norm, which coincides with the power distribution (6.8) with v = 1. One can see that the variants VAMPia and VAMPii perform exactly the same as VAMP. To understand why this is, regard

the end-to-end cascade \boldsymbol{K} from (4.44) with $\boldsymbol{\varPhi}_{\mathsf{x}} = \sigma_{\mathsf{x}}^2 \boldsymbol{I}_N$

$$\boldsymbol{K} = \sigma_{\mathsf{x}}^{2} \boldsymbol{A}^{\top} (\sigma_{\mathsf{x}}^{2} \boldsymbol{A} \boldsymbol{A}^{\top} + \sigma_{\mathsf{n}}^{2} \boldsymbol{I}_{M})^{-1} \boldsymbol{A} = (\boldsymbol{A}^{\top} \boldsymbol{A} + \frac{\sigma_{\mathsf{n}}^{2}}{\sigma_{\mathsf{x}}^{2}} \boldsymbol{I}_{N})^{-1} \boldsymbol{A}^{\top} \boldsymbol{A} .$$
(6.9)

First, note that the Gram matrix of \boldsymbol{A} can be written as

$$\boldsymbol{A}^{\top}\boldsymbol{A} = \begin{bmatrix} \sum_{i=1}^{M} a_{i1}^{2} & \sum_{i=1}^{M} a_{i1}a_{i2} & \dots & \sum_{i=1}^{M} a_{i1}a_{iN} \\ \sum_{i=1}^{M} a_{i1}a_{i2} & \sum_{i=1}^{M} a_{i2}^{2} & \dots & \sum_{i=1}^{M} a_{i2}a_{iN} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^{M} a_{i1}a_{iN} & \sum_{i=1}^{M} a_{i2}a_{iN} & \dots & \sum_{i=1}^{M} a_{iN}^{2} \end{bmatrix},$$
(6.10)

which shows that the diagonal contains the column energies (squared ℓ_2 -norms), which are in the case v = 1 all equal to 1, i.e., $\sum_{i=1}^{M} a_{ij}^2 = 1$ for $j \in \{1, \ldots, N\}$. Sensing matrices are usually chosen such that the measurements are as independent as possible from each other, in order to gain as much information as possible from each individual measurement. This means that the sensing matrices should have a low *coherence* [EK12, Def. 1.5]. Luckily, i.i.d. Gaussian matrices meet this requirement relatively well, if the dimensions are large enough [Don06]¹, leading to $\sum_{i=1}^{M} a_{ij}a_{ij'} \approx 0$ or at least $|\sum_{i=1}^{M} a_{ij}a_{ij'}| \ll 1$. With this knowledge, one finds that $M = \mathbf{A}^{\top}\mathbf{A} + \frac{\sigma_n^2}{\sigma_v^2}\mathbf{I}_N$ can in this case be written as

$$\boldsymbol{M} = \begin{bmatrix} 1 + \frac{\sigma_{n}^{2}}{\sigma_{x}^{2}} & \epsilon_{12} & \epsilon_{13} & \dots & \epsilon_{1N} \\ \epsilon_{12} & 1 + \frac{\sigma_{n}^{2}}{\sigma_{x}^{2}} & \epsilon_{23} & \dots & \epsilon_{2N} \\ \epsilon_{13} & \epsilon_{23} & 1 + \frac{\sigma_{n}^{2}}{\sigma_{x}^{2}} & \vdots \\ \vdots & \vdots & \ddots & \epsilon_{N-1N} \\ \epsilon_{1N} & \epsilon_{2N} & \dots & \epsilon_{N-1N} & 1 + \frac{\sigma_{n}^{2}}{\sigma_{x}^{2}} \end{bmatrix},$$
(6.11)

with $|\epsilon_{jj'}| \ll 1$ for $j' \neq j$, $j, j' \in \{1, \ldots, N\}$. Hence, one can deduce that in the inverse M^{-1} mostly the diagonal elements $1 + \frac{\sigma_2^2}{\sigma_x^2}$ contribute and, consequently, the diagonal entries of K are dominated by equal or close to equal scaling factors. Compensating the effects on average ignores the small variation but does not result in significant errors. Thus, the effect of average bias compensation is in this case very close to individual bias compensation. The other two algorithms, VAMPind and VAMPire, show a somewhat faster convergence due to their utilization of individual variances for the channel-constrained estimation.

By decreasing v, variations between the column energies appear due to the different scaling factors of the power distribution. So, the effect of individual and average bias compensation after the channel-constrained estimation differs. The convergence over iterations for v = 0.5 is shown on the right side of Fig. 6.8. Comparing to the left reveals that the steady-state performance is in general worse, which is due to the fact that the estimation problem overall gets harder under these conditions. It is clearly visible that the performance of VAMP degrades most for the considered algorithms. This can be explained by the stronger variations that appear in the diagonal entries of the end-to-end cascade K; the effect can no longer be compensated on average.

¹There are ways to generate matrices, that fulfill the property of low coherence even better [Spa19], but for our purposes, the Gaussian matrices suffice. For a definition of the coherence, see also Appendix A.



Figure 6.9: Example for variations in diagonal of K for an i.i.d. Gaussian matrix of size 125×250 , $\sigma_x^2 = 1/25$, $\sigma_n^2 = 0.2$, with (v = 0.5) and without (v = 1) power distribution. The diagonal entries are sorted in ascending order for the plot.

An example for the variations that result in the diagonal of K is displayed for an i.i.d. Gaussian matrix of size 125×250 , whose columns are normalized to unit ℓ_2 -norm in one case and compared to the case with applied power distribution (6.8) with v = 0.5. For the calculation of K, $\sigma_x^2 = 1/25$ is used, which corresponds to the initialization with s = 10, and $\sigma_n^2 = 0.2$, which represents an SNR of roughly $10 \log_{10}(1/\sigma_n^2) = 17$ dB. The diagonal entries are sorted in ascending order for better visibility in the plot. One can see that with employed power distribution, the variations between diagonal elements increase.

In order to see the entire picture, the simulations over the parameter v of the power distribution are shown in Fig. 6.10. One can see that the performance of VAMP degrades much more than the other algorithms when v decreases and, thus, more power variations are present. Considering the algorithms with individual bias compensation and average variance in the channel-constrained estimation, VAMPii is preferable over VAMPia, because VAMPii shows slightly better or equal steady-state error. This means that the usage of individual variances in the signal-constrained estimation can yield some insight to lower the steady-state error. For s = 20, the algorithms with individual variances are inferior to the ones with average variances when $v \rightarrow 1$, which coincides with the slightly worse compressibility of the algorithms for larger variances in Fig. 6.4(b). In this case, the individual variance algorithms show their benefit only for very small values of v, i.e., they are generally less impaired by large variations in the column powers, because of more insight to reliabilities of the channel-constrained estimates.

In Fig. 6.11 the steady-state error in terms of SER over the SNR, given by the inverse noise variance in dB, is considered. On the left, the results for v = 1 are shown, where the curves are very close. Noteworthy, the algorithms with individual variances are generally better for higher SNR, except for s = 20, where they generally deteriorate as already discussed. On the right, one can see the results for v = 0.5; it is again apparent that the performance of VAMP



Figure 6.10: Performance over parameter v of the power profile (6.8). N = 250, M = 125, i.i.d. Gaussian sensing matrices, DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 20 iterations.



Figure 6.11: Performance over SNR. N = 250, M = 125, i.i.d. Gaussian sensing matrices with power profile (6.8), DT prior for different sparsities, evaluation after 20 iterations.



Figure 6.12: Performance over FLOPs per iteration. N = 250, i.i.d. Gaussian sensing matrices with power profile (6.8), v = 0.5, DT prior with sparsity s = 5, $M \in \{50, 100, 125, 200\}$ with increasing mark size, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 20 iterations.

degrades more than all other algorithms with VAMPii being mostly the best.

Up to now, it has become clear that the algorithms VAMPia and VAMPii can very well compete with the algorithms with individual variances. In order to see the actual benefit over the strategies with individual variances, the computational complexity of the algorithms is compared by counting the floating-point operations (FLOPs) per iteration. The results are plotted in Fig. 6.12; the simulations and computations were executed for a sparsity s = 5 and observation dimension $M \in \{50, 100, 125, 200\}$. For the algorithms with average variance in the channel-constrained estimation, the SVD implementation is utilized for the inverse in the LMMSE [RSF19, Algorithm 2]. The benefit of the SVD implementation is obvious, as the individual variance algorithms need up to two orders of magnitude more FLOPs per iteration. Since the algorithms combine both, the low complexity of VAMP as well as a good performance due to the individual bias compensation after the channel-constrained estimation, they are the algorithms to be preferred in such a scenario. Since VAMPii is due to the usage of individual variances in the signal-constrained estimation usually slightly better than VAMPia, it is preferable over VAMPia.

VAMP is usually promoted by its stability in ill-conditioned sensing scenarios, where the sensing matrix has high condition number κ [RSF19], for the definition of the condition number in the form that it is used here, see Appendix A. The comparison of the algorithms in such a scenario², see Fig. 6.13, shows that the algorithms utilizing individual variances can improve over VAMP, as well as the algorithms with individual bias compensation and average variance at the channel-constrained estimation. Noteworthy, VAMPii supersedes partially even the algorithms with individual variances, which is remarkable considering the much lower complexity. The explanation for that is the average variance in the channel-constrained estimation, which brings more stability into the algorithm, thereby improving the overall convergence.

²The sensing matrix A is generated as in [RSF19, Sec. VI.A], which is explained in Appendix A.



Figure 6.13: Performance over condition number κ of sensing matrix **A**. N = 250, M = 125, DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 20 iterations.

6.5 Schedules for the Sequential Processing of Variables

In this section, simulation results for the algorithms that utilize a sequential processing of the signal components and were introduced in Sec. 5.1.4 and 5.1.5 are presented. They are compared to the Turbo-type algorithms VAMP, VAMPind, and VAMPire, since they are closest to them in terms of estimation strategy with VAMP being furthest apart, because it employs average variances as reliabilities for the estimations.

The scenario is the same as in the section above, i.e., i.i.d. Gaussian sensing matrices of size 125×250 are drawn and the power distribution (6.8) is applied. As signal model, the DT prior (2.4) with s = 5, if not stated otherwise, is assumed. The signal-to-noise ratio is set to $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$ throughout the simulations. Similar results have been shown in [SF21].

The simulations in Sec. 6.3 have already shown that seqVAMPire usually performs close to VAMPire and improves over seqVAMPind in the standard setting utilizing i.i.d. Gaussian matrices with columns normalized to unit ℓ_2 -norm. The comparison in the section at hand focuses mostly on the difference in convergence speed, but before turning to that, the performance over the parameter v of the power distribution (6.8) is considered, in order to get a picture of the performance. The results are shown in Fig. 6.14. For small values of v, the performance of the sequential algorithms is very close to the Turbo-type algorithms with individual variances from which they are derived. For $v \rightarrow 1$, the performance tends to get worse than the parallel algorithms, which is especially visible for the larger sparsity value s = 20. The sequential algorithms utilize the $\sigma^2 \nearrow$ schedule, which will be explained below. Noteworthy, this schedule was also used in the simulations over the observation dimensions in Sec. 6.3.

In the following, the simulations are restricted to s = 5 and $v \in \{0.5, 1\}$.

Opposing the Turbo-type algorithms, where each processing step is prescribed, the sequential processing of the signal components enables a degree of freedom in choosing among the signal components the one to be processed in the next step. This opens the possibility for



Figure 6.14: Performance over parameter v of the power profile (6.8). N = 250, M = 125, i.i.d. Gaussian sensing matrices, DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 20 iterations. Comparison between Turbo-type and sequential (processing variables) algorithms. The sequential algorithms utilize the $\sigma^2 \nearrow$ schedule, which will be explained below.

a variety of schedules. The schedules that are discussed here were first introduced in [SF21].

In order to be able to properly define the schedules, a distinction between *inner iterations* and *outer iterations* is employed. One *inner iteration* comprises the computation steps that are described in Fig. 5.2, i.e., the signal-constrained estimation of one signal component x_j and the subsequent broadcasting of the result to the other signal components. An *outer iteration* describes the processing of N inner iterations, i.e., the number of inner iterations that are necessary to process each signal component at least once. These outer iterations are used to compare to the Turbo-type algorithms, which process all signal components together.

At first, schedules that sweep entirely through the signal components are considered, so that no component is left out and none is processed twice in one outer iteration. Therefore, one can directly compare over the outer iterations. As schedules, a random choice and two deterministic choices are introduced. The random choice (random) draws the order to be processed before each outer iteration, i.e., permutes the sequence $[1, \ldots, N]$ and runs through it. For the deterministic choices, the variances $\sigma_{c,j}^2$ ($j \in \{1, \ldots, N\}$) that represent the reliability in the channel-constrained estimation of the respective signal component need to be considered. The variances are either sorted in ascending ($\sigma^2 \nearrow$), or in descending ($\sigma^2 \searrow$) order and the order is kept throughout the outer iteration.

The comparison of these schedules over the iteration is depicted in Fig. 6.15, on the left for v = 1, on the right for v = 0.5. In the unit-column-norm case (v = 1), there is almost no difference between the schedules. Nevertheless, it is visible that the convergence of the sequential algorithms is slightly faster than the Turbo-type algorithms in terms of outer iterations. Furthermore, in contrast to seqVAMPire, whose steady-state error ends up at the same level as VAMPire, the sequential version of VAMPind (called seqVAMPind) results in a steady-state error worse than all the compared Turbo-type algorithms. This shows that the estimation-theoretic view on bias compensation becomes even more important, when the view on the



Figure 6.15: Comparison of schedules that encompass all signal elements per outer iteration. N = 250, M = 125, i.i.d. Gaussian sensing matrices with power distribution (6.8), DT prior with sparsity s = 5, $10 \log_{10}(1/\sigma_n^2) = 16$ dB.

estimation is even more detailed as in this sequential processing case. The faster convergence of the sequential algorithms can be explained by the fact that the knowledge gained by processing a single variable is immediately fed back to the other variables, so that within an outer iteration already more knowledge can be gained as when processing all variables together.

For v = 0.5, one can see a difference between the schedules; especially after the first outer iteration. Later the schedules tend to the same level in terms of SER. The schedules that exceeds the other two in the first few iterations is $\sigma^2 \nearrow$, i.e., the processing in order of ascending variances. The schedule $\sigma^2 \searrow$ is worst and the random schedule is between the other two.

The reason for this order is as follows. Since the variance $\sigma_{c,j}^2$ represents the reliability in the estimates, starting with low variances (as in $\sigma^2 \nearrow$) means to start with the ones one is most sure about, which leads to more confidence in subsequent estimations of signal components and thereby to faster convergence. The second deterministic schedule $\sigma^2 \searrow$ works oppositely and therefore performs worst, although still being better than the Turbo-type algorithms.

The schedules examined so far followed the strategy of choosing the order once before an outer iteration and kept this order throughout the outer iteration. However, it might be worthwhile to reconsider and perhaps change the order of processing within one outer iteration. For that purpose, the index set $\mathcal{J} \subseteq \{1, \ldots, N\}$ of cardinality $|\mathcal{J}| = T$ is introduced, which specifies the T variables that are processed before reconsidering the order.

There are plenty of ways to choose the index set \mathcal{J} . Again two deterministic and one random strategy are examined. The deterministic approaches follow the deterministic choices introduced earlier, i.e., one picks the positions corresponding to the T smallest variances, the other goes for the T highest variances.

For the random approach, a strategy that prefers high variances is chosen. The motivation therefor is that the recovery algorithm needs to make sure that the reliability in the estimation increases over the iterations, i.e., the variances need to decrease. The variance of a signal component generally decreases most, when the signal component itself is processed in the signal-constrained estimation, because the signal-constrained (non-linear) estimation yields more accurate estimates than the channel-constrained (linear) one. Signal components with high variance $\sigma_{c,j}^2$ therefore indicate that they require to be processed via the signal-constrained estimation.

In order to focus on high variances, the positions are drawn according to a respective probability distribution. The distribution is created such that the probability for a position j to be drawn is proportional to its corresponding variance $\sigma_{c,j}^2$. Since unique positions in the set \mathcal{J} are necessary, it is built successively. Assume that $\ell < T$ positions have already been drawn, i.e., define $\mathcal{J}^{(\ell)} = \{j_1, \ldots, j_\ell\}$. Then, the probability for the other positions to be drawn is given by

$$\Pr\{\mathbf{j}_{\ell+1} = j\} = \frac{\sigma_{\mathbf{c},j}^2}{\sum_{j'=1,j'\notin\mathcal{J}^{(\ell)}}^N \sigma_{\mathbf{c},j'}^2}, \qquad j \in \{1, \dots, N\} \setminus \mathcal{J}^{(\ell)}.$$
(6.12)

This procedure is denoted by "rv".

In order to make sure that all the positions are processed equally often, one can track the number of processings conducted for each position j with a counter c_j . For later reference, the set of positions with minimal counters is denoted by

$$\mathcal{J}_{c} = \{ j \mid c_{j} \le c_{j'} \; \forall j, j' \in \{1, \dots, N\} \} .$$
(6.13)

This procedure is especially necessary for the procedure minv, which chooses always the positions with lowest variances (see below), because, as explained before, processing a signal component leads to a decrease of the respective variance, which means that this procedure is likely to pick signal components that it has chosen in the first iteration over and over again, without ever processing the high-variance components which actually needs to be processed. When considering the deterministic strategies, we will additionally always keep track of the processing by the counter.

Furthermore, one can deliberately switch between the strategies before choosing the next \mathcal{J} . When the index set \mathcal{J} is specified, the order of processing the set can be chosen as one of the previously considered strategies, namely, $\sigma^2 \nearrow$, $\sigma^2 \searrow$, and random.

In the simulations, we restrict to the extreme case, opposing the strategies shown before, i.e., T = 1 so that after each inner iteration, the processing order is reconsidered. In this case, there is of course no need to order the set, which consists of one position only, i.e., $\mathcal{J} = \{j\}$. In summary, the following strategies are considered

- maxv: choose $j = \operatorname{argmax}_{j' \in \mathcal{J}_c} \sigma_{c,j'}^2$ with \mathcal{J}_c from (6.13),
- minv: choose $j = \operatorname{argmin}_{j' \in \mathcal{J}_c} \sigma_{c,j'}^2$ with \mathcal{J}_c from (6.13),
- rv: draw position from $\{1, \ldots, N\}$ randomly with probability given by (6.12),
- minv + maxv,
- \blacksquare minv + rv,



Figure 6.16: Comparison of schedules with T = 1 over inner iterations. N = 250, M = 125, power distribution (6.8) with v = 0.5, DT prior with s = 5, $10 \log_{10}(1/\sigma_n^2) = 16$ dB.

where the plus means that the strategy is switched after every processed signal component. Because the schedules show differences, especially when the power distribution is applied, the focus lies now on simulations with v = 0.5.

The simulation results are shown in Fig. 6.16 over the inner iterations, where on the left side the individual variance version of VAMP (VAMPind) is compared to its sequential version (seqVAMPind) and on the right the respective algorithms improved by estimation-theoretic bias compensation (VAMPire and seqVAMPire) are depicted. The Turbo-type algorithms are evaluated at multiples of N = 250, because they process all N signal components at once. As can be seen, the general behavior of the schedules is for both sequential algorithms, seqVAMPind and seqVAMPire the same, only the steady-state performance differs, which is a property of the algorithm, as already discussed.

For better visibility, the results of estimation-theoretically improved algorithms is again shown in Fig. 6.17, so that one can describe and interpret the behavior of the schedules. Since the counter is used for the deterministic strategies, they process all signal components once between multiples of N, i.e., per outer iteration. When choosing always the largest variance (maxv) from the (per outer iteration) not yet processed signal components, a fast convergence at the beginning of the outer iteration results, because the respective positions require processing the most. Towards the end of the loop, the performance flattens out, because the positions that need less attention are addressed. This leads to a staircase behavior. The staircase is better visible in the opposite case (minv), which chooses the position with lowest variance (among the per outer iteration not yet processed signal components). The behavior is contrary to maxv in the sense that the performance stays equal at the beginning of the outer iteration and shows a drop at the end. Noteworthy, with this drop it gets below the maxv strategy, which is consistent with the results from Fig. 6.15, where the strategy going for low variances first $(\sigma^2 \nearrow)$ improved over the opposite strategy $(\sigma^2 \searrow)$.

The probabilistic approach does not make use of the counter and therefore converges



Figure 6.17: Comparison of schedules with T = 1 over inner iterations for the estimation-theoretically adapted algorithms. N = 250, M = 125, power distribution (6.8) with v = 0.5, DT prior with sparsity s = 5, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$.

smoothly, but the convergence is not as steep as for maxv at the beginning.

With the switching between opposing strategies, it is possible to combine the advantages of the respective strategies and yield both fast convergence and well performance at convergence; in the given comparison, the combination of rv and minv is the preferable schedule, since it combines the good convergence of the minv schedule with the smoothness (no steps) of the rv schedule.

6.6 Variable Noise Scenario

In this last scenario, the case of non-i.i.d. noise with independent but individual noise variances as introduced in Sec. 5.2 is considered, i.e.,

$$\mathbf{n} \sim \mathcal{N}\left(\mathbf{0}, \operatorname{diag}\left(\sigma_{\mathbf{n},i}^{2}\right)\right)$$
 (6.14)

This scenario was used to motivate the sequential processing of observations, which contrasts VAMP in the sense that VAMP considers all observations at once, i.e., parallel. The algorithm seqVAMPobs is therefore compared to two versions of VAMP. Naturally, the known noise variances $\sigma_{n,1}^2, \ldots, \sigma_{n,M}^2$ can simply be considered in the VAMP algorithm by computing the channel-constrained estimation by

$$\boldsymbol{m}_{\mathrm{c}} = \tilde{\boldsymbol{x}}_{\mathrm{c}} + \sigma_{\mathrm{x}}^{2} \boldsymbol{A}^{\top} \left(\sigma_{\mathrm{x}}^{2} \boldsymbol{A} \boldsymbol{A}^{\top} + \mathrm{diag} \left(\sigma_{\mathrm{n},i}^{2} \right) \right)^{-1} \left(\boldsymbol{y} - \boldsymbol{A} \tilde{\boldsymbol{x}}_{\mathrm{c}} \right) , \qquad (6.15)$$

with respective error covariance matrix

$$\boldsymbol{\Phi}_{c} = \left(\boldsymbol{A}^{\top} \operatorname{diag}\left(1/\sigma_{n,i}^{2}\right)\boldsymbol{A} + \frac{1}{\sigma_{x}^{2}}\boldsymbol{I}_{N}\right)^{-1} .$$
(6.16)

Note that our implementation uses the direct computation of the unbiased variable \tilde{x}_s by employing the scaling matrix W appropriately, so that Φ_c does not have to be computed directly and, thus, the computation of the $N \times N$ inverse can be spared. The other variant of VAMP averages the noise variances $\sigma_{\bar{n}}^2 = \frac{1}{N} \sum_{i=1}^M \sigma_{n,i}^2$ and uses this single value for the noise variance as usual.

If not stated otherwise, a scenario with N = 250, M = 125, the DT prior with various sparsities, and $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$ is simulated.

The different noise variances are obtained by scaling the noise variance σ_n^2 , which is specified by the SNR, with factors that are drawn uniformly from the interval [1 - b/2, 1 + b/2], i.e., the random variable ς_n^2 representing the noise variance σ_n^2 , is uniformly distributed around its mean value σ_n^2 with width $b\sigma_n^2$ by

$$\varsigma_{n}^{2} \sim \mathcal{U}\left(\left[\sigma_{n}^{2}(1-b/2), \sigma_{n}^{2}(1+b/2)\right]\right).$$
(6.17)

Since negative variances are not possible, this naturally yields a maximum of $b \le 2$ on the width parameter b of the uniform distribution. For b = 0, the i.i.d.-noise case is retrieved.

First, the standard case with i.i.d. noise is considered. As before, i.i.d. Gaussian sensing matrices with columns normalized to unit ℓ_2 -norm are used. VAMP and seqVAMPobs are compared over the observation dimension M in Fig. 6.18(a). The simulation shows that in this standard case, both algorithms perform exactly the same, which is no surprise since they are based on the same framework.

In contrast, the same plot is shown for non-i.i.d. noise with noise variance uniformly distributed around σ_n^2 with parameter b = 1 in Fig. 6.18(b), where one can see that seqVAM-Pobs improves over VAMP both, in terms of compressibility and resulting steady-state error for $M \rightarrow N$. Interestingly, VAMP performs better, when the average of the noise variances is used in the algorithms, but still not getting the performance of seqVAMPobs. The worse performance of VAMP with individual noise variances can be explained by the fact that the scaled identity matrix of using a single noise variance has an averaging effect in the processing, thereby making the inverse of the LMMSE much more stable, especially compared to the case when small values occur on the diagonal of the matrix that has to be inverted. The sequential approach benefits over these approaches because it can make use of the knowledge about the individual noise variances without getting impairments by the less stable matrix inverse.

In order to see the impact of the noise variance model on the behavior of the algorithms, they are compared over the value of b that parameterizes the width of the uniform distribution. Recall that for b = 0 the i.i.d.-noise case is recovered and the maximum is naturally given by b = 2. The results are shown in Fig. 6.19 for M = 125. On the left side, one can see that since b = 0 recovers the i.i.d. case, VAMP and seqVAMPobs perform the same. Towards the right, the performance of VAMP that utilizes the individual noise variances decreases, due to the



(a) Usual case with i.i.d. noise.





Figure 6.18: Comparison of VAMP and seqVAMPobs over the observation dimension M. N = 250, i.i.d. Gaussian sensing matrices with columns normalized to unit ℓ_2 -norm, DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 20 iterations.



Figure 6.19: Performance over width parameter *b* of the uniform distribution (6.17) for the noise variance. N = 250, M = 125, i.i.d. Gaussian sensing matrices with columns normalized to unit ℓ_2 -norm, DT prior, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$, evaluation after 20 iterations.

missing averaging effect inside the LMMSE inverse, as already mentioned. When the average of the noise variances is used, the performance of VAMP stays almost constant. In contrast, the sequential version seqVAMPobs can actually improve performance w.r.t. the i.i.d.-noise case, because it handles the observations individually, which only gets less stable for $b \approx 2$.

The sequential processing of the observations enables the use of certain schedules. As in the previous section, a random schedule is compared to two deterministic schedules with opposite ordering. For the deterministic sorting, a measure is needed that characterizes the reliability of the respective factor. In case of the channel-constrained factors $f_{y_i|x}(x)$, this is the ℓ_2 -norm of the row a_i , which tells us how much the signal is amplified by the measurement, divided by the respective noise variance $\sigma_{n,i}^2$. Additionally, there is the signal-constrained factor $f_x(x)$, which is characterized by $\sigma_x^2 = s/N$. Since these measures do not change throughout the iterations of the algorithm, the sorting is done once at the beginning and stays until the end. The random approach permutes the M + 1 possibilities before each iteration through the factors.

The sequential approach for the signal components from Sec. 5.1 showed that the processing of the signal-constrained estimation is very crucial for the progress of the algorithm. One may therefore emphasize on the factor $f_x(x)$, when designing schedules for the sequential processing of observations. In order to stay somewhat fair to VAMP, an iteration for the sequential algorithm seqVAMPobs is defined to be the processing of exactly M+1 factors. In the simulations below, versions of the previously introduced schedules are shown that process the factor $f_x(x)$ every second time; so after one channel-constrained factor $f_{y_i|x}(x)$ always $f_x(x)$ is processed, before the next channel-constrained factor is chosen. All in all, the following schedules are considered (the index s is used to denote the signal factor)

- \nearrow : sort $\{\|\boldsymbol{a}_i\|_2^2/\sigma_{\mathsf{n},i}^2 \mid i \in \{1, \ldots, M\}\} \cup \{s/N\}$ in ascending order
- \searrow : sort $\{\|\boldsymbol{a}_i\|_2^2/\sigma_{\mathsf{n},i}^2 \mid i \in \{1, \ldots, M\}\} \cup \{s/N\}$ in descending order



Figure 6.20: Performance over iterations with comparison of schedules. N = 250, M = 125, i.i.d. Gaussian sensing matrices with columns normalized to unit ℓ_2 -norm, non-i.i.d. noise with uniformly distributed variance according to (6.17) with b = 1, DT prior with sparsity s = 10, $10 \log_{10}(1/\sigma_n^2) = 16 \text{ dB}$.

- π : random permutation of $[s, 1, \ldots, M]$ before each iteration
- r^{n} : sort $\|\boldsymbol{a}_{i}\|_{2}^{2}/\sigma_{\mathsf{n},i}^{2} \nearrow$ and interleave indices with s
- Sum sort $\|\boldsymbol{a}_i\|_2^2/\sigma_{\mathsf{n},i}^2$ and interleave indices with s
- π_s : permutation of $[1, \ldots, M]$ interleaved with s, yielding $[s, \pi(1), s, \pi(2), \ldots, s, \pi(M)]$

The simulations shown so far used the scheduling that starts with the most reliable observation and proceeds towards the least reliable one, i.e., \nearrow .

The comparison of the schedules over the iterations are shown in Fig. 6.20 for b = 1. One can see that VAMP converges faster, but has higher steady-state error, as expected from the previous results. The schedules considering every factor once per iteration are shown with filled markers and the schedules emphasizing on the signal-constrained factor $f_x(x)$ are shown lighter and with non-filled markers. The sequential algorithm is evaluated each time when the schedule is run through entirely. Because the schedules emphasizing on s process 2M factors, instead of M + 1 as the others, they are only plotted every second iteration.

Except for the first iteration the two deterministic procedures yield very close results; also the emphasis on the signal-constrained estimation does not gain much, if at all. The approach with ascending variances yields the most benefit in the first iteration, which means the most insight from reliable observations is obtained at the beginning. The random and the descending variance schedule start similar in the first iteration, but the random approach, converges a little slower; a deterministic approach seems to be the means of choice. The steady-state error of the different schedules is the same.

7. Conclusions and Future Work

The thesis at hand considered recovery algorithms for compressed sensing (CS), i.e., the reconstruction of high-dimensional signals from linearly compressed measurements in noise. Turbo-type and message-passing-type algorithms were derived from the framework of minimizing the Kullback–Leibler divergence under expectation constraints tailored to the respective factor graph representation, which shows the common basis for both types of algorithms. Furthermore, sequential algorithms were developed under the same framework and comparisons of all algorithms were documented via numerical simulations.

7.1 Conclusions

In Chapter 2, the CS model was introduced and the reconstruction task was characterized by two different kinds of estimations that resulted from independently treating the two basic constraints, which are imposed by the problem. These constraints were distinguished as the *channel-constrained* part, resulting from the linear compression with known (sensing) matrix, and the *signal-constrained* part, which is defined by the inherent structure of the signal, i.e., its statistical properties (in form of a probability density function) and especially its *sparsity*. The exact form of the channel-constrained estimation depended on the considered factorizations of the problem. Different representations of the problem were introduced in the form of factor graphs and the approach of minimizing the Kullback–Leibler divergence was motivated.

The most fine-granular representation of the problem led in Chapter 3 to message-passing algorithms, which simplify the estimation to the treatment of scalars (instead of vectors). This simplification of the estimation required extensive exchange (messages) between the many factors and variables. In the exchange between the estimations, the well-known concept of *extrinsic* was found to be crucial for the processing of the algorithm. In its usual form, the resulting algorithm suffer in the CS scenario from its fine-granularity because of the high number of messages that have to be processed individually [KMS⁺12]. The expense could be reduced by exploiting the benefits of *post-processing* the extrinsic, which requires less estimations because the estimation becomes edge-independent, instead of *pre-processing* it. In

Chapter 3 it was shown that post-processing and pre-processing the extrinsic is mathematically the same for the channel-constrained estimation due to its linearity. Numerical results in Chapter 6 pointed out that the usage of post-processed extrinsic for the signal-constrained estimation do not yield a significant drawback.

The Turbo approach of Chapter 4 contrasted the message-passing approach as it denoted the factorization with least possible fragmentation, i.e., signal- and channel-constrained part oppose each other directly, resulting in two separated estimations. The derivation, again, motivated the use of extrinsic and led to the famous VAMP algorithm [RSF19] and a version of it utilizing individual variances. The estimation-theoretic analysis of the extrinsic computation revealed the connection to *bias compensation*. The optimization-based derivation linked the treatment of individual variances directly with individual bias compensation, as well as average variances and average bias compensation. For the bias compensation after the channelconstrained estimation it was proposed to cut the tie and combine individual bias compensation with the use of average variances. The algorithm with individual variances, which was directly derived from the optimization framework showed the drawback of not preventing the computation of negative variances. With the estimation-theoretic approach an adapted update was derived for the processing after the signal-constrained estimation, which guarantees positive unbiased variances. Utilizing these methods, several variants of the so far existing algorithms were proposed.

Both, the message-passing and the Turbo approach, elucidated the well-known importance of the computation of *extrinsic* in the exchange of parameters between different estimations.

The constrained optimization framework was also used to derive algorithms that process either signal components or measurements in a sequential manner in Chapter 5. Since the sequential processing of variables requires individual variances, the estimation-theoretic adaptions could be used to enhance the algorithm resulting directly from the optimization-based derivation.

In summary, four different factorizations of the posterior representing the CS problem were considered. One framework was used to provide derivations of algorithms for the different factorizations. Estimation-theoretical bias compensation was used to improve the derived algorithms.

The simulation results in Chapter 6 showed that the estimation-theoretically improved algorithms could generally beat the algorithms directly derived from the optimization-based approach in terms of compressibility and steady-state error. The model for wireless sensor networks showed a drawback of VAMP for unequally scaled sensing matrix columns, which could be overcome by individual bias compensation after the channel-constrained estimation without increasing the computational complexity significantly. The sequential algorithms have the freedom of varying the schedule of processing either the signal components or the observations. For the algorithms that process signal components sequentially, the possible schedules were examined in detail. It could be shown by numerical simulations that the sequential processing of measurements can defeat VAMP in scenarios with non-i.i.d. noise.

All in all, the alternative approaches (leading to sequential algorithms) and the estimationtheoretic enhancements turned out to be able to improve VAMP in certain specific scenarios. Generally, on the one hand, the algorithms with individual variances play to their strengths for small signal dimensions N, where the advance to VAMP is stronger and the restraint of their high complexity does not (yet) kick in. On the other hand, one could see that the average variance strategy of VAMP inherits useful stability properties on both sides of noise and prior variances. Nevertheless, the usage of average variance can be beneficially combined with individual bias compensation.

7.2 Future Work

In order to focus on the estimation-theoretic concepts within the algorithms, the thesis was restricted to relatively basic assumptions on the scenario, such as independent and identically distributed signal components with known sparsity and measurements obtained by a known sensing matrix in Gaussian noise. Possible future work could therefore encompass to test and adapt the proposed algorithms in less restrictive scenarios, closer to real-world applications; a keyword being *block sparsity* [BCDH10] or multidimensional priors [BG19, Bir19] that model dependencies between the signal components. It would also be interesting to see how robust the algorithms are when there is a mismatch to the prior assumption [Ver10, VS13]. Additionally, also the performance in the multiple measurement vector (MMV) problem [EK12, CREK05] or the distributed compressed sensing (DCS) [SBW⁺05] problem could be examined. Both problems work with multiple measurements, where DCS utilizes different sensing matrices per measurement vector; while the MMV problem considers the same sensing matrix in each measurement. The multiple measurements may be exploited to gain more knowledge from temporal or spatial correlations.

The algorithm, which processes measurements sequentially has not been investigated very thoroughly yet. Possible applications could be PAR reduction schemes for OFDM, which have to handle varying noise variances [KPN⁺14, LHD20, LHL⁺19]. Moreover, the trade-off between computing the channel-constrained estimate based on a single observation (as in the sequential algorithm), some part of the observation vector, or the entire observation vector (as in VAMP) might be worthwhile to consider, both in terms of performance and complexity.

Furthermore, the estimation-theoretic adaptions in the bias compensation strategies might be beneficially applicable in inference problems outside the CS scenario, as long as they inherit minimum mean-squared error estimators.

A. Vector Spaces

In this appendix, the notion of a vector space is discussed because the entire thesis is based on the treatment of vectors and matrices. Furthermore, the notion is used in Chapter 4 for a geometrical interpretation of estimation and bias compensation. The appendix comprises a summary of the necessary definitions, e.g., norms, and necessary properties of vectors and matrices.

A vector space \mathcal{V} over a certain (scalar) field \mathbb{F} is a set of objects that can be linearly combined, i.e., scaled and added, without leaving the set. Mathematically this means the set is closed under addition and multiplication with elements of the corresponding scalar field, i.e., with $a, b \in \mathbb{F}$, it holds for all $x, y \in \mathcal{V}$ that $ax + by \in \mathcal{V}$. The scalar field is in this thesis always the set of real numbers \mathbb{R} . For mathematical details and more properties, the reader is referred to [HJ09, Sec. 0.1.2].

A.1 Vectors

In this thesis, the objects are mostly vectors, which are represented by a list of entries, e.g.,

$$\boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0, 1 \end{bmatrix}^\top, \qquad \boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} -0.23 \\ 0.76 \end{bmatrix},$$
 (A.1)

as in Ex. 2.1. Note that vectors are represented as column vectors. Other examples for the objects are functions and thereby random variables, which is discussed in detail in Sec. 4.4.2.

The objects *span* the space. In the above example, the space is of dimension two, because there are two entries in the vectors and the vectors are linearly independent (one cannot get the other vector by simply scaling one of them). The same space is spanned by \boldsymbol{x} and $[1, 0]^{\top}$ as well, which together form an *orthonormal basis*. This means that the vectors are *orthogonal*, i.e., the *inner product* is zero. The inner product of vectors $\boldsymbol{x}, \boldsymbol{y} \in \mathcal{V}$ of dimension N is here given by the sum of the component-wise multiplied entries

$$\boldsymbol{x}^{\top}\boldsymbol{y} = \sum_{j=1}^{N} x_j \cdot y_j \;. \tag{A.2}$$

$$\boldsymbol{y} = \begin{bmatrix} -0.23, \ 0.76 \end{bmatrix}^{\top} \\ \|\boldsymbol{y}\|_{2} = 0.79 \\ \|\boldsymbol{y}\|_{2} = 0.79 \\ \|\boldsymbol{y}\|_{2} = 0.79 \\ \|\boldsymbol{y}\|_{1} = 1 \\ \|\boldsymbol{y}\|_{1} = 1 \\ \|\boldsymbol{y}\|_{1} = 1 \\ \|\boldsymbol{y}\|_{1} = |-0.23| + |0.76| = 0.99 \\ \|\boldsymbol{x}\|_{0} = 1 \\ \|\boldsymbol{y}\|_{0} = 2 \\ \end{bmatrix}$$

Figure A.1: Depiction of x and y from Ex. 2.1 in the space spanned by x and $[1, 0]^{\top}$.

In the example, it is $\boldsymbol{x}^{\top}\boldsymbol{y} = 0.76$ and $\boldsymbol{x}^{\top}[1, 0]^{\top} = 0 \cdot 1 + 1 \cdot 0 = 0$. The orthogonality of the basis makes a graphical representation convenient and the inner product allows for geometric notions as length and angles. In Fig. A.1 a coordinate system is used that is obtained by scaling the vectors \boldsymbol{x} and $[1, 0]^{\top}$ to represent vectors \boldsymbol{x} and \boldsymbol{y} from Ex. 2.1. The vectors are represented as arrows that point to the coordinates that define them. This shows, as already mentioned, that the vectors are defined by a direction and a length; \boldsymbol{x} points in the direction of one of the coordinate axes, because the axis is defined based on it. In the following, it is considered what is meant by "length".

If the inner product is specified as in (A.2), the vector space is called a *Euclidean vector* space, because the norm corresponding to the inner product is the so-called *Euclidean norm* or ℓ_2 -norm, which for is $x \in \mathcal{V}$ (dimensionality N) defined by [GV96, HJ09]

$$\|\boldsymbol{x}\|_2 \stackrel{\text{\tiny def}}{=} \sqrt{\boldsymbol{x}^\top \boldsymbol{x}} = \sqrt{\sum_{j=1}^N x_j^2} .$$
 (A.3)

Another example for a norm is the ℓ_1 -norm or Manhattan norm, which is defined as [GV96, HJ09]

$$\|\boldsymbol{x}\|_{1} \stackrel{\text{\tiny def}}{=} \sum_{j=1}^{N} |x_{j}| . \tag{A.4}$$

Both norms are special cases of the ℓ_p -norm with definition [GV96, HJ09]

1.0

$$\|\boldsymbol{x}\|_{p} \stackrel{\text{\tiny def}}{=} \left(\sum_{j=1}^{N} |x_{j}|^{p}\right)^{1/p} . \tag{A.5}$$

Noteworthy, for p < 1, the definition of the ℓ_p -norm does not yield a norm in the mathematical sense anymore and for $p \to 0$ the result yields the *sparsity*, i.e., the number of non-zero elements, of the vector [EK12]

$$\|\boldsymbol{x}\|_{0} \stackrel{\text{\tiny def}}{=} |\{x_{j} \neq 0 \mid j \in \{1, \dots, N\}\}|.$$
(A.6)

A comparison of the different norms is given in Fig. A.1 for x and y from Ex. 2.1. The Euclidean norm is usually associated as length of a vector.

Generally, the vectors of one vector space can be represented in other vector spaces, cf. [Kay93, PP02]. In Sec. 4.4.2, this is used to represent random variables in the Euclidean space. Note that the vector space of the random variables itself does not use the Euclidean norm as inner product.

A.2 Matrices

By appending vectors column-wise, one obtains matrices, e.g., in Ex. 2.1, it is

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = \begin{bmatrix} 0.75 & -0.25 \\ -0.5 & 0.75 \end{bmatrix} , \qquad (A.7)$$

which consists of the vectors $\tilde{a}_1 = [a_{11}, a_{21}]^{\top} = [0.75, -0.5]^{\top}$ and $\tilde{a}_2 = [a_{12}, a_{22}]^{\top} = [-0.25, 0.75]^{\top}$.

An important property of matrices that is used in this thesis is the so-called *Frobenius* norm [GV96], which is for $\mathbf{A} \in \mathbb{R}^{M \times N}$ given by

$$\|\boldsymbol{A}\|_{\mathrm{F}} \stackrel{\text{\tiny def}}{=} \sqrt{\sum_{i=1}^{M} \sum_{j=1}^{N} a_{ij}^{2}} .$$
(A.8)

The Frobenius norm of the matrix given above is, thus,

$$\|\boldsymbol{A}\|_{\rm F} = \sqrt{2 \cdot 0.75^2 + (-0.25)^2 + (-0.5)^2} = 0.9345$$
 (A.9)

A.2.1 Coherence

The coherence [EK12], which is mentioned in Chapters 2 and 6, is defined as

$$c(\boldsymbol{A}) = \max_{\substack{i,j \in \{1,\dots,N\}\\i \neq j}} \frac{|\tilde{\boldsymbol{a}}_i^\top \tilde{\boldsymbol{a}}_j|}{\|\tilde{\boldsymbol{a}}_i\|_2 \|\tilde{\boldsymbol{a}}_j\|_2}, \qquad (A.10)$$

where $|\cdot|$ denotes the absolute value. A matrix satisfying $c(\mathbf{A}) = 0$ is called an *orthogonal matrix*, because it implies that every column is orthogonal to each other column; noteworthy, an orthogonal matrix needs to be square, i.e., M = N.

A.2.2 Condition Number

Furthermore, the notion of a condition number of a matrix is required. The condition number of \boldsymbol{A} gives the chance to bound the relative error obtained when inverting the matrix to solve a system of equations $\boldsymbol{A}\boldsymbol{x} = \boldsymbol{b}$ for \boldsymbol{x} based on the relative error in \boldsymbol{b} [HJ09]. In this thesis, we restrict to the condition number based on the ℓ_2 -norm, so that the condition number is basically the ratio between largest and smallest singular value of the matrix. With matrix $\boldsymbol{A} \in \mathbb{R}^{M \times N}$ (where M < N) having singular value decomposition (SVD) given by $\boldsymbol{A} = \boldsymbol{U} \operatorname{diag}(s_j) \boldsymbol{V}^{\top}$, i.e., singular values $s_1, \ldots, s_M > 0$ and $s_{M+1} = \cdots = s_N = 0$, where the singular values are sorted in descending order, meaning $s_1 \ge s_2 \ge \cdots \ge s_M$, then the condition number of A is defined as [HJ09, Sec. 5.7 Problem 1]

$$\kappa(\boldsymbol{A}) = \frac{s_1}{s_M} \,. \tag{A.11}$$

If the condition number is large, A is said to be *ill conditioned* or *poorly conditioned*, meaning that errors are generally amplified. For $\kappa(A)$ close or equal to one, one may say that A is *well conditioned*.

When constructing sensing matrices with specified condition number κ , the procedure from [RSF19, Sec. VI. A] is chosen, where the matrix A is constructed via an SVD $A = U \operatorname{diag}(s_j) V^{\top}$ with non-zero singular values being specified by

$$s_j = \kappa^{(M-j+1)/(M-1)}, \quad j \in \{1, \dots, M\},$$
 (A.12)

as well as $U \in \mathbb{R}^{M \times M}$ and $V \in \mathbb{R}^{N \times N}$ being drawn randomly from the group of orthogonal matrices.

B. Convex Optimization

This appendix comprises the notions of optimization theory that are necessary to follow the derivation of the algorithms in Chapter 3 through 5.

The following definitions are according to [BV04]. A set \mathcal{A} is *convex* if for any two elements of the set, let's call them x_1 and x_2 , and $\gamma \in [0, 1]$, the convex combination is in the set, i.e.,

$$\gamma x_1 + (1 - \gamma) x_2 \in \mathcal{A} . \tag{B.1}$$

This holds as well for vector spaces, i.e., \mathcal{X} is convex if for $\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{X}, \gamma \in [0, 1]$, it is $\gamma \boldsymbol{x}_1 + (1 - \gamma) \boldsymbol{x}_2 \in \mathcal{X}$. This means that line segment between \boldsymbol{x}_1 and \boldsymbol{x}_2 is always inside the set \mathcal{X} . A convex function is a function with convex domain satisfying

$$f(\gamma \boldsymbol{x}_1 + (1-\gamma)\boldsymbol{x}_2) \le \gamma f(\boldsymbol{x}_1) + (1-\gamma)f(\boldsymbol{x}_2) , \qquad (B.2)$$

i.e., the line segment connecting the points $f(x_1)$ and $f(x_2)$ is always above the function, cf. Fig. B.1. A function f(x) is *concave*, if -f(x) is convex.

Constrained optimization problems are usually formulated in the following way

minimize
$$f(\boldsymbol{x})$$

subject to $f_k(\boldsymbol{x}) = 0$, $k \in \{1, \dots, K\}$. (B.3)



Figure B.1: Example for a convex function and the line segment according to the definition.

The function f(x) is called the *objective function* or *cost function*, x is the *optimization variable*, and $f_k(x) = 0$, $k \in \{1, ..., K\}$ are the *equality constraints*. We do not define inequality constraints since they do not appear in the thesis. In the given case, the optimization problem is called convex, if the cost function f(x) and the constraints are convex.

A point x that satisfies all constraints, while also being part of the respective domains of the functions, is called *feasible*. A point that minimizes the problem (in the sense of an infimum, if the minimum can not be attained) is called *optimal point* or *optimizer* x^* and $p^* = f(x^*)$ is the *optimal (minimal) value*.

In order to fulfill one equality constraint, an optimal point must lie where the contour line of the objective function and the respective constraint function is tangential, because otherwise that means the objective function changes along the constraint function, which in turn means the objective function can still further be minimized without violating the constraint [Gra14, HF14]. Mathematically, this can be formulated as collinearity of the gradients¹ [Gra14], meaning that the differentiation w.r.t. every possible direction of the constraint function and the objective function must be multiples of each other. If several constraints are imposed, the same reasoning leads to the insight that the gradient of the objective function must be a linear combination of the equality constraints. By introducing respective multipliers $\boldsymbol{\mu} = [\mu_1, \ldots, \mu_K]^{\top}$, this yields a system of equations

$$\nabla f(\boldsymbol{x}^*) \stackrel{!}{=} -\sum_{k=1}^{K} \mu_k \nabla f_k(\boldsymbol{x}^*) .$$
 (B.4)

The reason for placing the minus here, stems from alternatively writing the functions on one side and setting the sum to zero, i.e.,

$$\nabla f(\boldsymbol{x}^*) + \sum_{k=1}^{K} \mu_k \nabla f_k(\boldsymbol{x}^*) \stackrel{!}{=} 0$$
(B.5)

$$\Leftrightarrow \qquad \nabla\left(f(\boldsymbol{x}^*) + \sum_{k=1}^{K} \mu_k f_k(\boldsymbol{x}^*)\right) \stackrel{!}{=} 0. \tag{B.6}$$

Then, it is only one step left to summarizing the procedure by additively combining objective function and the (weighted) constraint functions into a new function. The result is the so-called *Lagrangian*

$$\mathcal{L}(\boldsymbol{x}, \boldsymbol{\mu}) = f(\boldsymbol{x}) + \sum_{k=1}^{K} \mu_k f_k(\boldsymbol{x})$$
(B.7)

with $\boldsymbol{\mu} = [\mu_1, \ldots, \mu_K]^\top$ being called *Lagrange multipliers*.

This way, the optimization problem is converted to finding the stationary points of the Lagrangian. An example in two dimensions is shown for illustration purposes.

Example B.1: _

Here, Ex. 2.1 is extended for the purpose of showing the idea of the Lagrangian. Let $p_{\boldsymbol{y}}(\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi\sigma_n^2)^M}} \exp\left(-\frac{1}{2\sigma_n^2} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2\right)$ with the values of \boldsymbol{A} , \boldsymbol{y} , and σ_n^2 taken from Ex. 2.1

¹Gradients summarize the differentiation with respect to the elements of the argument, i.e., if $\boldsymbol{x} \in \mathbb{R}^N$, then $\nabla f(\boldsymbol{x}) = [\frac{\partial}{\partial x_1} f(\boldsymbol{x}), \dots, \frac{\partial}{\partial x_N} f(\boldsymbol{x})]^\top$. We may use $\nabla f(\boldsymbol{x}^*)$ as shorthand for $\nabla f(\boldsymbol{x})|_{\boldsymbol{x}=\boldsymbol{x}^*}$.



Figure B.2: Contour plot of $p_y(x)$ from Ex. 2.1 with constraints $x_1 = 0$ and $x_2 = 0$ as white lines and the solution $x^* = [-0.68, 0]^{\top}$ of the minimization with constraint $x_2 = 0$.

be the negative cost function, i.e., $f(\mathbf{x}) = -p_{\mathbf{y}}(\mathbf{x})$ should be minimized. Further, assume that it is known that $x_2 \stackrel{!}{=} 0$, so an equality constraint $f_1(\mathbf{x}) = x_2$ is added. $p_{\mathbf{y}}(\mathbf{x})$ is shown as contour plot in Fig. B.2, the equality constraint is the horizontal white line. The Lagrangian reads $\mathcal{L}(\mathbf{x}, \mu_1) = -p_{\mathbf{y}}(\mathbf{x}) + \mu_1 f_1(\mathbf{x})$. Differentiation w.r.t. x_1, x_2 , and μ_1 and subsequent zeroing yields three equations, which can be use to solve for the variables.

Differentiation w.r.t. μ_1 retrieves the constraint $f_1(x)$, i.e.,

$$\frac{\partial}{\partial \mu_1} \mathcal{L}(\boldsymbol{x}, \mu_1) = \frac{\partial}{\partial \mu_1} (-\mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) + \mu_1 f_1(\boldsymbol{x})) = f_1(\boldsymbol{x}) = x_2 \stackrel{!}{=} 0.$$
(B.8)

Together with the other equations, the missing values of x_1 and μ_1 can be obtained.

First note that

$$\frac{\partial}{\partial x_j} \mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) = \mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) \cdot \left(-\frac{1}{\sigma_{\mathsf{n}}^2} \sum_{i=1}^M (y_i - \boldsymbol{a}_i \boldsymbol{x}) \cdot (-a_{ij}) \right) .$$
(B.9)

Therefore,

$$\frac{\partial}{\partial x_1} \mathcal{L}(\boldsymbol{x}, \mu_1) = -\frac{\partial}{\partial x_1} \mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) + \mu_1 \frac{\partial}{\partial x_1} x_2$$
$$= -\mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) \cdot \left(\frac{1}{\sigma_n^2} \sum_{i=1}^M a_{i1}(y_i - \boldsymbol{a}_i \boldsymbol{x})\right) + 0 \stackrel{!}{=} 0 \qquad (B.10)$$

$$\sum_{i=1}^{M} a_{i1}(y_i - a_i x) \stackrel{!}{=} 0 , \qquad (B.11)$$

 \Rightarrow

which yields together with $x_2 = 0$

 \Rightarrow

$$a_{11}(y_1 - a_{11}x_1) + a_{21}(y_2 - a_{21}x_1) \stackrel{!}{=} 0$$
(B.12)

$$\Rightarrow \qquad x_1 = \frac{a_{11}y_1 + a_{21}y_2}{a_{11}^2 + a_{21}^2} = -0.68 . \tag{B.13}$$

With the values $x_1 = -0.68$ and $x_2 = 0$, the last equation yields the value for μ_1

$$\frac{\partial}{\partial x_2} \mathcal{L}(\boldsymbol{x}, \mu_1) = -\frac{\partial}{\partial x_2} \mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) + \mu_1 \frac{\partial}{\partial x_2} x_2$$
$$= -\mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) \cdot \left(\frac{1}{\sigma_{\mathsf{n}}^2} \sum_{i=1}^M a_{i2}(y_i - \boldsymbol{a}_i \boldsymbol{x})\right) + \mu_1 \stackrel{!}{=} 0 \qquad (B.14)$$

$$\mu_1 = \mathbf{p}_{\boldsymbol{y}}(\boldsymbol{x}) \cdot \left(\frac{1}{\sigma_n^2} \sum_{i=1}^M a_{i2}(y_i - \boldsymbol{a}_i \boldsymbol{x})\right) . \tag{B.15}$$

Hence, the solution is given by $x^* = [-0.68, 0]^\top$, which fulfills the equality constraint $x_2 = 0$ and maximizes $p_y(x)$, given the constraint. As one can see, the constraint $x_2 = 0$ and the contour lines are tangential in the solution point.

Note that \boldsymbol{y} was generated by $\boldsymbol{x} = [0, 1]^{\top}$, so the constraint $x_1 = 0$ is much closer to the actual maximum of $p_{\boldsymbol{y}}(\boldsymbol{x})$.

It can be shown that the optimal point *minimizes* the Lagrangian. This property is used in the *Lagrange dual function*. It is defined as

$$\mathcal{L}_{\mathrm{D}}(\boldsymbol{\mu}) = \inf_{\boldsymbol{x}} \mathcal{L}(\boldsymbol{x}, \, \boldsymbol{\mu}) \,. \tag{B.16}$$

Noteworthy, the dual function is concave, even if the problem itself is not convex. This concavity ensures that the dual function yields lower bounds on the optimal value p^* , i.e.,

$$\mathcal{L}_{\mathrm{D}}(\boldsymbol{\mu}) \le p^* . \tag{B.17}$$

If the optimal value can be attained, i.e.,

$$\max_{\boldsymbol{\mu}} \mathcal{L}_{\mathrm{D}}(\boldsymbol{\mu}) = p^* , \qquad (B.18)$$

one speaks of *strong duality*, otherwise of *weak duality*. Note that this leads the way for an approach to solve the optimization problem from a different point of view. The resulting formulation of the problem is called *Lagrange dual problem*

maximize
$$\mathcal{L}_{\mathrm{D}}(\boldsymbol{\mu})$$
 .

Noteworthy, *Slater's condition* states that the optimal value can be attained via this approach, i.e., (B.18) holds in this case, if there exists a feasible point.

C. Functional Derivatives

A *functional* describes a function which has itself functions as arguments; the differential entropy (2.23) is an example for a functional, since it depends on the form of the density q(x). The term *functional derivative* or *variational derivative* stems from the *calculus of variations* [GFS63], which addresses the search for extrema in such functionals.

Here, functionals of the form

$$F(\mathbf{q}(\boldsymbol{x})) = \int f(\mathbf{q}(\boldsymbol{x}), \, \boldsymbol{x}) \, \mathrm{d}\boldsymbol{x}$$
(C.1)

are considered, with $f(\cdot, \cdot)$ being a function of the two variables. The task of variational calculus is to find the function q(x) that maximizes or minimizes the functional. For the derivation of the form of the functional derivative in this case, the necessary conditions to obtain an extremum in such a functional are described. The explanations follow the exposition given in [Gra14] as it focuses on a rather intuitive understanding. An entirely mathematical description is given in [GFS63].

In order to derive conditions for an extremum of such a functional, one considers the optimal function $q^*(x)$ and a *variation* $\delta q(x)$ thereof, i.e., any solution q(x) may be represented as

$$q(\boldsymbol{x}) = q^*(\boldsymbol{x}) + \epsilon \delta q(\boldsymbol{x}) , \qquad (C.2)$$

with scalar parameter ϵ quantifying the amount of variation. Obviously, for $\epsilon \to 0$, the optimal solution is obtained. This also means that at an extremum, the derivative w.r.t. ϵ must vanish, as otherwise $q^*(x)$ would not be the optimal solution. Mathematically, it is

$$\delta F \stackrel{\text{\tiny def}}{=} \frac{\mathrm{d}F(\mathbf{q}^*(\boldsymbol{x}) + \epsilon \delta \mathbf{q}(\boldsymbol{x}))}{\mathrm{d}\epsilon} \Big|_{\epsilon=0} \stackrel{!}{=} 0 , \qquad (C.3)$$

where δF is called the *(first) variation* or *differential*. By use of the chain rule¹ the differentia-

¹The chain rule states that a derivative of nested functions $f(a(\epsilon))$ w.r.t. the variable ϵ can be computed by differentiating the outer function $f(\cdot)$ w.r.t. the inner function as argument and multiplying it to the derivative of the inner function $a(\epsilon)$, i.e., $\frac{d}{d\epsilon}f(a(\epsilon)) = \frac{d}{da}f(a) \cdot \frac{d}{d\epsilon}a(\epsilon)$.

tion w.r.t. ϵ yields in the given case

$$\delta F = \int \frac{\mathrm{d}}{\mathrm{d}\epsilon} f(\mathbf{q}^*(\boldsymbol{x}) + \epsilon \delta \mathbf{q}(\boldsymbol{x}), \, \boldsymbol{x}) \Big|_{\epsilon=0} \mathrm{d}\boldsymbol{x}$$

=
$$\int \frac{\partial}{\partial \mathbf{q}} f(\mathbf{q}(\boldsymbol{x}), \, \boldsymbol{x}) \cdot \frac{\mathrm{d}}{\mathrm{d}\epsilon} (\epsilon \delta \mathbf{q}(\boldsymbol{x})) \mathrm{d}\boldsymbol{x}$$

=
$$\int \frac{\partial}{\partial \mathbf{q}} f(\mathbf{q}(\boldsymbol{x}), \, \boldsymbol{x}) \delta \mathbf{q}(\boldsymbol{x}) \mathrm{d}\boldsymbol{x} .$$
 (C.4)

At this point, one needs to make use of the fundamental lemma of the calculus of variations [GFS63, Lemma 1], which states that for every continuous function $\delta q(x)$ from

$$\int_{\mathcal{X}} f(\boldsymbol{x}) \delta \boldsymbol{q}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = 0 \tag{C.5}$$

it follows that $f(\boldsymbol{x}) = 0$ for any $\boldsymbol{x} \in \mathcal{X}$.

Since $\delta F = 0$ is required for all variations $\delta q(x)$, it can be deduced that

$$\frac{\partial}{\partial \mathbf{q}} f(\mathbf{q}(\boldsymbol{x}), \, \boldsymbol{x}) \stackrel{!}{=} 0 \tag{C.6}$$

must hold at an extremum, i.e., the integral can be dropped for these considerations. Ultimately, the left hand-side of (C.6) yields the desired way of computing the functional derivative.

Noteworthy, the differentiation in (C.6) is carried out as if q(x) were a simple variable instead of a function. Hence, the classical differentiation rules as sum, product, and chain rule, hold as usual. All in all, the procedure to differentiate a functional as in (C.1) w.r.t. q(x) is to drop the integral and differentiate the function(al) inside as if q(x) were a variable (not a function). In order to get familiar with this type of differentiation, the functional derivative of the differential entropy (2.23) and the Kullback–Leibler divergence (2.21) are computed exemplary.

Example C.1: _

Recall that the differential entropy is defined as

$$h(\mathbf{q}(\boldsymbol{x})) = -\int \mathbf{q}(\boldsymbol{x}) \log \mathbf{q}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} , \qquad (C.7)$$

i.e., it is
$$f(\mathbf{q}(\boldsymbol{x}),\,\boldsymbol{x}) = -\mathbf{q}(\boldsymbol{x})\log\mathbf{q}(\boldsymbol{x})$$
. Hence

$$\frac{\partial \mathbf{h}}{\partial \mathbf{q}} = -\frac{\partial}{\partial \mathbf{q}} \mathbf{q}(\boldsymbol{x}) \log \mathbf{q}(\boldsymbol{x}) = -1 - \log \mathbf{q}(\boldsymbol{x}) .$$
(C.8)

Example C.2: ____

The Kullback-Leibler divergence is defined as

$$D_{KL}(\mathbf{q}(\boldsymbol{x}) \mid| f_{\mathbf{x}}(\boldsymbol{x})) = \int \mathbf{q}(\boldsymbol{x}) \log \frac{\mathbf{q}(\boldsymbol{x})}{f_{\mathbf{x}}(\boldsymbol{x})} d\boldsymbol{x} , \qquad (C.9)$$

i.e.,
$$f(\mathbf{q}(\mathbf{x}), \mathbf{x}) = \mathbf{q}(\mathbf{x}) \log \frac{\mathbf{q}(\mathbf{x})}{\mathbf{f}_{\mathbf{x}}(\mathbf{x})}$$
 and, thus,
 $\frac{\partial}{\partial \mathbf{q}} D_{\mathrm{KL}}(\mathbf{q}(\mathbf{x}) || \mathbf{f}_{\mathbf{x}}(\mathbf{x})) = \frac{\partial}{\partial \mathbf{q}} (\mathbf{q}(\mathbf{x}) \log \mathbf{q}(\mathbf{x}) - \mathbf{q}(\mathbf{x}) \log \mathbf{f}_{\mathbf{x}}(\mathbf{x})) = 1 + \log \mathbf{q}(\mathbf{x}) - \log \mathbf{f}_{\mathbf{x}}(\mathbf{x})$. (C.10)

D. Exponential Families

Exponential families denote a class of distributions which are entirely described by certain moments. The moments that specify the respective distribution are collected in the so-called *sufficient statistics* g(x). One can see the sufficient statistic as measure that needs to be known in order to fully characterize a given distribution [Fis22, PP02]. To point out the importance of the sufficient statistics, we derive the form of an exponential family as result of the minimization of the Kullback–Leibler divergence under a moment constraint.

D.1 Projection Property of the Kullback–Leibler Divergence

Throughout the thesis, the minimization of the Kullback–Leibler divergence plays an important role. In this section, we show the strong connection between exponential families and the minimization of the Kullback–Leibler divergence under a moment constraint by considering a constrained optimization problem and especially its Lagrangian. The task is to approximate an arbitrary distribution $f_x(x)$ by a substitute distribution q(x), while being consistent to the moments $E_{x \sim f_x} \{g(x)\}$, i.e., a constrained optimization problem

$$\begin{array}{ll} \text{minimize} & \mathrm{D}_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) \mid\mid \mathsf{f}_{\mathbf{x}}(\boldsymbol{x})) \\ \text{subject to} & \mathrm{E}_{\mathbf{x} \sim \mathsf{q}}\{\boldsymbol{g}(\mathbf{x})\} = \mathrm{E}_{\mathbf{x} \sim \mathsf{f}_{\mathbf{x}}}\{\boldsymbol{g}(\mathbf{x})\} \\ & \int \mathsf{q}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = 1 \end{array}$$
 (D.1)

is given. The respective Lagrangian with multipliers θ for the moment constraint and ν for the normalization constraint reads

$$\mathcal{L}(\mathbf{q}(\boldsymbol{x}), \boldsymbol{\theta}, \nu) = D_{\mathrm{KL}}(\mathbf{q}(\boldsymbol{x}) \mid\mid f_{\mathbf{x}}(\boldsymbol{x})) + \boldsymbol{\theta}^{\top} (E_{\mathbf{x} \sim q} \{ \boldsymbol{g}(\mathbf{x}) \} - E_{\mathbf{x} \sim f_{\mathbf{x}}} \{ \boldsymbol{g}(\mathbf{x}) \}) + \nu \left(\int \mathbf{q}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - 1 \right) .$$
(D.2)

For the functional derivative with respect to q(x), one gets

$$\frac{\partial \mathcal{L}}{\partial \mathbf{q}} = \log \mathbf{q}(\mathbf{x}) + 1 - \log f_{\mathbf{x}}(\mathbf{x}) + \boldsymbol{\theta}^{\top} \boldsymbol{g}(\mathbf{x}) + \nu \stackrel{!}{=} 0.$$
(D.3)

Hence, the form of the substitute distribution that minimizes the Kullback-Leibler divergence under the moment constraint is given by

$$\mathbf{q}(\boldsymbol{x}) = \mathbf{f}_{\mathbf{x}}(\boldsymbol{x}) \exp\left(\boldsymbol{\theta}^{\top} \boldsymbol{g}(\boldsymbol{x}) - 1 - \nu\right) . \tag{D.4}$$

For q(x) to fulfill the normalization constraint, the constant terms together with the multiplier ν must form the following equality

$$Z(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \exp(1+\nu) = \int f_{\mathbf{x}}(\boldsymbol{x}) \exp\left(\boldsymbol{\theta}^{\top} \boldsymbol{g}(\boldsymbol{x})\right) \, \mathrm{d}\boldsymbol{x}$$
(D.5)

and is called *partition function*.

All in all, an exponential family $f_{\theta}(x)$ is represented as¹

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} \exp\left(\boldsymbol{\theta}^{\top} \boldsymbol{g}(\boldsymbol{x}) + k(\boldsymbol{x})\right) .$$
(D.6)

The elements of the Lagrangian multiplier θ are called *natural parameters*. The pdf $f_x(x)$ is called *base measure* or *carrier density*, $\log f_x(x)$ is called *carrier measure* [NG09].

This also means that minimizing the Kullback–Leibler divergence of q(x) with respect to $f_x(x)$ forces the substitute distribution to be a member of an exponential family, if a moment constraint is given; the family is specified by the respective moments.

Note that the Boltzmann distribution (2.32) is also an exponential family with sufficient statistics $-U(\mathbf{x})$ and natural parameter β .

D.2 Sufficient Statistics

It has been shown above that the sufficient statistic plays an important role for an exponential family, as it specifies the measure that needs to be known for full characterization of the distribution. Since the thesis mainly works with the exponential family of Gaussian distributions, the corresponding specifications of the sufficient statistics are displayed here. For the first-order moments, we define

$$\boldsymbol{g}_{\boldsymbol{\lambda}}(\boldsymbol{x}) = [x_1, \dots, x_N]^\top$$
, i.e., $g_{\lambda_j}(x_j) = x_j$, (D.7)

with corresponding natural parameters

$$\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_N]^\top, \qquad (D.8)$$

while for the second-order moment, either a vector parameter $\mathbf{\Lambda} = [\Lambda_1, \ldots, \Lambda_N]^{\top}$ with sufficient statistics

$$\boldsymbol{g}_{\boldsymbol{\Lambda}}(\boldsymbol{x}) = -\frac{1}{2} [x_1^2, \dots, x_N^2]^{\top} = -\frac{1}{2} \mathbf{diag} (\boldsymbol{x} \boldsymbol{x}^{\top}), \qquad \text{i.e.,} \quad g_{\Lambda_j}(x_j) = -\frac{1}{2} x_j^2 \qquad (D.9)$$

¹Use the definition $k(\boldsymbol{x}) \stackrel{\text{def}}{=} \log f_{\boldsymbol{x}}(\boldsymbol{x})$ for consistency to the last representation
	individual variances	average variances
$egin{array}{c} egin{array}{c} egin{array}$	$[oldsymbol{x}^{ op},-rac{1}{2} extbf{diag}(oldsymbol{x}oldsymbol{x}^{ op})^{ op}]^{ op}$	$[oldsymbol{x}^ op,-rac{1}{2}oldsymbol{x}^ opoldsymbol{x}]^ op$
λ	$oldsymbol{\Phi}_{x}^{-1}oldsymbol{m}_{x}$	$oldsymbol{m}_{x}/\sigma_{x}^2$
λ_j	$m_{{f x},j}/\sigma_{{f x},j}^2$	$m_{{\sf x},j}/\sigma_{\sf x}^2$
Λ_j	$1/\sigma_{x,j}^2$	_
Λ	_	$1/\sigma_{x}^2$

Table D.1: Summary of the parameter connections for Gaussian densities with individual and average variances.

or, alternatively, a single parameter Λ with corresponding statistics function

$$g_A(\boldsymbol{x}) = -\frac{1}{2} \sum_{j=1}^N x_j^2 = -\frac{1}{2} \boldsymbol{x}^\top \boldsymbol{x}$$
 (D.10)

is utilized.

Together, the individual variances case is given by

$$\boldsymbol{g}(\boldsymbol{x}) = [\boldsymbol{g}_{\boldsymbol{\lambda}}(\boldsymbol{x})^{\top}, \, \boldsymbol{g}_{\boldsymbol{\Lambda}}(\boldsymbol{x})^{\top}]^{\top}, \quad \text{with} \quad \boldsymbol{\theta} = [\boldsymbol{\lambda}^{\top}, \, \boldsymbol{\Lambda}^{\top}]^{\top}$$
(D.11)

and in the average variance case

$$\boldsymbol{g}(\boldsymbol{x}) = [\boldsymbol{g}_{\boldsymbol{\lambda}}(\boldsymbol{x})^{\top}, g_{\boldsymbol{\Lambda}}(\boldsymbol{x})]^{\top}, \text{ with } \boldsymbol{\theta} = [\boldsymbol{\lambda}^{\top}, \boldsymbol{\Lambda}]^{\top}.$$
 (D.12)

The representation of Gaussian distributions as exponential families and the connection between the different parameterizations are explained in the examples below. The results are summarized in Table D.1.

Example D.1: _

Consider the *N*-dimensional Gaussian random vector $\mathbf{x} \sim \mathcal{N}(\boldsymbol{m}_x, \sigma_x^2 \boldsymbol{I})$ with mean \boldsymbol{m}_x and (isotropic) variance σ_x^2 , i.e., the elements of \mathbf{x} are independent but share the same variance. The probability density function reads

$$f_{\mathbf{x}}(\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi\sigma_{\mathbf{x}}^2)^N}} \exp\left(-\frac{1}{2\sigma_{\mathbf{x}}^2}(\boldsymbol{x} - \boldsymbol{m}_{\mathbf{x}})^\top (\boldsymbol{x} - \boldsymbol{m}_{\mathbf{x}})\right) .$$
(D.13)

The density can be expressed as exponential family with sufficient statistics $g(x) = [x^{\top}, -\frac{1}{2}x^{\top}x]^{\top}$ and natural parameters $\lambda = m_x/\sigma_x^2$ and $\Lambda = \frac{1}{\sigma_x^2}$, as well as partition function $Z(\lambda, \Lambda) =$

$$\sqrt{(2\pi\sigma_{\mathbf{x}}^{2})^{N}} \exp\left(\frac{1}{2\Lambda}\boldsymbol{\lambda}^{\top}\boldsymbol{\lambda}\right), \text{ since}$$

$$-\frac{1}{2\sigma_{\mathbf{x}}^{2}}(\boldsymbol{x}-\boldsymbol{m}_{\mathbf{x}})^{\top}(\boldsymbol{x}-\boldsymbol{m}_{\mathbf{x}}) = \boldsymbol{m}_{\mathbf{x}}^{\top}\boldsymbol{x}/\sigma_{\mathbf{x}}^{2} - \frac{1}{2\sigma_{\mathbf{x}}^{2}}\left(\boldsymbol{m}_{\mathbf{x}}^{\top}\boldsymbol{m}_{\mathbf{x}} + \boldsymbol{x}^{\top}\boldsymbol{x}\right)$$

$$= \frac{1}{\sigma_{\mathbf{x}}^{2}}\boldsymbol{m}_{\mathbf{x}}^{\top}\boldsymbol{x} - \frac{1}{2}\left(\frac{1}{\sigma_{\mathbf{x}}^{2}}\boldsymbol{m}_{\mathbf{x}}^{\top}\frac{\sigma_{\mathbf{x}}^{2}}{\sigma_{\mathbf{x}}^{2}}\boldsymbol{m}_{\mathbf{x}} + \frac{1}{\sigma_{\mathbf{x}}^{2}}\boldsymbol{x}^{\top}\boldsymbol{x}\right)$$

$$= \boldsymbol{\lambda}^{\top}\boldsymbol{x} - \frac{1}{2\Lambda}\boldsymbol{\lambda}^{\top}\boldsymbol{\lambda} - \frac{\Lambda}{2}\boldsymbol{x}^{\top}\boldsymbol{x}.$$
(D.14)

Using the denomination for exponential families, the pdf becomes

$$\mathbf{f}_{\boldsymbol{\lambda},\Lambda}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\lambda},\Lambda)} \exp\left(\boldsymbol{\lambda}^{\top} \boldsymbol{x} - \frac{\Lambda}{2} \boldsymbol{x}^{\top} \boldsymbol{x}\right) . \tag{D.15}$$

Example D.2:

Here, the *N*-dimensional Gaussian random vector $\mathbf{x} \sim \mathcal{N}(\mathbf{m}_x, \mathbf{\Phi}_x)$ with mean \mathbf{m}_x and covariance matrix $\mathbf{\Phi}_x = \operatorname{diag}(\sigma_{x,j}^2)$ is considered, i.e., the elements of \mathbf{x} are independent but not identically distributed, since each has an individual variance. The probability density function reads

$$f_{\mathbf{x}}(\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi)^N \det \boldsymbol{\Phi}_{\mathbf{x}}}} \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{m}_{\mathbf{x}})^\top \boldsymbol{\Phi}_{\mathbf{x}}^{-1}(\boldsymbol{x} - \boldsymbol{m}_{\mathbf{x}})\right) .$$
(D.16)

With sufficient statistics $\boldsymbol{g}(\boldsymbol{x}) = [\boldsymbol{x}^{\top}, -\frac{1}{2}\operatorname{diag}(\boldsymbol{x}\boldsymbol{x}^{\top})^{\top}]^{\top}$ and natural parameters $\boldsymbol{\lambda} = \boldsymbol{\Phi}_{\mathsf{x}}^{-1}\boldsymbol{m}_{\mathsf{x}}$ and $\boldsymbol{\Lambda} = \operatorname{diag}(\boldsymbol{\Phi}_{\mathsf{x}}^{-1}) = [(\sigma_{\mathsf{x},1}^2)^{-1}, \ldots, (\sigma_{\mathsf{x},N}^2)^{-1}]^{\top}$, i.e., $\lambda_j = m_{\mathsf{x},j}/\sigma_{\mathsf{x},j}^2$, respectively $\Lambda_j = 1/\sigma_{\mathsf{x},j}^2$. The partition function reads in this case $Z(\boldsymbol{\lambda}, \boldsymbol{\Lambda}) = \sqrt{(2\pi)^N \det \boldsymbol{\Phi}_{\mathsf{x}}} \exp\left(\frac{1}{2}\boldsymbol{m}_{\mathsf{x}}^{\top}\boldsymbol{\lambda}\right)$. The reformulation of the exponent reads in this case

$$-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{m}_{\mathsf{x}})^{\top} \boldsymbol{\varPhi}_{\mathsf{x}}^{-1}(\boldsymbol{x} - \boldsymbol{m}_{\mathsf{x}}) = \boldsymbol{m}_{\mathsf{x}}^{\top} \boldsymbol{\varPhi}_{\mathsf{x}}^{-1} \boldsymbol{x} - \frac{1}{2} \left(\boldsymbol{m}_{\mathsf{x}}^{\top} \boldsymbol{\varPhi}_{\mathsf{x}}^{-1} \boldsymbol{m}_{\mathsf{x}} + \boldsymbol{x}^{\top} \boldsymbol{\varPhi}_{\mathsf{x}}^{-1} \boldsymbol{x} \right)$$
$$= \boldsymbol{\lambda}^{\top} \boldsymbol{x} - \frac{1}{2} \left(\boldsymbol{m}_{\mathsf{x}}^{\top} \boldsymbol{\lambda} + \boldsymbol{\Lambda}^{\top} \operatorname{diag} \left(\boldsymbol{x} \boldsymbol{x}^{\top} \right) \right) .$$
(D.17)

D.3 Derivative of the Log-Partition Function

One very helpful property of exponential families is that the differentiation of the logarithm of partition function $Z(\theta)$ (usually called log-partition function) with respect to the natural parameters θ yields the moments (specified by the sufficient statistics) of the exponential family [WJ08]

$$\nabla \log Z(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \log \int \exp\left(\boldsymbol{\theta}^{\top} \boldsymbol{g}(\boldsymbol{x}) + k(\boldsymbol{x})\right) \, \mathrm{d}\boldsymbol{x}$$
$$= \frac{1}{Z(\boldsymbol{\theta})} \int \boldsymbol{g}(\boldsymbol{x}) \exp\left(\boldsymbol{\theta}^{\top} \boldsymbol{g}(\boldsymbol{x}) + k(\boldsymbol{x})\right) \, \mathrm{d}\boldsymbol{x}$$
$$= \int \boldsymbol{g}(\boldsymbol{x}) f_{\boldsymbol{\theta}}(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = \mathrm{E}_{\mathbf{x} \sim f_{\boldsymbol{\theta}}} \{\boldsymbol{g}(\mathbf{x})\} \,.$$
(D.18)

E. Reconstruction Algorithms

On the following pages, all algorithms that are used within the thesis are summarized in pseudo code notation. The stopping criterion, mentioned in the algorithms below, can, e.g., be realized by comparing the squared (Euclidean) distance between the estimates m_c and m_s to a given threshold ϵ , i.e.,

$$\|\boldsymbol{m}_{\rm c} - \boldsymbol{m}_{\rm s}\|_2^2 \le \epsilon \; . \tag{E.1}$$

In the simulations for this thesis, the algorithms always run for a given number of iterations before the result is evaluated.

E.1 Gaussian Message Passing (GMP)

This algorithm is close to an algorithm stated in [KMS⁺12] and results from the fine-granular view on the factor graph. Note that the extrinsic is here pre-processed. The notation here utilizes means and variances for comprehensibility; using natural parameters (scaled means and inverse variances) where appropriate instead can lead to better computational complexity, because some (scalar) inversions can be saved.

Algorithm E.1: $\boldsymbol{m}_{\rm s} = \text{GMP}(\boldsymbol{y}, \boldsymbol{A}, \sigma_{\rm n}^2)$ $\begin{array}{ll} \mathbf{1} \ m_{\mathrm{s},j}^{\backslash i} = 0 \ , \ \sigma_{\mathrm{s},j}^{2,\backslash i} = s/N \ \forall \ j \in \{1, \ \dots, \ N\}, \ i \in \{1, \ \dots, \ M\} \\ \mathbf{2} \ m_{\mathrm{c},i,j}^{\backslash j} = 0 \ , \ \sigma_{\mathrm{c},i,j}^{2,\backslash j} = s/N \ \forall \ j \in \{1, \ \dots, \ N\}, \ i \in \{1, \ \dots, \ M\} \end{array}$ 3 while stopping criterion not met do Choose $j \in \{1, ..., N\}$ and $i \in \{1, ..., M\}$ 4 $m_{\mathrm{c},i,j}^{\backslash j} = (y_i - \sum_{\substack{j'=1\\j'\neq j}}^{N} a_{ij'} m_{\mathrm{s},j'}^{\backslash i})/a_{ij}$ $\sigma_{\mathrm{c},i,j}^{2,\backslash j} = (\sigma_{\mathrm{n}}^2 + \sum_{\substack{j'=1\\j'\neq j}}^{N} \sigma_{\mathrm{s},j'}^{2,\backslash i} a_{ij'}^2)/a_{ij}^2$ // 'upward' messages 5 6 $\left| \begin{array}{c} \tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i} = \left(\sum_{\substack{i'=1\\i'\neq i}}^{M} 1/\sigma_{\mathrm{c},i',j}^{2,\backslash j} \right)^{-1} \\ \tilde{x}_{\mathrm{s},j}^{\backslash i} = \tilde{\sigma}_{\mathrm{s},j}^{2,\backslash i} \sum_{\substack{i'=1\\i'\neq i}}^{M} m_{\mathrm{c},i',j}^{\backslash j}/\sigma_{\mathrm{c},i',j}^{2,\backslash j} \end{array} \right|$ 7 // input to NLMMSE (edge-dependent) 8 9 $\begin{array}{|c|c|c|c|c|c|} & m_{\mathbf{s},j}^{\backslash i} = \mathbf{E}_{\mathsf{x}_j} \{ \mathsf{x}_j \mid \tilde{x}_{\mathbf{s},j}^{\backslash i}, \, \tilde{\sigma}_{\mathbf{s},j}^{2,\backslash i} \} \\ & \sigma_{\mathbf{s},j}^{2,\backslash i} = \mathbf{E}_{\mathsf{x}_j} \{ (\mathsf{x}_j - m_{\mathbf{s},j}^{\backslash i})^2 \mid \tilde{x}_{\mathbf{s},j}^{\backslash i}, \, \tilde{\sigma}_{\mathbf{s},j}^{2,\backslash i} \} \end{array}$ // ''downward'' messages $\begin{array}{ll} \mathbf{n} & \tilde{\sigma}_{\mathrm{s},j}^{2} = \left(\sum_{i=1}^{M} 1/\sigma_{\mathrm{c},i',j}^{2,\backslash j}\right)^{-1} & \forall \ j \in \{1, \ldots, N\} \\ \mathbf{n} & \tilde{x}_{\mathrm{s},j} = \tilde{\sigma}_{\mathrm{s},j}^{2} \sum_{i=1}^{M} m_{\mathrm{c},i',j}^{\backslash j} / \sigma_{\mathrm{c},i',j}^{2,\backslash j} & \forall \ j \in \{1, \ldots, N\} \\ \mathbf{n} & \mathbf{m}_{\mathrm{s}} = \mathrm{E}_{\mathbf{x}} \{\mathbf{x} \mid \tilde{\mathbf{x}}_{\mathrm{s}}, \tilde{\sigma}_{\mathrm{s},1}^{2}, \ldots, \tilde{\sigma}_{\mathrm{s},N}^{2}\} \end{array}$ // NLMMSE

E.2 GMP with Post-processed Extrinsic

In contrast to the previous algorithm, an algorithm that computes the extrinsics via postprocessing is shown, which is computationally much simpler because the estimates do not have to be calculated per edge $(M \cdot N)$, but per factor node (M + N). Here, scaled means $\lambda_{ij}^{\setminus j}$ and inverse variances $\Lambda_{ij}^{\setminus j}$ are used as upward messages, because they are summed up in the next step, i.e., a conversion to mean and variance would have to be converted back immediately.

Algorithm E.2: $m{m}_{
m s}={ t GMP}{ t postExt}(m{y},\,m{A},\,\sigma_{
m p}^2)$ $\begin{array}{l} \mathbf{1} \ m_{\mathbf{s},j}^{\backslash i} = 0 \ , \ \sigma_{\mathbf{s},j}^{2,\backslash i} = s/N \ \forall \ j \in \{1, \ \dots, \ N\}, \ i \in \{1, \ \dots, \ M\} \\ \mathbf{2} \ m_{\mathbf{c},i,j}^{\backslash j} = 0 \ , \ \sigma_{\mathbf{c},i,j}^{2,\backslash j} = s/N \ \forall \ j \in \{1, \ \dots, \ N\}, \ i \in \{1, \ \dots, \ M\} \end{array}$ // initialization 3 while stopping criterion not met do for $i \in \{1, ..., M\}$ do 4 $\boldsymbol{m}_{\mathrm{c},i} = \boldsymbol{m}_{\mathrm{s}}^{\backslash i} + \mathrm{diag}(\sigma_{\mathrm{s},j}^{2,\backslash i}) \boldsymbol{a}_{i}(y_{i} - \boldsymbol{a}_{i}^{\top}\boldsymbol{m}_{\mathrm{s}}^{\backslash i}) / (\sigma_{\mathrm{n}}^{2} + \sum_{j'=1}^{N} \sigma_{\mathrm{s},j'}^{2,\backslash i} a_{ij'}^{2})$ 5 $\begin{array}{l} \mathbf{for} \; j \in \{1, \, \dots, \, N\} \; \mathbf{do} \\ \left| \begin{array}{c} \sigma_{\mathrm{c},i,j}^2 = \sigma_{\mathrm{s},j}^{2,\backslash i} (1 - \sigma_{\mathrm{s},j}^{2,\backslash i} a_{ij}^2 / (\sigma_{\mathrm{n}}^2 + \sum_{j'=1}^N \sigma_{\mathrm{s},j'}^{2,\backslash i} a_{ij'}^2)) \end{array} \right. \end{array}$ 6 7 $\begin{array}{|c|c|c|c|} & \Lambda_{ij}^{\backslash j} = 1/\sigma_{\mathrm{c},i,j}^2 - 1/\sigma_{\mathrm{s},j}^{2,\backslash i} \\ & \lambda_{ij}^{\backslash j} = m_{\mathrm{c},i,j}/\sigma_{\mathrm{c},i,j}^2 - m_{\mathrm{s},j}^{\backslash i}/\sigma_{\mathrm{s},j}^{2,\backslash i} \end{array}$ // 'upward' messages 8 9 for $j \in \{1, ..., N\}$ do 10 $\tilde{\sigma}_{\mathrm{s},j}^2 = \left(\sum_{i'=1}^M \Lambda_{i'j}^{\setminus j}\right)^{-1}$ // input to NLMMSE (not edge-dependent) 11 $\tilde{x}_{\mathrm{s},j} = \tilde{\sigma}_{\mathrm{s},j}^2 \sum_{i'=1}^M \lambda_{i'}^{\lambda_{i'}}$ 12 $m_{\mathbf{s},j} = \mathbf{E}_{\mathbf{x}_j} \{ \mathbf{x}_j \mid \tilde{x}_{\mathbf{s},j}, \, \tilde{\sigma}_{\mathbf{s},j}^2 \}$ $\sigma_{\mathbf{s},j}^2 = \mathbf{E}_{\mathbf{x}_j} \{ (\mathbf{x}_j - m_{\mathbf{s},j})^2 \mid \tilde{x}_{\mathbf{s},j}, \, \tilde{\sigma}_{\mathbf{s},j}^2 \}$ // NLMMSE 13 14 $\left| \begin{array}{c} \mathbf{for} \ i \in \{1, \dots, M\} \ \mathbf{do} \\ \left| \begin{array}{c} \sigma_{\mathrm{s},j}^{2,\backslash i} = (1/\sigma_{\mathrm{s},j}^2 - \Lambda_{ij}^{\backslash j})^{-1} \\ m_{\mathrm{s},j}^{\backslash i} = \sigma_{\mathrm{s},j}^{2,\backslash i} (m_{\mathrm{s},j}/\sigma_{\mathrm{s},j}^2 - \lambda_{ij}^{\backslash j}) \end{array} \right|$ 15 // ''downward'' messages 16 17 18 $\boldsymbol{m}_{\mathrm{s}} = [m_{\mathrm{s},1},\,\ldots,\,m_{\mathrm{s},N}]^{ op}$

E.3 Approximate Message Passing (AMP)

The AMP algorithm was first introduced in [DMM09] and rigorously analyzed in [BM11]. Here, the "Bayesian version" [DMM10] is displayed, i.e., the case with non-linear MMSE estimator as so-called *denoiser*, i.e., for the signal-constrained estimation. The algorithm is stated in the notation used throughout the thesis. In AMP, the residuum with the so-called *Onsager correction term* plays an important role, since it removes the correlations between estimates of different iterations, which is in the other algorithms done by bias compensation or computing extrinsic. Consequently, the processing omits the unbiased estimates, i.e., when comparing to VAMP one may set $\tilde{x}_c = m_s$ and $\tilde{x}_s = m_c$, as well as $\tilde{\sigma}_c^2 = \sigma_s^2$ and $\tilde{\sigma}_s^2 = \sigma_c^2$ before the estimations.

Algorithm E.3: $m{m}_{ m s}= extsf{AMP}(m{y},m{A},\sigma_{ m n}^2)$	
1 $oldsymbol{m}_{ m s}=oldsymbol{0},oldsymbol{r}=oldsymbol{y},\sigma_{ m s}^2=s/N$	<pre>// initialization</pre>
2 while stopping criterion not met do	
3 $oldsymbol{m}_{ m c}=oldsymbol{m}_{ m s}+oldsymbol{A}^{ op}oldsymbol{r}$	// MF
$4 \qquad \mathbf{\sigma}_{\mathrm{c}}^2 = \sigma_{\mathrm{n}}^2 + \frac{N}{M}\sigma_{\mathrm{s}}^2$	
5 $\boldsymbol{m}_{s} = \mathrm{E}\{\boldsymbol{x} \mid \boldsymbol{m}_{c}, \sigma_{c}^{2}\}$	// NLMMSE
$\sigma_{\rm s}^2 = \frac{1}{N} \sum_{j=1}^{N} E_{\rm x} \{ ({\rm x} - m_{{\rm s},j})^2 \mid m_{{\rm c},j}, \sigma_{\rm c}^2 \}$	
7 $\left[{\left. {\left. {{f{r}} = {m{y}} - {m{A}{m{m}_{ m{s}}}} + rac{N}{M}rac{{\sigma _{ m{s}}^2}}{{\sigma _{ m{c}}^2}} \cdot {m{r}} ight.} ight.} ight.$	// residuum with Onsager term
8 $oldsymbol{m}_{\mathrm{s}} = [m_{\mathrm{s},1},\ldots,m_{\mathrm{s},N}]^{ op}$	

E.4 Generalized Approximate Message Passing (GAMP)

When referring to GAMP, what is meant is the algorithm stated in [Ran11], but restricted to additive Gaussian noise in the measurements. Closer to the notation here, is the same algorithm stated in [KMS⁺12, Sec. 3.2] or [MKTZ15, Sec. 2.2], where it is called relaxed Belief Propagation.

Algorithm E.4: $m_{\rm s} = \text{GAMP}(\boldsymbol{y}, \boldsymbol{A}, \sigma_{\rm n}^2)$ 1 $m_{\rm s} = 0, \sigma_{A,i}^2 = \frac{s}{N} \boldsymbol{a}_i^{\mathsf{T}} \boldsymbol{a}_i \ i \in \{1, \ldots, M\}, \ \boldsymbol{r} = \text{diag}(\frac{1}{\sigma_{\rm n}^2 + \sigma_{A,i}^2})\boldsymbol{y}$ // initialization2 while stopping criterion not met do3 $\left| \begin{array}{c} \sigma_{\rm c,j}^2 = \left(\sum_{i=1}^M \frac{a_{ij}^2}{\sigma_{\rm n}^2 + \sigma_{A,i}^2}\right)^{-1} \quad j \in \{1, \ldots, N\}$ 4 $m_{\rm c} = \boldsymbol{m}_{\rm s} + \text{diag}(\sigma_{\rm c,j}^2)\boldsymbol{A}^{\mathsf{T}}\boldsymbol{r}$ // MF with bias compensation5 $\left| \begin{array}{c} \text{for } j \in \{1, \ldots, N\} \ \boldsymbol{do} \\ m_{\rm s,j} = E_{\rm x}\{{\rm x} \mid m_{\rm c,j}, \sigma_{\rm c,j}^2\} \\ \sigma_{\rm s,j}^2 = E_{\rm x}\{({\rm x} - m_{\rm s,j})^2 \mid m_{\rm c,j}, \sigma_{\rm c,j}^2\} \\ \end{array} \right| \left| \begin{array}{c} \sigma_{A,i}^2 = \boldsymbol{a}_i^{\mathsf{T}} \text{diag}(\sigma_{\rm s,j}^2) \boldsymbol{a}_i \\ r = \text{diag}(\frac{1}{\sigma_{n}^2 + \sigma_{A,i}^2})(\boldsymbol{y} - \boldsymbol{A}\boldsymbol{m}_{\rm s} + \text{diag}(\sigma_{A,i}^2)\boldsymbol{r}) \\ \end{array} \right| \left| \begin{array}{c} r \in \{1, \ldots, M\} \\ r \in [m_{\rm s,1}, \ldots, m_{\rm s,N}]^{\mathsf{T}} \end{array} \right|$

Note that GAMP imitates the LMMSE inverse from (2.10) by scaling the residuum with $\operatorname{diag}(1/(\sigma_n^2 + \sigma_{A,i}^2))$, which is the inverse of the diagonal of $\boldsymbol{M} = \boldsymbol{A} \boldsymbol{\Phi}_{\mathbf{x}} \boldsymbol{A}^{\top} + \sigma_n^2 \boldsymbol{I}_M$. The scaling with $\operatorname{diag}(\sigma_{c,j}^2)$ in the linear estimation (Line 4) performs the respective individual bias compensation, since the entries of the end-to-end cascade between \boldsymbol{x} and $\boldsymbol{m}_s + \boldsymbol{A}^{\top} \boldsymbol{r}$ are here given by

$$[\mathbf{A}^{\top} \mathbf{diag} \left(\frac{1}{\sigma_{\mathsf{n}}^{2} + \sigma_{\mathbf{A},i}^{2}} \right) \mathbf{A}]_{jj} = \sum_{i=1}^{M} \frac{a_{ij}^{2}}{\sigma_{\mathsf{n}}^{2} + \sigma_{\mathbf{A},i}^{2}} = \frac{1}{\sigma_{\mathsf{c},j}^{2}}, \quad j \in \{1, \dots, N\}.$$
(E.2)

Note that the usual unbiased variance is given by (4.56)

$$\begin{aligned} \sigma_{\mathrm{c},\mathrm{u},j}^{2} &= \sigma_{\mathrm{x},j}^{2} ([\boldsymbol{W}]_{jj} - 1) = \sigma_{\mathrm{x},j}^{2} (\frac{1}{[\boldsymbol{K}]_{jj}} - 1) \\ &= \sigma_{\mathrm{x},j}^{2} \left(\frac{1}{[\boldsymbol{\Phi}_{\mathrm{x}} \boldsymbol{A}^{\top} \boldsymbol{M}^{-1} \boldsymbol{A}]_{jj}} - 1 \right) \\ &= \frac{1}{[\boldsymbol{A}^{\top} \boldsymbol{M}^{-1} \boldsymbol{A}]_{jj}} - \sigma_{\mathrm{x},j}^{2} , \end{aligned}$$
(E.3)

i.e., the variance $\sigma_{c,j}^2$ in the algorithm is an approximation of an unbiased variance, which should work at least for small variances $\sigma_{x,j}^2$. Compared to VAMP, one may interpret the Onsager correction term in Line 9 as a compensation for the neglection of the off-diagonal elements in the LMMSE inverse and the missing of the bias compensation procedure after the signal-constrained estimation.

E.5 Vector Approximate Message Passing (VAMP)

The VAMP algorithm was first introduced in [RSF17], derived from the expectation-consistent (EC) approximate inference framework [OW05]. A different approach leading to the same algorithm is shown in [Spa19, SF17].

Alg	gorithm E.5: $m{m}_{ m s}=$ VAMP $(m{y},m{A},\sigma_{ m n}^2)$	
ı õ	1 $ ilde{\sigma}_{ m c}^2 = s/N, \; ilde{x}_{ m c} = 0$ // initialization	
2 W	while stopping criterion not met do	
3	$oldsymbol{m}_{ ext{c}} = ilde{oldsymbol{x}}_{ ext{c}} + (oldsymbol{A}^{ op}oldsymbol{A} + rac{\sigma_{ ext{n}}^2}{ ilde{\sigma}_{ ext{c}}^2}oldsymbol{I}_N)^{-1}oldsymbol{A}^{ op}(oldsymbol{y} - oldsymbol{A} ilde{oldsymbol{x}}_{ ext{c}})$	// LMMSE
4	$\sigma_{\mathrm{c}}^2 = \mathrm{trace}(\sigma_{n}^2 (\boldsymbol{A}^{\top} \boldsymbol{A} + rac{\sigma_{\mathrm{n}}^2}{\tilde{\sigma}_{\mathrm{c}}^2} \boldsymbol{I}_N)^{-1})/N$	
5	$ ilde{\sigma}_{ m s}^2 = (rac{1}{\sigma_{ m c}^2} - rac{1}{ ilde{\sigma}_{ m c}^2})^{-1}$	<pre>// unbiasing</pre>
6	$ ilde{m{x}}_{ m s} = ilde{\sigma}_{ m s}^2 (rac{m_{ m c}}{\sigma_{ m c}^2} - rac{ ilde{m{x}}_{ m c}}{ ilde{\sigma}_{ m c}^2})$	
7	$oldsymbol{m}_{\mathrm{s}} = \mathrm{E}_{x}\{x \mid ilde{x}_{\mathrm{s}}, ilde{\sigma}_{\mathrm{s}}^2\}$	// NLMMSE
8	$\sigma_{\mathrm{s}}^2 = rac{1}{N} \sum_{j=1}^{N} \mathrm{E}_{x} \{ (x - m_{\mathrm{s},j})^2 \mid \tilde{x}_{\mathrm{s},j}, \tilde{\sigma}_{\mathrm{s}}^2 \}$	
9	$ ilde{\sigma}_{ m c}^2 = (rac{1}{\sigma_{ m s}^2} - rac{1}{ ilde{\sigma}_{ m s}^2})^{-1}$	<pre>// unbiasing</pre>
10	$\sum ilde{oldsymbol{x}}_{ m c} = ilde{\sigma}_{ m c}^2 (rac{m_{ m s}}{\sigma_{ m s}^2} - rac{ ilde{oldsymbol{x}}_{ m s}}{ ilde{\sigma}_{ m s}^2})$	

The representation here, shall show the connection to extrinsic calculation as subtraction of natural parameters. For the implementation, it is better suited to directly compute the *unbiased* channel-constrained estimate via

$$w = \left(\frac{1}{N} \operatorname{trace} \left(\boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{\mathsf{n}}^{2}}{\tilde{\sigma}_{\mathsf{c}}^{2}} \boldsymbol{I}_{N})^{-1} \boldsymbol{A} \right) \right)^{-1} , \qquad (E.4)$$

$$\tilde{\boldsymbol{x}}_{s} = \tilde{\boldsymbol{x}}_{c} + w\boldsymbol{A}^{\top} (\boldsymbol{A}\boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A}\tilde{\boldsymbol{x}}_{c}) , \qquad (E.5)$$

$$\tilde{\sigma}_{\rm s}^2 = \tilde{\sigma}_{\rm c}^2(w-1) , \qquad (E.6)$$

because it involves only the inversion of an $M \times M$ matrix, instead of an $N \times N$ matrix, which is beneficial, since M < N. Furthermore, one can simplify the algorithm by utilizing the SVD of sensing matrix \boldsymbol{A} for the computation of the inverse. Let the SVD be given by $\boldsymbol{A} = \boldsymbol{U} \operatorname{diag}(s_i) \boldsymbol{V}^{\top}$, where $\boldsymbol{U} \in \mathbb{R}^{M \times M}$ is an orthonormal matrix (satisfying $\boldsymbol{U}^{-1} = \boldsymbol{U}^{\top}$), $s_1 \geq s_2 \geq \cdots \geq s_M > 0$, and $\boldsymbol{V} \in \mathbb{R}^{N \times M}$ is a matrix with orthonormal column vectors, which means $\boldsymbol{V}^{\top} \boldsymbol{V} = \boldsymbol{I}_M$. Then, the estimate $\tilde{\boldsymbol{x}}_s$ can be written as

$$\tilde{\boldsymbol{x}}_{s} = \tilde{\boldsymbol{x}}_{c} + w \boldsymbol{V} \operatorname{diag}(s_{i}) \left(\operatorname{diag}(s_{i}^{2}) + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M} \right)^{-1} \left(\boldsymbol{U}^{\top} \boldsymbol{y} - \operatorname{diag}(s_{i}) \boldsymbol{V}^{\top} \tilde{\boldsymbol{x}}_{c} \right), \quad (E.7)$$

where w is computed as

$$w = \left(\frac{1}{N} \operatorname{trace}\left(\boldsymbol{V}\operatorname{diag}(s_i) \left(\operatorname{diag}(s_i^2) + \frac{\sigma_n^2}{\tilde{\sigma}_c^2}\boldsymbol{I}_M\right)^{-1} \operatorname{diag}(s_i)\boldsymbol{V}^{\top}\right)\right)^{-1} .$$
 (E.8)

The important simplification lies in the fact, that the matrix inverse is replaced by M scalar inverses, since the matrix to be inverted is diagonal. For details see [RSF19, Spa19].

E.6 VAMP with Individual Variances (VAMPind)

This algorithm results from VAMP, when simply replacing the averaged variances by individual ones without changing anything else in the processing. It was first proposed in [FSRS16]. For comparability, the same notation as in the VAMP algorithm above is used. Mathematically the same but better for the implementation of the LMMSE is the processing used in Algorithm E.7 below.

Algorithm E.6: $m{m}_{ m s}= extsf{VAMPind}(m{y},m{A},\sigma_{ extsf{n}}^2)$		
1 $ ilde{m{ extbf{\Phi}}}_{ ext{c}} = rac{s}{N} m{I}_N, \; ilde{m{x}}_{ ext{c}} = m{0}$ // initialization		
² while stopping criterion not met do		
3 $\boldsymbol{m}_{\mathrm{c}} = \tilde{\boldsymbol{x}}_{\mathrm{c}} + (\boldsymbol{A}^{\top}\boldsymbol{A} + \sigma_{n}^{2}\tilde{\boldsymbol{\varPhi}}_{\mathrm{c}}^{-1})^{-1}\boldsymbol{A}^{\top}(\boldsymbol{y} - \boldsymbol{A}\tilde{\boldsymbol{x}}_{\mathrm{c}})$	// LMMSE	
4 $\boldsymbol{\varPhi}_{\mathrm{c}} = \sigma_{n}^2 (\boldsymbol{A}^{ op} \boldsymbol{A} + \sigma_{n}^2 \tilde{\boldsymbol{\varPhi}}_{\mathrm{c}}^{-1})^{-1}$		
5 for $i \in \{1,, N\}$ do		
$6 \sigma_{ci}^2 = [\boldsymbol{\Phi}_c]_{ii}$		
7 $\tilde{\sigma}_{s,j}^{2} = \left(\frac{1}{\sigma_{c,j}^{2}} - \frac{1}{\tilde{\sigma}_{c,j}^{2}}\right)^{-1}$	<pre>// unbiasing</pre>	
8 $\tilde{x}_{\mathrm{s},j} = \tilde{\sigma}_{\mathrm{s},j}^2 \left(\frac{m_{\mathrm{c},j}}{\sigma_{\mathrm{c},j}^2} - \frac{\tilde{x}_{\mathrm{c},j}}{\tilde{\sigma}_{\mathrm{c},j}^2} \right)$		
9 $m_{\mathrm{s},j} = \mathrm{E}_{x}\{x \mid \tilde{x}_{\mathrm{s},j}, \tilde{\sigma}_{\mathrm{s},j}^2\}$	// NLMMSE	
10 $\sigma_{s,j}^2 = E_x \{ (x - m_{s,j})^2 \tilde{x}_{s,j}, \tilde{\sigma}_{s,j}^2 \}$		
11 $\tilde{\sigma}_{\mathrm{c},j}^2 = (\frac{1}{\sigma_{\mathrm{s},j}^2} - \frac{1}{\tilde{\sigma}_{\mathrm{s},j}^2})^{-1}$, $[\tilde{\boldsymbol{\Phi}}_{\mathrm{c}}]_{jj} = \tilde{\sigma}_{\mathrm{c},j}^2$	// unbiasing	
12 $\qquad \qquad \qquad$		
13 $\boldsymbol{m}_{\mathrm{s}} = [m_{\mathrm{s},1},\ldots,m_{\mathrm{s},N}]^{ op}$		

E.7 VAMPire

The Algorithm E.6 shows significant drawbacks in the processing of the variances, since it allows for negative variances. By an estimation-theoretic point of view, this treatment can be corrected and leads to, what we call VAMPire; for the derivation see Sec. 4.4.2.1. The algorithm was first published in [FSG20].

Algorithm E.7: $m{m}_{
m s}=$ VAMPire $(m{y},\,m{A},\,\sigma_{
m n}^2)$ 1 $\tilde{\boldsymbol{\Phi}}_{\mathrm{c}} = rac{s}{N} \boldsymbol{I}_{N}, \; \tilde{\boldsymbol{x}}_{\mathrm{c}} = \boldsymbol{0}$ // initialization ² while stopping criterion not met do $\boldsymbol{K} = \boldsymbol{\tilde{\tilde{\boldsymbol{\Phi}}}_{\mathrm{c}}}\boldsymbol{A}^{\top} (\boldsymbol{A}\boldsymbol{\tilde{\boldsymbol{\Phi}}}_{\mathrm{c}}\boldsymbol{A}^{\top} + \sigma_{\mathrm{n}}^{2}\boldsymbol{I}_{N})^{-1}\boldsymbol{A}$ 3 $[\mathbf{W}]_{jj} = 1/[\mathbf{K}]_{jj} \qquad j \in \{1, \dots, N\}$ $\tilde{\mathbf{x}}_{s} = \tilde{\mathbf{x}}_{c} + \mathbf{W}\tilde{\mathbf{\Phi}}_{c}\mathbf{A}^{\top}(\mathbf{A}\tilde{\mathbf{\Phi}}_{c}\mathbf{A}^{\top} + \sigma_{n}^{2}\mathbf{I}_{N})^{-1}(\mathbf{y} - \mathbf{A}\tilde{\mathbf{x}}_{c})$ 4 // unbiased LMMSE 5 $\begin{array}{l} \mathbf{for} \; j \in \{1, \, \dots, \, N\} \; \mathbf{do} \\ \big| \; \; \tilde{\sigma}_{\mathrm{s},j}^2 = \tilde{\sigma}_{\mathrm{c},j}^2([\boldsymbol{W}]_{jj} - 1) \end{array}$ 6 7 // NLMMSE 8 9 // $\tilde{\mathbf{x}}_{s} = \mathbf{x} + \mathbf{n}_{s}, \ \mathbf{n}_{s} \sim \mathcal{N}(0, \tilde{\sigma}_{s,i}^{2})$ 10 $\begin{bmatrix} \tilde{\sigma}_{c,j}^2 = \sigma_{s,j}^2 + \left(\frac{\sigma_j^2}{\tilde{\sigma}_{s,j}^2 - \sigma_j^2} (m_{s,j} - \tilde{x}_{s,j})\right)^2 \\ \tilde{x}_{c,j} = (1/\sigma_j^2 - 1/\tilde{\sigma}_{s,j}^2)^{-1} (\frac{m_{s,j}}{\sigma_j^2} - \frac{\tilde{x}_{s,j}}{\tilde{\sigma}_{s,j}^2}) \end{bmatrix}$ // unbiasing 11 12 $\tilde{\boldsymbol{\Phi}}_{\mathrm{c}} = \mathbf{diag}(\tilde{\sigma}_{\mathrm{c},j}^2), \, \tilde{\boldsymbol{x}}_{\mathrm{c}} = [\tilde{x}_{\mathrm{c},1}, \, \dots, \, \tilde{x}_{\mathrm{c},N}]$ 13 14 $\boldsymbol{m}_{\mathrm{s}} = [m_{\mathrm{s},1}, \, \ldots, \, m_{\mathrm{s},N}]^{\top}$

E.8 VAMP with Individually Unbiased LMMSE Estimator

The following algorithms result from VAMP by replacing the average bias compensation that is performed after the channel-constrained estimation with individual unbiasing. They differ in the treatment of the variances before the signal-constrained estimation. VAMPia averages the variances before the signal-constrained estimation; VAMPii uses the individual variances directly. The algorithms were first published in [SF22].

	gorithm E.8: $m{m}_{ m s}=$ VAMPia $(m{y},m{A},\sigma_{ m n}^{ m z})$		
ı õ	1 $ ilde{\sigma}_{ m c}^2=s/N,\; ilde{m{x}}_{ m c}=m{0}$ // initialization		
2 V	while stopping criterion not met do		
3	$oldsymbol{K} = oldsymbol{A}^ op (oldsymbol{A}oldsymbol{A}^ op + rac{\sigma_n^2}{ ilde\sigma_c^2}oldsymbol{I}_M)^{-1}oldsymbol{A}$		
4	$oldsymbol{W} = \mathbf{diag}(1/[oldsymbol{K}]_{jj})$		
5	$ ilde{oldsymbol{x}}_{ ext{s}} = ilde{oldsymbol{x}}_{ ext{c}} + oldsymbol{W}oldsymbol{A}^ op (oldsymbol{A}oldsymbol{A}^ op + rac{\sigma_n^2}{ ilde{\sigma}_z^2}oldsymbol{I}_M)^{-1}(oldsymbol{y} - oldsymbol{A} ilde{oldsymbol{x}}_{ ext{c}})$	//	unbiased LMMSE
6	$ ilde{\sigma}_{ m s}^2 = ilde{\sigma}_{ m c}^2 (rac{1}{N} \sum_{j=1}^N rac{1}{ m{K} _{ij}} - 1)$		
7	$oldsymbol{m}_{\mathrm{s}} = \mathrm{E}_{x}\{x \mid ilde{x}_{\mathrm{s}}, ilde{\sigma}_{\mathrm{s}}^2\}$		// NLMMSE
8	$\sigma_{\rm s}^2 = \frac{1}{N} \sum_{j=1}^{N} E_{\sf x} \{ ({\sf x} - m_{{ m s},j})^2 \mid \tilde{x}_{{ m s},j}, \tilde{\sigma}_{ m s}^2 \}$		
_	~ 2 (1 1)-1		//
9	$O_{\mathbf{c}} = \left(\frac{\sigma_{\mathbf{s}}^2}{\sigma_{\mathbf{s}}^2} - \frac{\tilde{\sigma}_{\mathbf{s}}^2}{\tilde{\sigma}_{\mathbf{s}}^2}\right)^{-1}$		// unblasing
10	$\sum egin{array}{c} m{x}_{ m c} = \sigma_{ m c}^2 (rac{m_{ m s}}{\sigma_{ m s}^2} - rac{\omega_{ m s}}{ ilde{\sigma}_{ m s}^2}) \end{array}$		
	sorithm F 9: $m - VAMPii(u \ d \ \sigma^2)$		
	\mathbf{g}_{S}		
1 ã	$\tilde{c}^2 - a/N \tilde{c} = 0$		initialization
1 õ	$\hat{x}_{c}^{2} = s/N, \ \tilde{x}_{c} = 0$	//	initialization
1 õ 2 W	$\tilde{\boldsymbol{x}}_{c}^{2} = s/N, \ \tilde{\boldsymbol{x}}_{c} = \boldsymbol{0}$ while stopping criterion not met do $\boldsymbol{K} = \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\sigma_{n}} \boldsymbol{I}_{M})^{-1} \boldsymbol{A}$	//	initialization
1 $\tilde{\sigma}$ 2 W 3	$\mathbf{\hat{x}}_{c}^{2} = s/N, \ \tilde{\mathbf{x}}_{c} = 0$ while stopping criterion not met do $\mathbf{K} = \mathbf{A}^{\top} (\mathbf{A}\mathbf{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}}\mathbf{I}_{M})^{-1}\mathbf{A}$ $\mathbf{W} = \operatorname{diag}(1/[\mathbf{K}]_{W})$	//	initialization
1 Õ 2 W 3 4	$\mathbf{\hat{x}}_{c}^{2} = s/N, \ \tilde{\mathbf{x}}_{c} = 0$ while stopping criterion not met do $\mathbf{K} = \mathbf{A}^{\top} (\mathbf{A}\mathbf{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}}\mathbf{I}_{M})^{-1}\mathbf{A}$ $\mathbf{W} = \operatorname{diag}(1/[\mathbf{K}]_{jj})$ $\tilde{\mathbf{x}}_{m} = \tilde{\mathbf{x}}_{m} + \mathbf{W}\mathbf{A}^{\top} (\mathbf{A}\mathbf{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}}\mathbf{I}_{m})^{-1}(\mathbf{x}_{m} - \mathbf{A}\tilde{\mathbf{x}}_{m})$	//	initialization
1 Õ 2 W 3 4 5	$\begin{split} \mathbf{\tilde{x}}_{c}^{2} &= s/N, \ \tilde{\boldsymbol{x}}_{c} = 0 \\ \text{while stopping criterion not met do} \\ \mathbf{K} &= \mathbf{A}^{\top} (\mathbf{A}\mathbf{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \mathbf{I}_{M})^{-1} \mathbf{A} \\ \mathbf{W} &= \text{diag}(1/[\mathbf{K}]_{jj}) \\ \tilde{\boldsymbol{x}}_{s} &= \tilde{\boldsymbol{x}}_{c} + \mathbf{W}\mathbf{A}^{\top} (\mathbf{A}\mathbf{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \mathbf{I}_{M})^{-1} (\boldsymbol{y} - \mathbf{A}\tilde{\boldsymbol{x}}_{c}) \\ \tilde{\boldsymbol{x}}_{s}^{2} &= \tilde{\boldsymbol{x}}_{c}^{2} (1/[\mathbf{K}] = 1) \end{split}$	//	initialization unbiased LMMSE
$\begin{array}{c} 1 \widetilde{O} \\ 2 \mathbf{W} \\ 3 \\ 4 \\ 5 \\ 6 \end{array}$	$\begin{split} \mathbf{\hat{x}}_{c}^{2} &= s/N, \ \tilde{\boldsymbol{x}}_{c} = \boldsymbol{0} \\ \textbf{while stopping criterion not met do} \\ \boldsymbol{K} &= \boldsymbol{A}^{\top} (\boldsymbol{A}\boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}}\boldsymbol{I}_{M})^{-1}\boldsymbol{A} \\ \boldsymbol{W} &= \textbf{diag}(1/[\boldsymbol{K}]_{jj}) \\ \tilde{\boldsymbol{x}}_{s} &= \tilde{\boldsymbol{x}}_{c} + \boldsymbol{W}\boldsymbol{A}^{\top} (\boldsymbol{A}\boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}}\boldsymbol{I}_{M})^{-1}(\boldsymbol{y} - \boldsymbol{A}\tilde{\boldsymbol{x}}_{c}) \\ \tilde{\sigma}_{s,j}^{2} &= \tilde{\sigma}_{c}^{2}(1/[\boldsymbol{K}]_{jj} - 1) \qquad j \in \{1, \ldots, \end{split}$	// // N}	initialization unbiased LMMSE
$\begin{array}{c} 1 \widetilde{O} \\ 2 \mathbf{W} \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \end{array}$	$\begin{split} \mathbf{\tilde{x}}_{c}^{2} &= s/N, \ \tilde{\boldsymbol{x}}_{c} = 0 \\ \mathbf{\tilde{x}}_{c}^{2} &= \mathbf{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} \boldsymbol{A} \\ \mathbf{W} &= \mathbf{diag} (1/[\boldsymbol{K}]_{jj}) \\ \tilde{\boldsymbol{x}}_{s} &= \tilde{\boldsymbol{x}}_{c} + \boldsymbol{W} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \tilde{\boldsymbol{x}}_{c}) \\ \tilde{\sigma}_{s,j}^{2} &= \tilde{\sigma}_{c}^{2} (1/[\boldsymbol{K}]_{jj} - 1) \qquad j \in \{1, \dots, n\} \\ m_{c,i} &= \operatorname{E}_{v} \{ \mathbf{x} \mid \tilde{\boldsymbol{x}}_{n,i}, \ \tilde{\boldsymbol{\sigma}}^{2}_{s} \} \qquad i \in \{1, \dots, n\} \end{split}$	// N} N}	initialization unbiased LMMSE
$\begin{array}{c} 1 \widetilde{O} \\ 2 \mathbf{W} \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \end{array}$	$\begin{split} \mathbf{\tilde{c}}_{c}^{2} &= s/N, \ \tilde{\boldsymbol{x}}_{c} = 0 \\ \textbf{while stopping criterion not met do} \\ \boldsymbol{K} &= \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} \boldsymbol{A} \\ \boldsymbol{W} &= \textbf{diag}(1/[\boldsymbol{K}]_{jj}) \\ \tilde{\boldsymbol{x}}_{s} &= \tilde{\boldsymbol{x}}_{c} + \boldsymbol{W} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \tilde{\boldsymbol{x}}_{c}) \\ \tilde{\sigma}_{s,j}^{2} &= \tilde{\sigma}_{c}^{2} (1/[\boldsymbol{K}]_{jj} - 1) \qquad j \in \{1, \dots, n\} \\ m_{s,j} &= \mathbf{E}_{s} \{ \mathbf{x} \mid \tilde{\boldsymbol{x}}_{s,j}, \ \tilde{\sigma}_{s,j}^{2} \} \qquad j \in \{1, \dots, n\} \\ \sigma_{s}^{2} &= \frac{1}{2^{1}} \sum_{i=1}^{N} \mathbf{E}_{s} \{ (\mathbf{x} - m_{s,i})^{2} \mid \tilde{\boldsymbol{x}}_{s,i}, \ \tilde{\sigma}_{s,j}^{2} \} \end{split}$	// N} N}	initialization unbiased LMMSE // NLMMSE
$\begin{array}{c} 1 \widetilde{\mathcal{O}} \\ 2 \mathbf{W} \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \end{array}$	$\begin{split} & \overset{(\mathbf{r}^{*}) \to \mathbf{W}}{\overset{(\mathbf{r}^{*}) \to \mathbf{W}}}} \\ & \mathbf{K} = \mathbf{A}^{\top} (\mathbf{A} \mathbf{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \mathbf{I}_{M})^{-1} \mathbf{A} \\ & \mathbf{W} = \mathbf{diag}(1/[\mathbf{K}]_{jj}) \\ & \tilde{\mathbf{x}}_{s} = \tilde{\mathbf{x}}_{c} + \mathbf{W} \mathbf{A}^{\top} (\mathbf{A} \mathbf{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \mathbf{I}_{M})^{-1} (\mathbf{y} - \mathbf{A} \tilde{\mathbf{x}}_{c}) \\ & \tilde{\sigma}_{s,j}^{2} = \tilde{\sigma}_{c}^{2} (1/[\mathbf{K}]_{jj} - 1) & j \in \{1, \dots, \\ & m_{s,j} = \mathbf{E}_{s} \{\mathbf{x} \mid \tilde{\mathbf{x}}_{s,j}, \tilde{\sigma}_{s,j}^{2}\} & j \in \{1, \dots, \\ & m_{s,j} = \mathbf{E}_{s} \{\mathbf{x} \mid \tilde{\mathbf{x}}_{s,j}, \tilde{\sigma}_{s,j}^{2}\} & j \in \{1, \dots, \\ & \sigma_{s}^{2} = \frac{1}{N} \sum_{j=1}^{N} \mathbf{E}_{s} \{(\mathbf{x} - m_{s,j})^{2} \mid \tilde{\mathbf{x}}_{s,j}, \tilde{\sigma}_{s,j}^{2}\} \\ & \tilde{\sigma}_{s}^{2} = \frac{1}{N} \sum_{i=1}^{N} \tilde{\sigma}_{s,i}^{2} \end{split}$	// N} N}	initialization unbiased LMMSE // NLMMSE
1 Õ 2 W 3 4 5 6 7 8 9	$\begin{split} \mathbf{\hat{c}}_{c}^{2} &= s/N, \ \tilde{\boldsymbol{x}}_{c} = 0 \\ \mathbf{\hat{v}} \mathbf{hile \ stopping \ criterion \ not \ met \ do} \\ & \boldsymbol{K} = \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} \boldsymbol{A} \\ & \boldsymbol{W} = \mathbf{diag}(1/[\boldsymbol{K}]_{jj}) \\ & \tilde{\boldsymbol{x}}_{s} &= \tilde{\boldsymbol{x}}_{c} + \boldsymbol{W} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \tilde{\boldsymbol{x}}_{c}) \\ & \tilde{\sigma}_{s,j}^{2} &= \tilde{\sigma}_{c}^{2} (1/[\boldsymbol{K}]_{jj} - 1) \qquad j \in \{1, \dots, n, n_{s,j} = \mathbf{E}_{s} \{\mathbf{x} \mid \tilde{\boldsymbol{x}}_{s,j}, \tilde{\sigma}_{s,j}^{2}\} \qquad j \in \{1, \dots, n_{s,j}\} \\ & \sigma_{s}^{2} &= \frac{1}{N} \sum_{j=1}^{N} \mathbf{E}_{s} \{(\mathbf{x} - m_{s,j})^{2} \mid \tilde{\boldsymbol{x}}_{s,j}, \tilde{\sigma}_{s,j}^{2}\} \\ & \tilde{\sigma}_{s}^{2} &= \frac{1}{N} \sum_{j=1}^{N} \tilde{\sigma}_{s,j}^{2} \end{split}$	// N} N}	initialization unbiased LMMSE // NLMMSE
1 $\tilde{\sigma}$ 2 W 3 4 5 6 7 8 9 10	$\begin{split} \tilde{\boldsymbol{c}}_{c}^{2} &= s/N, \ \tilde{\boldsymbol{x}}_{c} = \boldsymbol{0} \\ \text{while stopping criterion not met do} \\ \boldsymbol{K} &= \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} \boldsymbol{A} \\ \boldsymbol{W} &= \operatorname{diag}(1/[\boldsymbol{K}]_{jj}) \\ \tilde{\boldsymbol{x}}_{s} &= \tilde{\boldsymbol{x}}_{c} + \boldsymbol{W} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \tilde{\boldsymbol{x}}_{c}) \\ \tilde{\sigma}_{s,j}^{2} &= \tilde{\sigma}_{c}^{2} (1/[\boldsymbol{K}]_{jj} - 1) \qquad j \in \{1, \dots, n, n\} \\ m_{s,j} &= \operatorname{E}_{x} \{ \mathbf{x} \mid \tilde{\boldsymbol{x}}_{s,j}, \ \tilde{\sigma}_{s,j}^{2} \} \qquad j \in \{1, \dots, n\} \\ \sigma_{s}^{2} &= \frac{1}{N} \sum_{j=1}^{N} \operatorname{E}_{x} \{ (\mathbf{x} - m_{s,j})^{2} \mid \tilde{\boldsymbol{x}}_{s,j}, \ \tilde{\sigma}_{s,j}^{2} \} \\ \tilde{\sigma}_{s}^{2} &= (\frac{1}{\sigma_{s}^{2}} - \frac{1}{\tilde{\sigma}_{s,j}^{2}})^{-1} \end{split}$	// N} N}	initialization unbiased LMMSE // NLMMSE // unbiasing
1 $\tilde{\sigma}$ 2 W 3 4 5 6 7 8 9 10 11	$\begin{split} \vec{c}_{c} &= s/N, \ \tilde{\boldsymbol{x}}_{c} = \boldsymbol{0} \\ \text{while stopping criterion not met do} \\ & \boldsymbol{K} = \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} \boldsymbol{A} \\ & \boldsymbol{W} = \text{diag}(1/[\boldsymbol{K}]_{jj}) \\ & \tilde{\boldsymbol{x}}_{s} = \tilde{\boldsymbol{x}}_{c} + \boldsymbol{W} \boldsymbol{A}^{\top} (\boldsymbol{A} \boldsymbol{A}^{\top} + \frac{\sigma_{n}^{2}}{\tilde{\sigma}_{c}^{2}} \boldsymbol{I}_{M})^{-1} (\boldsymbol{y} - \boldsymbol{A} \tilde{\boldsymbol{x}}_{c}) \\ & \tilde{\sigma}_{s,j}^{2} = \tilde{\sigma}_{c}^{2} (1/[\boldsymbol{K}]_{jj} - 1) \qquad j \in \{1, \dots, n, n\} \\ & \boldsymbol{m}_{s,j} = \mathbf{E}_{s} \{ \mathbf{x} \mid \tilde{\boldsymbol{x}}_{s,j}, \ \tilde{\sigma}_{s,j}^{2} \} \qquad j \in \{1, \dots, n\} \\ & \sigma_{s}^{2} = \frac{1}{N} \sum_{j=1}^{N} \mathbf{E}_{s} \{ (\mathbf{x} - m_{s,j})^{2} \mid \tilde{\boldsymbol{x}}_{s,j}, \ \tilde{\sigma}_{s,j}^{2} \} \\ & \tilde{\sigma}_{s}^{2} = (\frac{1}{\sigma_{s}^{2}} - \frac{1}{\tilde{\sigma}_{s}^{2}})^{-1} \\ & \tilde{\boldsymbol{x}}_{c} = \tilde{\sigma}_{c}^{2} (\frac{m_{s}}{\sigma_{s}^{2}} - \frac{\tilde{\boldsymbol{x}}_{s}}{\tilde{\sigma}_{s}^{2}}) \end{split}$	// N} N}	<pre>initialization unbiased LMMSE // NLMMSE // unbiasing</pre>

E.9 VAMPind with Sequentially Processed Variables

The usage of individual variances enables the possibility for a sequential processing of the signal components to be recovered. The algorithm seqVAMPind is therefore a sequential version (in terms of processing the variables) of VAMPind. The idea for the algorithm was given by the algorithm stated in [OW05, Appendix D.], which differs from the algorithm here by the utilization of a binary prior. In [BMPP20, Algorithm 1], the algorithm is stated for Bernoulli-Gaussian prior.

Algorithm E.10: $m{m}_{ m s}={\tt seqVAMPind}(m{y},m{A},\sigma_{\sf n}^2)$			
1 /	1 $\boldsymbol{\Lambda}_{\mathrm{c}} = [\frac{N}{s}, \ldots, \frac{N}{s}]_{N}^{T}, \ \boldsymbol{\lambda}_{\mathrm{c}} = \boldsymbol{0}$ // initialization		
2 4	2 $\boldsymbol{\varPhi}_{\mathrm{c}} = (\frac{1}{\sigma_{\mathrm{c}}^2} \boldsymbol{A}^{\top} \boldsymbol{A} + \boldsymbol{\Lambda}_{\mathrm{c}})^{-1} = [\boldsymbol{\phi}_{\mathrm{c},1}, \dots, \boldsymbol{\phi}_{\mathrm{c},N}], \sigma_{\mathrm{c},j}^2 = [\boldsymbol{\varPhi}_{\mathrm{c}}]_{jj}$ // LMMSE		
3 N	$oldsymbol{n}_{ m c} = oldsymbol{\varPhi}_{ m c}^{''} (rac{1}{\sigma_{ m n}^2}oldsymbol{A}^{ op}oldsymbol{y} + oldsymbol{\lambda}_{ m c})$		
4 V	while stopping criterion not met do		
5	Choose index set $\mathcal{J}\subseteq\{1,\ldots,N\}$ with $ \mathcal{J} =T$		
6	Specify sequence $s_{\mathcal{J}} = [j_1, j_2, \ldots, j_T] : j_i \in \mathcal{J} \ \forall i \in \{1, \dots, j_T\}$	$\ldots, T\}$	
7	for $j = s_{\mathcal{J}}$ do		
8	$ ilde{\sigma}_{\mathrm{s},j}^2 = (1/\sigma_{\mathrm{c},j}^2 - \Lambda_{\mathrm{c},j})^{-1}$	<pre>// unbiasing</pre>	
9	$ ilde{x}_{\mathrm{s},j} = ilde{\sigma}_{\mathrm{s},j}^2(m_{\mathrm{c},j}/\sigma_{\mathrm{c},j}^2 - \lambda_{\mathrm{c},j})$		
10	$m_{\mathbf{s},j} = \mathbf{E}_{x} \{ x \mid \tilde{x}_{\mathbf{s},j}, \tilde{\sigma}_{\mathbf{s},j}^2 \},$	// NLMMSE	
11	$\sigma_{\mathbf{s},j}^2 = \mathbf{E}_{x}\{(x - m_{\mathbf{s},j})^2 \mid \tilde{x}_{\mathbf{s},j}, \tilde{\sigma}_{\mathbf{s},j}^2\}$		
12	$\lambda_{c,i} = m_{s,i} / \sigma_{c,i}^2 - \tilde{x}_{s,i} / \tilde{\sigma}_{c,i}^2$	<pre>// unbiasing</pre>	
13	$\Lambda_{{ m c},j} = 1/\sigma_{{ m s},j}^2 - 1/\tilde{\sigma}_{{ m s},j}^2$		
14	$\boldsymbol{m} - \boldsymbol{m} + \frac{1}{2} (m \cdot - m \cdot) \boldsymbol{\phi}$	// rank_one undate	
14	$m_{\rm c} = m_{\rm c} + \frac{1}{\sigma_{\rm c,j}^2} (m_{\rm s,j} - m_{\rm c,j}) \psi_{\rm c,j}$	// Tank-one update	
15	$oldsymbol{\Phi}_{\mathrm{c}} = oldsymbol{\Phi}_{\mathrm{c}} + rac{1}{(\sigma_{\mathrm{c},j}^2)^2} \left(\sigma_{\mathrm{s},j}^2 - \sigma_{\mathrm{c},j}^2 ight) oldsymbol{\phi}_{\mathrm{c},j} oldsymbol{\phi}_{\mathrm{c},j}^{+}$		
16	for $j' = 1,, N$ do		
17	$\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ $		
18 N	$ar{m{n}}_{ m s} = [m_{{ m s},1},\ldots,m_{{ m s},N}]^{ op}$		

E.10 VAMPire with Sequentially Processed Variables

The algorithm shown here improves upon seqVAMPind by employing estimation-theoretic bias compensation after the signal-constrained estimation, i.e., it is the sequential version of VAMPire. It was introduced by the author in [SF21].

Algorithm E.11: $m_{\rm s} = \operatorname{seqVAMPire}(\boldsymbol{y}, \boldsymbol{A}, \sigma_{\rm n}^2)$ 1 $\boldsymbol{\Lambda}_{\rm c} = [\frac{N}{s}, \ldots, \frac{N}{s}]_N^{\top}, \boldsymbol{\lambda}_{\rm c} = \boldsymbol{0}, // \text{ initialization}$ 2 $\boldsymbol{\Phi}_{\rm c} = (\frac{1}{\sigma_{\rm n}^2} \boldsymbol{A}^{\top} \boldsymbol{A} + \boldsymbol{\Lambda}_{\rm c})^{-1} = [\boldsymbol{\phi}_{{\rm c},1}, \ldots, \boldsymbol{\phi}_{{\rm c},N}], \sigma_{{\rm c},j}^2 = [\boldsymbol{\Phi}_{\rm c}]_{jj} \ j \in \{1, \ldots, N\} // \text{ LMMSE}$ з $\boldsymbol{m}_{\mathrm{c}} = \boldsymbol{\varPhi}_{\mathrm{c}}(rac{1}{\sigma^{2}}\boldsymbol{A}^{ op}\boldsymbol{y} + \boldsymbol{\lambda}_{\mathrm{c}})$ 4 while stopping criterion not met do Choose index set $\mathcal{J} \subseteq \{1, \ldots, N\}$ with $|\mathcal{J}| = T$ 5 Specify sequence $s_{\mathcal{J}} = [j_1, j_2, \dots, j_T] : j_i \in \mathcal{J} \ \forall i \in \{1, \dots, T\}$ 6 for $j = s_{\mathcal{J}} \operatorname{do}$ 7 $\tilde{\sigma}_{{\rm s},j}^2 = (1/\sigma_{{\rm c},j}^2 - \Lambda_{{\rm c},j})^{-1}$ // unbiasing 8 $\tilde{x}_{\mathrm{s},j} = \tilde{\sigma}_{\mathrm{s},j}^2 (m_{\mathrm{c},j} / \sigma_{\mathrm{c},j}^2 - \lambda_{\mathrm{c},j})$ 9 $ilde{x}_{\mathrm{s},j} = \lambda_{\mathrm{s},j} / \Lambda_{\mathrm{s},j}, \, ilde{\sigma}_{\mathrm{s},j}^2 = 1 / \Lambda_{\mathrm{s},j}$ 10 $m_{\mathbf{s},j} = \mathbf{E}_{\mathsf{x}} \{ \mathsf{x} \mid \tilde{x}_{\mathbf{s},j}, \, \tilde{\sigma}_{\mathbf{s},j}^2 \},$ // NLMMSE 11 $\sigma_{\mathbf{s},j}^2 = \mathbf{E}_{\mathsf{x}}\{(\mathbf{x} - m_{\mathbf{s},j})^2 \mid \tilde{x}_{\mathbf{s},j}, \, \tilde{\sigma}_{\mathbf{s},i}^2\}$ 12 // $\tilde{\mathbf{x}}_{\mathbf{s}} = \mathbf{x} + \mathbf{n}_{\mathbf{s}}, \, \mathbf{n}_{\mathbf{s}} \sim \mathcal{N}(0, \tilde{\sigma}_{\mathbf{s},i}^2)$ $\sigma_j^2 = \mathcal{E}_{\tilde{x}_s} \{ \sigma_{s,j}^2 \}$ 13
$$\begin{split} \tilde{\sigma}_{\text{c},j}^2 &= \sigma_{\text{s},j}^2 + (\frac{\sigma_j^2}{\tilde{\sigma}_{\text{s},j}^2 - \sigma_j^2} (m_{\text{s},j} - \tilde{x}_{\text{s},j}))^2 \\ \tilde{\sigma}_j^2 &= (1/\sigma_j^2 - 1/\tilde{\sigma}_{\text{s},j}^2)^{-1} \end{split}$$
14 15 $\Delta \Lambda_{\mathrm{c},j} = 1/\tilde{\sigma}_{\mathrm{c},j}^2 - \Lambda_{\mathrm{c},j}^2$ 16
$$\begin{split} \Lambda_{\mathrm{c},j} &= 1/\tilde{\sigma}_{\mathrm{c},j}^2 \\ \Delta\lambda_{\mathrm{c},j} &= \tilde{\sigma}_j^2/\tilde{\sigma}_{\mathrm{c},j}^2 \cdot (m_{\mathrm{s},j}/\sigma_j^2 - \lambda_{\mathrm{s},j}) - \lambda_{\mathrm{c},j} \end{split}$$
17 18 $\lambda_{\mathrm{c},j} = ilde{\sigma}_j^2 / ilde{\sigma}_{\mathrm{c},j}^2 \cdot (m_{\mathrm{s},j} / \sigma_j^2 - \lambda_{\mathrm{s},j})$ $m_{\mathrm{c}} = m_{\mathrm{c}} + rac{\Delta\lambda_{\mathrm{c},j} - \Delta\Lambda_{\mathrm{c},j}m_{\mathrm{c},j}}{1 + \Delta\Lambda_{\mathrm{c},j}\sigma_{\mathrm{c},j}^2} \phi_{\mathrm{c},j}$ $oldsymbol{\Phi}_{\mathrm{c}} = oldsymbol{\Phi}_{\mathrm{c}} - rac{\Delta\Lambda_{\mathrm{c},j}}{1 + \Delta\Lambda_{\mathrm{c},j}\sigma_{\mathrm{c},j}^2} \phi_{\mathrm{c},j} \phi_{\mathrm{c},j}^{\top}$ 19 // rank-one update 20 21 for j' = 1, ..., N do 22 $\sigma_{\mathbf{c},j'}^2 = [\boldsymbol{\Phi}_{\mathbf{c}}]_{j'j'}$ 23 24 $\boldsymbol{m}_{\mathrm{s}} = [m_{\mathrm{s},1}, \ldots, m_{\mathrm{s},N}]^{\top}$

E.11 Sequential Processing of Observations

Here, the algorithm based on VAMP, that processes the observations in a sequential way, is stated. As in VAMP, an average variance is used to represent the reliability of the entire vector. Note that a_i is the *i*th row of sensing matrix A.

The initialization is stated as in the implementation used for the simulations.

Algorithm E.12:
$$m = \operatorname{seqVAMPobs}(y, A, \sigma_n^2)$$

1 for $i = s, 1, ..., M$ do
2 $\lfloor \tilde{x}^{(i)} = 0, \sigma_{(i)}^2 = 10^8$ // initialization
3 $m = 0, \sigma^2 = 5 \cdot 10^7$
4 while stopping criterion not met do
5 Choose $i \in \{s, 1, ..., M\}$
6 $\sigma_{\langle i}^2 = (1/\sigma^2 - 1/\sigma_{(i)}^2)^{-1}, \tilde{x}^{\langle i} = \sigma_{\langle i}^2(m/\sigma^2 - \tilde{x}^{(i)}/\sigma_{(i)}^2)$
7 if $i \in \{1, ..., M\}$ then
8 $\left| \begin{array}{c} w = \frac{1}{N} \operatorname{trace}(a_i \left(a_i^{\top} a_i + \frac{\sigma_{n,i}^2}{\sigma_{\langle i \rangle}^2}\right)^{-1} a_i^{\top})$
9 $\left| \begin{array}{c} \tilde{x}^{(i)} = \tilde{x}^{\backslash i} + \frac{1}{w} a_i \left(a_i^{\top} a_i + \frac{\sigma_{n,i}^2}{\sigma_{\langle i \rangle}^2}\right)^{-1} (y_i - a_i^{\top} \tilde{x}^{\backslash i}) // \text{ unbiased LMMSE} \right|$
10 $\sigma_{\langle i \rangle}^2 = \sigma_{\langle i}^2 (\frac{1}{w} - 1)$
11 $\sigma^2 = \left(1/\sigma_{\langle i \rangle}^2 + 1/\sigma_{\langle i \rangle}^2\right)^{-1}$
12 $\left| \begin{array}{c} m = \sigma^2 \left(\tilde{x}^{(i)}/\sigma_{\langle i \rangle}^2 + \tilde{x}^{\backslash i}/\sigma_{\langle i \rangle}^2\right) \\ \sigma^2 = \frac{1}{N} \sum_{j=1}^{N} \operatorname{E}_x\{(x - m_j)^2 \mid \tilde{x}_j^{\backslash s}, \sigma_{\langle s \rangle}^2\} \\ \sigma_{\langle s \rangle}^2 = (1/\sigma^2 - 1/\sigma_{\langle s \rangle}^{-1}, \tilde{x}^{(s)} = \sigma_{\langle s \rangle}(m/\sigma^2 - \tilde{x}^{\backslash s}/\sigma_{\langle s \rangle}^2) \end{array} \right|$

_

F. Notation

F.1 Abbreviations

Acronym	Meaning
a.k.a.	<u>a</u> lso <u>k</u> nown <u>a</u> s
AMP	<u>a</u> pproximate <u>m</u> essage <u>p</u> assing
avg.	average
GAMP	generalized <u>AMP</u>
BCJR	<u>B</u> ahl- <u>C</u> ocke-Jelinek- <u>R</u> aviv (algorithm)
BG	<u>B</u> ernoulli- <u>G</u> auss (prior)
BP	<u>b</u> elief <u>p</u> ropagation
CDMA	<u>c</u> ode- <u>d</u> ivision <u>m</u> ultiple <u>a</u> ccess
cf.	<u>c</u> on <u>f</u> eratur (compare)
CoSaMP	<u>co</u> mpressive <u>sa</u> mpling <u>m</u> atching <u>p</u> ursuit
CS	compressed sensing, also compressive sampling
dB	decibel
DT	<u>d</u> iscrete <u>t</u> ernary (prior)
EC	expectation-consistent (approximate inference framework)
e.g.	exempli gratia (for example)
EP	expectation propagation
FLOP	<u>fl</u> oating-point <u>op</u> eration
GBP	\underline{G} aussian \underline{b} elief propagation
GMP	<u>G</u> aussian <u>m</u> essage <u>p</u> assing

Acronym	Meaning
GMPpostExt	<u>GMP</u> with <u>post</u> -processed <u>ext</u> rinsic
GMPpostExtLin	<u>GMPpostExt</u> at <u>lin</u> ear estimation
i.e.	$\underline{i}d \underline{e}st$ (it is)
i.i.d.	independent and identically distributed
ind.	individual
IHT	<u>i</u> terative <u>h</u> ard <u>t</u> hresholding
ISF	<u>i</u> terative <u>s</u> oft <u>f</u> eedback
IST	<u>i</u> terative <u>s</u> oft <u>t</u> hresholding
LMMSE	<u>l</u> inear <u>m</u> inimum <u>m</u> ean- <u>s</u> quared <u>e</u> rror
LDPC	<u>l</u> ow- <u>d</u> ensity <u>p</u> arity- <u>c</u> heck
LLR	<u>l</u> og- <u>l</u> ikelihood <u>r</u> atio
LUT	<u>l</u> ook- <u>u</u> p <u>t</u> able
MAP	<u>m</u> aximum <u>a</u> -posteriori
MF	<u>m</u> atched <u>f</u> ilter
MMSE	<u>m</u> inimum <u>m</u> ean- <u>s</u> quared <u>e</u> rror
MP	<u>m</u> essage <u>p</u> assing
OAMP	<u>o</u> rthogonal <u>approximate message passing</u>
OFDM	<u>o</u> rthogonal <u>f</u> requency- <u>d</u> ivision <u>m</u> ultiplexing
OMP	<u>o</u> rthogonal <u>m</u> atching pursuit
PAR	peak-to- <u>a</u> verage power <u>r</u> atio
pdf	probablity density function
SER	<u>s</u> ymbol <u>e</u> rror <u>r</u> ate
seqVAMPind	sequential VAMPind (sequential variable processing)
seqVAMPire	<u>seq</u> uential <u>VAMPire</u> (sequential variable processing)
seqVAMPobs	sequential <u>VAMP</u> observation processing
SNR	<u>s</u> ignal-to- <u>n</u> oise <u>r</u> atio
SVD	<u>s</u> ingular <u>v</u> alue <u>d</u> ecomposition
VAMP	<u>v</u> ector <u>a</u> pproximate <u>m</u> essage <u>p</u> assing
VAMPia	<u>VAMP</u> with individual unbiased LMMSE and average variance
VAMPii	<u>VAMP</u> with individual unbiased LMMSE and individual variance
VAMPind	<u>VAMP</u> with <u>ind</u> ividual variances
VAMPire	<u>VAMP</u> with individual reliabilities enhanced

Acronym	Meaning
w.r.t.	with respect to
WSN	<u>w</u> ireless <u>s</u> ensor <u>n</u> etwork

F.2 Mathematical Symbols

Throughout the thesis, scalars are denoted by lower-case letters, e.g., x, vectors by bold ones, e.g., x, matrices by upper-case bold, e.g., X, random variables in sans-serif font, e.g., x, random vectors in sans-serif bold, e.g., x, and random matrices in upper-case bold font, e.g., X. Vectors are always defined as column vectors.

F.2.1 System Parameters

Symbol	Meaning
i	row index of sensing matrix $oldsymbol{A}$
j	column index of sensing matrix $oldsymbol{A}$
M	observation dimension
N	signal dimension
s	sparsity, number of non-zero entries
$\sigma_{\sf n}^2$	noise variance

F.2.2 System Model Variables

Symbol	Meaning
A	sensing matrix
$oldsymbol{a}_i^ op$	i th row of sensing matrix $oldsymbol{A}$
$ ilde{m{a}}_j$	j th column of sensing matrix $oldsymbol{A}$
a_{ij}	entry in sensing matrix $oldsymbol{A}$ in i th row and j th column
n	noise vector
n_i	i th element of noise vector $oldsymbol{n}$
${m x}$	signal vector
x_j	j th element of signal vector $oldsymbol{x}$
\boldsymbol{y}	observation vector
y_i	i th element of observation vector $oldsymbol{y}$

F.2.3 Estimation Variables

Symbol	Meaning
$m_{ m c}$	channel-constrained estimate / linear MMSE estimate
$m_{\mathrm{c},j}$	$j\mathrm{th}$ element of linear MMSE estimate $oldsymbol{m}_{\mathrm{c}}$
$m{m}_{ m c}^{\circ}$	updated linear MMSE estimate estimate
$m^{\circ}_{{ m c},j}$	$j { m th}$ element of updated linear MMSE estimate $oldsymbol{m}_{ m c}^{ m o}$
$oldsymbol{m}_{\mathrm{c},i}$	linear MMSE estimate based on $i {\rm th}$ observation y_i
$m_{\mathrm{c},i,j}$	$j { m th}$ element of linear MMSE estimate $oldsymbol{m}_{{ m c},i}$
$m_{ m s}$	signal-constrained estimate / non-linear MMSE estimate
$m_{\mathrm{s},j}$	j th element of $oldsymbol{m}_{ m s}$
$oldsymbol{\Phi}_{ m c}$	(conditional) covariance corresponding to $m{m}_{ m c}$
$oldsymbol{\varPhi}_{\mathrm{c},i}$	(conditional) covariance corresponding to $oldsymbol{m}_{\mathrm{c},i}$
$\boldsymbol{\phi}_{\mathrm{c},j}$	j th column of covariance matrix $oldsymbol{arPsi}_{ m c}$
$\sigma^2_{\mathrm{c},j}$	$j{\rm th}$ diagonal entry of covariance matrix $oldsymbol{\varPhi}_{ m c}$
$\sigma_{ m c}^2$	average variance corresponding to $oldsymbol{m}_{ m c}$
$\sigma^2_{{ m s},j}$	$j{\rm th}$ (conditional) variance corresponding to $m_{{\rm s},j}$
$\sigma_{ m s}^2$	average (conditional) variance corresponding to $oldsymbol{m}_{ m s}$
m_{x}	mean vector of signal \mathbf{x}
$m_{x,j}$	j th element of $m{m}_{x}$
$oldsymbol{\Phi}_{ imes}$	(co)variance matrix of signal \mathbf{x}
$\sigma^2_{x,j}$	variance corresponding to x_j
σ_{x}^2	average variance corresponding to ${f x}$
K	end-to-end cascade between $oldsymbol{x}$ and $oldsymbol{m}_{ m c}$
W	scaling matrix for bias compensation
w	scaling factor for bias compensation
$ ilde{oldsymbol{x}}_{ ext{c}}$	input mean vector for channel-constrained estimation
$ ilde{x}_{\mathrm{c},j}$	j th element of $ ilde{m{x}}_{ m c}$
$ ilde{oldsymbol{x}}_{\mathrm{s}}$	input mean vector for signal-constrained estimation
$\tilde{x}_{\mathrm{s},j}$	j th element of $ ilde{m{x}}_{ m s}$
$oldsymbol{ ilde{oldsymbol{\Phi}}}_{ ext{c}}$	input (co)variance matrix corresponding to $\tilde{x}_{\mathrm{c},j}$
$ ilde{\sigma}^2_{\mathrm{c},j}$	$j \mathrm{th}$ input variance corresponding to $\tilde{x}_{\mathrm{c},j}$
$ ilde{\sigma}_{ m c}^2$	average input variance corresponding to $ ilde{m{x}}_{ ext{c}}$

Symbol	Meaning
$ ilde{\sigma}_{{ m s},j}^2$	j th input variance corresponding to $ ilde{x}_{\mathrm{s},j}$
$\sigma_{ m s}^2$	average input variance corresponding to $ ilde{m{x}}_{ m s}$

F.2.4 Other Variables

Symbol	Meaning
const.	arbitrary constant
$oldsymbol{e}_j$	jth unit vector with zeros everywhere, except a 1 at j th position
$oldsymbol{I}_m$	identity matrix of size $m \times m$
0	null vector of respective dimension

F.2.5 Stochastics

Symbol	Meaning
$f_x(x)$	probability density function of random variable \times
$f_{x y}(x)$	conditional probability density function of random variable x given a realization of random variable y
$\mathcal{N}(m_{x},\sigma_{x}^2)$	Gaussian distribution with mean $m_{\rm x}$ and variance $\sigma_{\rm x}^2$
$\mathcal{N}(oldsymbol{m}_{x}, oldsymbol{\varPhi}_{x})$	multi-dimensional Gaussian distribution with mean $m_{ imes}$ and covariance matrix $oldsymbol{\Phi}_{ imes}$
$\mathcal{U}([a,b])$	uniform distribution over the interval $[a,b]$
$E\{\cdot\}$	expectation
$\mathrm{E}\{\cdot\mid \cdot\}$	conditional expectation
$\mathrm{E}_{x\sim f_x}\{\cdot\}$	expectation of random variable x w.r.t. pdf $f_x(x)$
$\Pr\{\cdot\}$	probability

F.2.6 Operators

Symbol	Meaning
$ \cdot $	absolute value of a scalar or cardinality of a set
$\left\ \cdot\right\ _{p}$	ℓ_p -norm
$\left\ \cdot\right\ _{\mathrm{F}}$	Frobenius norm
$[oldsymbol{M}]_{jj}$	j th diagonal entry of matrix $oldsymbol{M}$
$D_{KL}(\cdot \cdot)$	Kullback–Leibler divergence
$\mathbf{diag}(x_j)$	diagonal matrix with x_1, \ldots, x_N as entries
$\mathbf{diag}(\boldsymbol{x})$	diagonal matrix with entries of vector $oldsymbol{x}$ as entries
$\mathbf{diag}(\boldsymbol{M})$	vector with diagonal entries of matrix $oldsymbol{M}$
$\delta(\cdot)$	Dirac delta function
$\mathrm{h}(\cdot)$	differential entropy of a distribution
$\kappa(oldsymbol{A})$	condition number of matrix $oldsymbol{A}$
$\mathcal{L}(\cdot)$	Lagrangian function
$\mathcal{L}_{\mathrm{D}}(\cdot)$	Lagrangian dual
$\log(\cdot)$	logarithmus naturalis (natural logarithm)
$M_A(x_j)$	arithmetic mean $M_A(x_j) = \frac{1}{N} \sum_{j=1}^N x_j$
$M_{\rm H}(x_j)$	harmonic mean $\mathrm{M}_\mathrm{H}(x_j) = N / \sum_{j=1}^N 1 / x_j$
$ abla f(oldsymbol{x})$	gradient of $f({m x})$: $ abla f({m x}) = [rac{\partial}{\partial x_1} f({m x}), \dots, rac{\partial}{\partial x_N} f({m x})]^ op$
$\mathcal{O}(\cdot)$	Landau notation, bound on growth of a function
$\pi(\cdot)$	permutation
$\mathrm{sgn}(\cdot)$	signum function
$\operatorname{trace}(\boldsymbol{M})$	trace of $N \times N$ matrix \boldsymbol{M} : trace $(\boldsymbol{M}) = \sum_{j=1}^{N} [\boldsymbol{M}]_{jj}$

F.2.7 Exponential Families

Symbol	Meaning
$oldsymbol{g}(oldsymbol{x})$	sufficient statistic
heta	natural parameters
$Z(oldsymbol{ heta})$	partition function
$k(oldsymbol{x})$	carrier measure

F.2.8 Representatives

Symbol	Meaning
$q(oldsymbol{x})$	substitute distribution for posterior $f_{x y}(m{x})$
$q_{\mathrm{c}}(oldsymbol{x})$	representative for channel factor $f_{y x}(\boldsymbol{x})$
$q_{\mathrm{c},i}(oldsymbol{x})$	representative for $f_{y_i x}(x)$
${f q}_{ m s}({m x})$	representative for signal factor $f_x(x)$
$q_{\mathrm{s},j}(x_j)$	representative for $f_x(x_j)$, marginal of $q_s(\boldsymbol{x})$
$q_{\mathrm{v}}(oldsymbol{x})$	representative for variables $oldsymbol{x}$
$q_{\mathrm{v},j}(x_j)$	representative for variable x_j , marginal of $q_{\mathrm{v}}(oldsymbol{x})$

Bibliography

- [Ban13] Paolo Banelli. Non-Linear Transformations of Gaussians and Gaussian-Mixtures with Implications on Estimation and Information Theory. *arXiv preprint arXiv:1111.5950*, 2013.
- [BCDH10] Richard G. Baraniuk, Volkan Cevher, Marco F. Duarte, Chinmay Hegde. Model-Based Compressive Sensing. *IEEE Transactions on Information Theory*, vol. 56, no. 4, pp. 1982–2001, 2010.
- [BD09] Thomas Blumensath, Mike E. Davies. Iterative Hard Thresholding for Compressed Sensing. *Applied and Computational Harmonic Analysis*, vol. 27, no. 3, pp. 265–274, 2009.
- [Ber20] Ali Bereyhi. *Statistical Mechanics of Regularized Least Squares*. PhD thesis, Friedrich-Alexander-Universität Erlangen-Nürnberg, June 2020.
- [Bet35] Hans A. Bethe. Statistical Theory of Superlattices. *Proceedings of the Royal Society of London A*, vol. 150, no. 871, pp. 552–575, 1935.
- [BG96] Claude Berrou, Alain Glavieux. Near Optimum Error Correcting Coding and Decoding: Turbo-Codes. *IEEE Transactions on Communications*, vol. 44, no. 10, pp. 1261–1271, October 1996.
- [BG18] Stefan C. Birgmeier, Norbert Goertz. Optimizing Approximate Message Passing for Variable Measurement Noise. In European Signal Processing Conference (EU-SIPCO), pp. 484–488, 2018.
- [BG19] Stefan C. Birgmeier, Norbert Goertz. Exploiting General Multi-Dimensional Priors in Compressed-Sensing Reconstruction. In 12th International ITG Conference on Systems, Communications and Coding (SCC), pp. 1–6, 2019.
- [BGT93] Claude Berrou, Alain Glavieux, Punya Thitimajshima. Near Shannon Limit Error-Correcting Coding and Decoding: Turbo-Codes. In Proceedings IEEE International Conference on Communications (ICC), pp. 1064–1070, Geneva, Switzerland, May 1993.

[Bir19]	Stefan Birgmeier. Message Passing For Multidimensional Inverse Problems. PhD thesis, TU Wien, 2019.
[BM11]	Mohsen Bayati, Andrea Montanari. The Dynamics of Message Passing on Dense Graphs, with Applications to Compressed Sensing. <i>IEEE Transactions on Information Theory</i> , vol. 57, no. 2, pp. 764–785, 2011.
[BMPP20]	Alfredo Braunstein, Anna Paola Muntoni, Andrea Pagnani, Mirko Pieropan. Com- pressed Sensing Reconstruction using Expectation Propagation. <i>Journal of Physics</i> <i>A: Mathematical and Theoretical</i> , vol. 53, no. 18, 2020.
[Bro86]	Lawrence D. Brown. <i>Fundamentals of Statistical Exponential Families: with Applications in Statistical Decision Theory</i> . Institute of Mathematical Statistics, 1986.
[BV04]	Stephen P. Boyd, Lieven Vandenberghe. <i>Convex Optimization</i> . Cambridge University Press, 2004.
[CDS01]	Scott Shaobing Chen, David L. Donoho, Michael A. Saunders. Atomic Decompo- sition by Basis Pursuit. <i>SIAM review</i> , vol. 43, no. 1, pp. 129–159, 2001.
[ÇLO22]	Burak Çakmak, Yue M. Lu, Manfred Opper. Analysis of Random Sequen- tial Message Passing Algorithms for Approximate Inference. <i>arXiv preprint</i> <i>arXiv:2202.08198</i> , 2022.
[CREK05]	Shane F. Cotter, Bhaskar D. Rao, Kjersti Engan, Kenneth Kreutz-Delgado. Sparse Solutions to Linear Inverse Problems With Multiple Measurement Vectors. <i>IEEE</i> <i>Transactions on Signal Processing</i> , vol. 53, no. 7, pp. 2477–2488, 2005.
[CRT06]	Emmanuel J. Candès, Justin Romberg, Terence Tao. Robust Uncertainty Principles: Exact Signal Reconstruction from Highly Incomplete Frequency Information. <i>IEEE Transactions on Information Theory</i> , vol. 52, no. 2, pp. 489–509, 2006.
[CSD ⁺ 17]	Jun Won Choi, Byonghyo Shim, Yacong Ding, Bhaskar Rao, Dong In Kim. Com- pressed Sensing for Wireless Communications: Useful Tips and Tricks. <i>IEEE Com-</i> <i>munications Surveys & Tutorials</i> , vol. 19, no. 3, pp. 1527–1550, 2017.
[CZK14]	Francesco Caltagirone, Lenka Zdeborová, Florent Krzakala. On Convergence of Approximate Message Passing. In <i>Proceedings IEEE International Symposium on Information Theory (ISIT)</i> , pp. 1812–1816. IEEE, 2014.
[DB18]	Guillaume Dehaene, Simon Barthelmé. Expectation Propagation in the Large Data Limit. <i>Journal of the Royal Statistical Society: Series B (Statistical Methodology)</i> , vol. 80, no. 1, pp. 199–217, 2018.
[DDD04]	Ingrid Daubechies, Michel Defrise, Christine De Mol. An Iterative Thresholding Algorithm for Linear Inverse Problems with a Sparsity Constraint. <i>Communica-</i> <i>tions on Pure and Applied Mathematics: A Journal Issued by the Courant Institute</i> <i>of Mathematical Sciences</i> , vol. 57, no. 11, pp. 1413–1457, 2004.

- [DJB⁺95] Catherine Douillard, Michel Jézéquel, Claude Berrou, Annie Picart, Pierre Didier, Alain Glavieux. Iterative Correction of Intersymbol Interference: Turbo-Equalization. *European Transactions on Telecommunications*, vol. 6, no. 5, pp. 507– 511, 1995.
- [DMM09] David L. Donoho, Arian Maleki, Andrea Montanari. Message-Passing Algorithms for Compressed Sensing. *Proceedings of the National Academy of Sciences*, vol. 106, no. 45, pp. 18914–18919, 2009.
- [DMM10] David L. Donoho, Arian Maleki, Andrea Montanari. Message Passing Algorithms for Compressed Sensing: I. Motivation and Construction. In Proceedings IEEE Information Theory Workshop (ITW), pp. 1–5, 2010.
- [Don06] David L. Donoho. Compressed Sensing. *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [DT97] G. B. Dantzig, M. N. Thapa. *Linear Programming 1: Introduction*. Springer Series in Operations Research. Springer, 1997.
- [EK12] Yonina C. Eldar, Gitta Kutyniok. Compressed Sensing: Theory and Applications. Cambridge University Press, 2012.
- [Fis22] Ronald A. Fisher. On the Mathematical Foundations of Theoretical Statistics. Philosophical Transactions of the Royal Society of London. Series A, containing papers of a mathematical or physical character, vol. 222, no. 594-604, pp. 309–368, 1922.
- [Fis02] Robert F. H. Fischer. *Precoding and Signal Shaping for Digital Transmission*. John Wiley & Sons, 2002.
- [Fis16] Robert F. H. Fischer. Notes on MMSE Estimation and Bias Compensation. Unpublished Manuscript, 2016.
- [Fis18] Robert F. H. Fischer. Signal Theory. Lecture Notes, Ulm University, 2018.
- [For01] G. David Forney. Codes on Graphs: Normal Realizations. *IEEE Transactions on Information Theory*, vol. 47, no. 2, pp. 520–548, 2001.
- [For03] G. David Forney. On the Role of MMSE Estimation in Approaching the Information-Theoretic Limits of Linear Gaussian Channels: Shannon meets Wiener. In Allerton Conference on Communication, Control, and Computing, pp. 430-439, 2003.
- [For18] G. David Forney. Codes on graphs: Models for elementary algebraic topology and statistical physics. *IEEE Transactions on Information Theory*, vol. 64, no. 12, pp. 7465–7487, 2018.

[FS07]	Robert F.H. Fischer, Cristian Siegl. Inflated Lattice Precoding, Bias Compensation,
	and the Uplink/Downlink Duality: The Connection. IEEE Communications Letters,
	vol. 11, no. 2, pp. 185–187, 2007.

- [FS22] Robert F. H. Fischer, Carmen Sippel. Unbiasing in Iterative Reconstruction Algorithms for Discrete Compressed Sensing. In Gitta Kutyniok, Holger Rauhut, Robert J. Kunsch, editors, *Compressed Sensing in Information Processing*, Applied and Numerical Harmonic Analysis, chapter 6. Birkhäuser, 2022.
- [FSG20] Robert F. H. Fischer, Carmen Sippel, Norbert Goertz. VAMP with Vector-Valued Diagonalization. In Proceedings IEEE International Conference on Acoustics, Speech and Signal Processing, 2020.
- [FSRS16] Alyson K. Fletcher, Mojtaba Sahraee-Ardakan, Sundeep Rangan, Philip Schniter. Expectation Consistent Approximate Inference: Generalizations and Convergence. In Proceedings IEEE International Symposium on Information Theory (ISIT), pp. 190–194, 2016.
- [Gal62] Robert Gallager. Low-Density Parity-Check Codes. *IRE Transactions on Information Theory*, vol. 8, no. 1, pp. 21–28, 1962.
- [GFS63] Izrail Moiseevitch Gelfand, Sergej Vasilevitch Fomin, Richard A. Silverman. *Calculus of Variations*. Prentice-Hall, Inc., 1963.
- [GH11] Qinghua Guo, Defeng David Huang. A Concise Representation for the Soft-In Soft-Out LMMSE Detector. *IEEE Communications Letters*, vol. 15, no. 5, pp. 566– 568, 2011.
- [Gra14] Knut Graichen. Methoden der Optimierung und optimalen Steuerung. Vorlesungsskript (in German), Universität Ulm, 2014.
- [GSV05] Dongning Guo, Shlomo Shamai, Sergio Verdú. Mutual Information and Minimum Mean-Square Error in Gaussian Channels. *IEEE Transactions on Information The*ory, vol. 51, no. 4, pp. 1261–1282, 2005.
- [GV96] Gene H. Golub, Charles F. Van Loan. *Matrix Computations*. John Hopkins University Press, 3 edition, 1996.
- [GW06] Dongning Guo, Chih-Chun Wang. Asymptotic Mean-Square Optimality of Belief Propagation for Sparse Linear Systems. In Proceedings IEEE Information Theory Workshop (ITW), pp. 194–198. IEEE, 2006.
- [GWSV11] Dongning Guo, Yihong Wu, Shlomo Shamai Shitz, Sergio Verdú. Estimation in Gaussian Noise: Properties of the Minimum Mean-Square Error. IEEE Transactions on Information Theory, vol. 57, no. 4, pp. 2371–2385, 2011.

- [HF14] Johannes B. Huber, Robert F.H. Fischer. Informationstheorie und deren Anwendungen zur Nachrichtenübertragung. Vorlesungsskript (in German), Friedrich-Alexander-Universität Erlangen-Nürnberg, 2014.
- [HJ09] Roger A. Horn, Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 2009.
- [HOP96] Joachim Hagenauer, Elke Offer, Lutz Papke. Iterative Decoding of Binary Block and Convolutional Codes. *IEEE Transactions on Information Theory*, vol. 42, no. 2, pp. 429–445, 1996.
- [HOW⁺05] Tom Heskes, Manfred Opper, Wim Wiegerinck, Ole Winther, Onno Zoeter. Approximate Inference Techniques with Expectation Constraints. *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2005, no. 11, 2005.
- [HTP12] Jinping Hao, Filippo Tosato, Robert J. Piechocki. Sequential Compressive Sensing in Wireless Sensor Networks. In IEEE 75th Vehicular Technology Conference (VTC Spring), pp. 1–5, 2012.
- [Jay57] Edwin T. Jaynes. Information Theory and Statistical Mechanics. *Physical Review*, vol. 106, no. 4, pp. 620–630, 1957.
- [Jay82] Edwin T. Jaynes. On the Rationale of Maximum-Entropy Methods. *Proceedings* of the IEEE, vol. 70, no. 9, pp. 939–952, 1982.
- [Kay93] Steven M. Kay. Fundamentals of Statistical Signal Processing: Estimation Theory. Prentice-Hall, Inc., 1993.
- [KF98] Frank R. Kschischang, Brendan J. Frey. Iterative Decoding of Compound Codes by Probability Propagation in Graphical Models. *IEEE Journal on Selected Areas in Communications*, vol. 16, no. 2, pp. 219–230, 1998.
- [KFL01] Frank R. Kschischang, Brendan J. Frey, H.-A. Loeliger. Factor Graphs and the Sum-Product Algorithm. *IEEE Transactions on Information Theory*, vol. 47, no. 2, pp. 498–519, 2001.
- [KL51] Solomon Kullback, Richard A. Leibler. On Information and Sufficiency. The Annals of Mathematical Statistics, vol. 22, no. 1, pp. 79–86, March 1951.
- [KMS⁺12] Florent Krzakala, Marc Mézard, Francois Sausset, Yifan Sun, Lenka Zdeborová. Probabilistic Reconstruction in Compressed Sensing: Algorithms, Phase Diagrams, and Threshold Achieving Matrices. *Journal of Statistical Mechanics: Theory* and Experiment, vol. 2012, no. 08, 2012.
- [KPN⁺14] Kee-Hoon Kim, Hosung Park, Jong-Seon No, Habong Chung, Dong-Joon Shin. Clipping Noise Cancelation for OFDM Systems Using Reliable Observations

Based on Compressed Sensing. *IEEE Transactions on Broadcasting*, vol. 61, no. 1, pp. 111–118, 2014.

- [Lan09] Edmund Landau. Handbuch der Lehre von der Verteilung der Primzahlen, volume 1.B. G. Teubner, 1909.
- [LCJ15] Markus Leinonen, Marian Codreanu, Markku Juntti. Sequential Compressed Sensing With Progressive Signal Reconstruction in Wireless Sensor Networks. IEEE Transactions on Wireless Communications, vol. 14, no. 3, pp. 1622–1635, 2015.
- [LDH⁺07] Hans-Andrea Loeliger, Justin Dauwels, Junli Hu, Sascha Korl, Li Ping, Frank R. Kschischang. The Factor Graph Approach to Model-Based Signal Processing. *Proceedings of the IEEE*, vol. 95, no. 6, pp. 1295–1322, 2007.
- [LHD20] Ning Li, Mengjia Huang, Zhongliang Deng. Elimination of Noise Distortion for OFDM Systems by Compressed Sensing Based on Distance Metric. *IEEE Access*, vol. 8, pp. 223700–223707, 2020.
- [LHL⁺19] Xinghui Liu, Lingna Huy, Yue Liz, Lianghui Ding, Feng Yang, Cheng Zhi. Modified Peak Clipping and Compressed Sensing Recovery Scheme in OFDM System. In 2019 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB), pp. 1–5. IEEE, 2019.
- [LHLT15] Yingzhen Li, José Miguel Hernández-Lobato, Richard E. Turner. Stochastic Expectation Propagation. Advances in Neural Information Processing Systems, vol. 28, 2015.
- [Loe04] H.-A. Loeliger. An Introduction to Factor Graphs. *IEEE Signal Processing Magazine*, vol. 21, no. 1, pp. 28–41, 2004.
- [Mal11] Arian Maleki. Approximate Message Passing Algorithms for Compressed Sensing. PhD thesis, Stanford University, September 2011.
- [MG98] Michael L. Moher, T. Aaron Gulliver. Cross-Entropy and Iterative Decoding. *IEEE Transactions on Information Theory*, vol. 44, no. 7, pp. 3097–3104, 1998.
- [Min01] Thomas Peter Minka. Expectation Propagation for Approximate Bayesian Inference. In *Proceedings UAI-2001*, pp. 362–369, 2001.
- [Min05] Thomas Peter Minka. Divergence Measures and Message Passing. Technical report, Microsoft Research, 2005.
- [MKTZ15] Andre Manoel, Florent Krzakala, Eric Tramel, Lenka Zdeborovà. Swept Approximate Message Passing for Sparse Estimation. In International Conference on Machine Learning, pp. 1123–1132, 2015.

- [MMC98] Robert J. McEliece, David J. C. MacKay, Jung-Fu Cheng. Turbo Decoding as an Instance of Pearl's "Belief Propagation" Algorithm. *IEEE Journal on Selected Areas* in Communications, vol. 16, no. 2, pp. 140–152, 1998.
- [Moh97a] Michael L. Moher. *Cross-Entropy and Iterative Detection*. PhD thesis, Carleton University, 1997.
- [Moh97b] Michael L. Moher. Turbo-based Multiuser Detection. In *Proceedings IEEE International Symposium on Information Theory (ISIT)*, pp. 195–195, 1997.
- [MP17] Junjie Ma, Li Ping. Orthogonal AMP. IEEE Access, vol. 5, pp. 2020–2033, 2017.
- [MSW10] Dmitry M. Malioutov, Sujay R. Sanghavi, Alan S. Willsky. Sequential Compressed Sensing. IEEE Journal of Selected Topics in Signal Processing, vol. 4, no. 2, pp. 435– 444, 2010.
- [NG09] Frank Nielsen, Vincent Garcia. Statistical Exponential Families: A Digest with Flash Cards. *arXiv preprint arXiv:0911.4863*, 2009.
- [NT09] Deanna Needell, Joel A. Tropp. CoSaMP: Iterative Signal Recovery from Incomplete and Inaccurate Samples. Applied and Computational Harmonic Analysis, vol. 26, no. 3, pp. 301–321, 2009.
- [OW05] Manfred Opper, Ole Winther. Expectation Consistent Approximate Inference. *Journal of Machine Learning Research*, vol. 6, pp. 2177–2204, December 2005.
- [Pea88] Judea Pearl. Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. Morgan Kaufmann Publishers, INC., 1988.
- [PP02] Athanasios Papoulis, S. Unnikrishna Pillai. Probability, Random Variables, and Stochastic Processes. McGraw-Hill, 2002.
- [PRK93] Yagyensh Chandra Pati, Ramin Rezaiifar, Perinkulam Sambamurthy Krishnaprasad. Orthogonal Matching Pursuit: Recursive Function Approximation with Applications to Wavelet Decomposition. In Proceedings of 27th Asilomar Conference on Signals, Systems and Computers, pp. 40–44, 1993.
- [Pro00] John G. Proakis. *Digital Communications*. McGraw-Hill, New York, NY, USA, 4th edition, 2000.
- [PSC⁺15] Marcelo Pereyra, Philip Schniter, Emilie Chouzenoux, Jean-Christophe Pesquet, Jean-Yves Tourneret, Alfred O. Hero, Steve McLaughlin. A Survey of Stochastic Simulation and Optimization Methods in Signal Processing. *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 2, pp. 224–241, 2015.
- [Ran11] Sundeep Rangan. Generalized Approximate Message Passing for Estimation With Random Linear Mixing. In Proceedings IEEE International Symposium on Information Theory (ISIT), pp. 2168–2172, 2011.

[RFSK16]	Sundeep Rangan, Alyson K. Fletcher, Philip Schniter, Ulugbek S. Kamilov. In- ference for Generalized Linear Models via Alternating Directions and Bethe Free Energy Minimization. <i>IEEE Transactions on Information Theory</i> , vol. 63, no. 1, pp. 676–697, 2016.
[RSF17]	Sundeep Rangan, Philip Schniter, Alyson K. Fletcher. Vector Approximate Mes- sage Passing. In <i>Proceedings IEEE International Symposium on Information Theory</i> <i>(ISIT)</i> , pp. 1588–1592, 2017.
[RSF19]	Sundeep Rangan, Philip Schniter, Alyson K. Fletcher. Vector Approximate Mes- sage Passing. <i>IEEE Transactions on Information Theory</i> , vol. 65, no. 10, pp. 6664– 6684, 2019.
[SBW ⁺ 05]	Shriram Sarvotham, Dror Baron, Michael Wakin, Marco F. Duarte, Richard G. Baraniuk. Distributed Compressed Sensing of Jointly Sparse Signals. In <i>Asilomar Conference on Signals, Systems, and Computers</i> , pp. 1537–1541, 2005.
[SD11]	Henning F. Schepker, Armin Dekorsy. Sparse Multi-User Detection for CDMA Transmission Using Greedy Algorithms. In <i>8th International Symposium on Wire-</i> <i>less Communication Systems</i> , pp. 291–295, 2011.
[SF14]	Susanne Sparrer, Robert F. H. Fischer. Adapting Compressed Sensing Algorithms to Discrete Sparse Signals. In <i>18th International ITG Workshop on Smart Antennas (WSA)</i> , 2014.
[SF16]	Susanne Sparrer, Robert F. H. Fischer. Enhanced Iterative Hard Thresholding for the Estimation of Discrete-Valued Sparse Signals. In <i>European Signal Processing</i> <i>Conference (EUSIPCO)</i> , pp. 71–75, 2016.
[SF17]	Susanne Sparrer, Robert F. H. Fischer. Unveiling Bias Compensation in Turbo- Based Algorithms for (Discrete) Compressed Sensing. In <i>European Signal Process-</i> <i>ing Conference (EUSIPCO)</i> , pp. 2091–2095, 2017.
[SF18]	Susanne Sparrer, Robert F. H. Fischer. Bias Compensation in Iterative Soft-Feedback Algorithms With Application to (Discrete) Compressed Sensing. In <i>IEEE Statistical Signal Processing Workshop (SSP)</i> , pp. 35–39, 2018.
[SF21]	Carmen Sippel, Robert F. H. Fischer. VAMP with Individual Variances and Sequen- tial Processing for Compressed Sensing. In <i>European Signal Processing Conference</i> <i>(EUSIPCO)</i> , Dublin, Ireland, August 2021.
[SF22]	Carmen Sippel, Robert F. H. Fischer. Variants of VAMP for Signal Recovery in Wireless Sensor Networks. In <i>IEEE International Mediterranean Conference on Communications and Networking (MeditCom)</i> , Athens, Greece, September 2022.
[Sha48]	Claude E. Shannon. A Mathematical Theory of Communication. <i>Bell System Technical Journal</i> , vol. 27, pp. 379–423, July 1948.

[SJ81]

- [SJ80] John E. Shore, Rodney W. Johnson. Axiomatic Derivation of the Principle of Maximum Entropy and the principle of Minimum Cross-Entropy. IEEE Transactions on Information Theory, vol. 26, no. 1, pp. 26-37, 1980. John Shore, Rodney Johnson. Properties of Cross-Entropy Minimization. IEEE Transactions on Information Theory, vol. 27, no. 4, pp. 472-482, 1981. [SN08] Matthias W. Seeger, Hannes Nickisch. Compressed Sensing and Bayesian Experimental Design. In Proceedings of the 25th International Conference on Machine Learning, pp. 912-919, 2008. [SN11] Matthias Seeger, Hannes Nickisch. Fast Convergent Algorithms for Expectation Propagation Approximate Bayesian Inference. In Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, pp. 652–660, 2011. [Spa19] Susanne Sparrer. Algorithms for (Discrete) Compressed Sensing-A Communications Engineering Perspective. PhD thesis, Ulm University, September 2019. [SSF13] Susanne Sparrer, Andreas Schenk, Robert F.H. Fischer. Communication over Impulsive Noise Channels: Channel Coding vs. Compressed Sensing. In 9th International ITG Conference on Systems, Communication and Coding (SCC), 2013. [Tan81] R. Tanner. A Recursive Approach to Low Complexity Codes. IEEE Transactions on Information Theory, vol. 27, no. 5, pp. 533-547, 1981. [Tib96] Robert Tibshirani. Regression Shrinkage and Selection via the LASSO. Journal of the Royal Statistical Society: Series B (Methodological), vol. 58, no. 1, pp. 267–288, 1996.
- [TS11] Michael Tüchler, Andrew C. Singer. Turbo Equalization: An Overview. IEEE Transactions on Information Theory, vol. 57, no. 2, pp. 920-952, 2011.
- [Ver10] Sergio Verdú. Mismatched estimation and relative entropy. IEEE Transactions on Information Theory, vol. 56, no. 8, pp. 3712–3720, 2010.
- [VS13] Jeremy P. Vila, Philip Schniter. Expectation-Maximization Gaussian-Mixture Approximate Message Passing. IEEE Transactions on Signal Processing, vol. 61, no. 19, pp. 4658-4672, 2013.
- [Wib96] Niclas Wiberg. Codes and Decoding on General Graphs. PhD thesis, Department of Electrical Engineering, Linköping University Sweden, 1996.
- [WJ08] Martin J. Wainwright, Michael I. Jordan. Graphical Models, Exponential Families, and Variational Inference. Foundations and Trends(R) in Machine Learning, vol. 1, no. 1-2, pp. 1-305, 2008.
- [Wu12] Nailong Wu. The Maximum Entropy Method, volume 32. Springer Verlag, 2012.

[WV12]	Yihong Wu, Sergio Verdú. Optimal Phase Transitions in Compressed Sensing. <i>IEEE Transactions on Information Theory</i> , vol. 58, no. 10, pp. 6241–6263, 2012.
[XDS13]	Tong Xue, Xiaodai Dong, Yi Shi. Multiple Access and Data Reconstruction in Wireless Sensor Networks Based on Compressed Sensing. <i>IEEE Wireless Communications (Magazine)</i> , vol. 12, no. 7, pp. 3399–3411, 2013.
[YFW05]	Jonathan S. Yedidia, William T. Freeman, Yair Weiss. Constructing Free-Energy Approximations and Generalized Belief Propagation Algorithms. <i>IEEE Transac-</i> <i>tions on Information Theory</i> , vol. 51, no. 7, pp. 2282–2312, 2005.

[YMG08] Jennifer Yick, Biswanath Mukherjee, Dipak Ghosal. Wireless Sensor Network Survey. *Computer Networks*, vol. 52, no. 12, pp. 2292–2330, 2008.

List of Scientific Publications and Presentations

- Jonathan Bechter, Carmen Sippel, Christian Waldschmidt. Bats-inspired Frequency Hopping for Mitigation of Interference between Automotive Radars. In *IEEE MTT-S International Conference on Microwaves for Intelligent Mobility (ICMIM)*, San Diego, CA, USA, May 2016.
- [2] Carmen Sippel, Sven Puchinger, Robert F.H. Fischer. Differential Linear Network Coding in Slowly-Varying Networks. Presentation at 33. Sitzung der ITG-Fachgruppe "Angewandte Informationstheorie", Codes, Lattices, and Decoding, Ulm, Germany, March 2019.
- [3] Carmen Sippel, Sven Puchinger, Robert F.H. Fischer. Differential Linear Network Coding in Slowly-Varying Networks. Poster at *the 19th Joint Workshop on Communications and Coding (JWCC)*, Kühtai, Austria, March 2019.
- [4] Carmen Sippel, Sven Puchinger, Robert F.H. Fischer. Differential Linear Network Coding in Slowly-Varying Networks. Presentation at *Stuttgart Workshop on Coding and Communications (SWCC)*, Stuttgart, Germany, July 2019.
- [5] Carmen Sippel, Cornelia Ott, Sven Puchinger, Martin Bossert. Reed–Solomon Codes over Fields of Characteristic Zero. In *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, pp. 1537–1541, July 2019.
- [6] Carmen Sippel, Robert F.H. Fischer. On the Invariance of Recovery Algorithms for Compressed Sensing based on Expectation-Consistent Approximate Inference. In Proceedings of the 24th International ITG Workshop on Smart Antennas (WSA), Hamburg, Germany, February 2020.
- [7] Robert F.H. Fischer, Carmen Sippel, Norbert Goertz. VAMP with Vector-Valued Diagonalization. In Proceedings of the IEEE International Conference on Acoustic, Speech and Signal Processing (ICASSP), Barcelona, Spain, May 2020.
- [8] Carmen Sippel, Robert F.H. Fischer. VAMP with Individual Variances and Sequential Processing for Compressed Sensing. In *Proceedings of the 29th European Signal Processing Conference (EUSIPCO)*, pp. 1935–1939, Dublin, Ireland, August 2021.

- [9] Carmen Sippel, Robert F.H. Fischer. Stabilization Techniques for Iterative Algorithms in Compressed Sensing. Online available at arxiv, https://arxiv.org/abs/2109.14917, 2021.
- [10] Carmen Sippel, Robert F.H. Fischer. Variants of VAMP for Signal Recovery in Wireless Sensor Networks. In Proceedings of the IEEE International Mediterranean Conference on Communications and Networking (MeditCom), Athens, Greece, September 2022.
- [11] Carmen Sippel, Robert F.H. Fischer. Comparison of Parallel and Sequential Algorithms for Compressed Sensing. In *Proceedings of the 7th International Conference on Frontiers of Signal Processing (ICFSP)*, Paris, France, September 2022.
- [12] Carmen Sippel, Robert F.H. Fischer. Evaluation of the Fractional Approach for Iterative Algorithms in Compressed Sensing. In Proceedings of the International ITG 26th Workshop on Smart Antennas (WSA) and 13th Conference on Systems, Communications, and Coding (SCC), Braunschweig, Germany, February 2023.
- [13] Elena Sterk, Carmen Sippel, Robert F.H. Fischer. Comparison of Damping Approaches for AMP. In Proceedings of the International ITG 26th Workshop on Smart Antennas (WSA) and 13th Conference on Systems, Communications, and Coding (SCC), Braunschweig, Germany, February 2023.

Curriculum Vitæ

For data protection reasons, the curriculum vitæ has been removed from the online version.

ISBN: 978-3-948303-38-9