

Solving Linear and Bilinear Inverse Problems using
Approximate Message Passing Methods

Dissertation

Presented in Partial Fulfillment of the Requirements for the Degree
Doctor of Philosophy in the Graduate School of The Ohio State
University

By

Subrata Sarkar, B.Tech., M.S.

Graduate Program in Electrical and Computer Engineering

The Ohio State University

2020

Dissertation Committee:

Dr. Philip Schniter, Advisor

Dr. Lee C. Potter

Dr. Rizwan Ahmad

© Copyright by

Subrata Sarkar

2020

Abstract

Recovering an unknown vector from its noisy linear measurements is an important problem that arises in many fields. In linear inverse problems, the forward operator is perfectly known, whereas in bilinear inverse problems, the forward operator has some unknown parameters. Most existing recovery algorithms are either slow to converge or give a final solution that is not accurate. In this dissertation, we develop algorithms using Approximate Message Passing (AMP) methods to solve linear and bilinear inverse problems.

First, we consider the computationally efficient Vector Approximate Message Passing (VAMP) algorithm, which is based on the Expectation Propagation framework. VAMP minimizes a cost function, also known as the Gibbs free energy in statistical physics, under moment-matching constraints. It iteratively calls a denoiser that is chosen based on prior knowledge about the unknown vector that we want to recover. VAMP has a remarkable property, that when the sensing matrix is a typical instance of a large right-rotationally invariant random matrix, the per-iteration macroscopic behaviour of VAMP can be exactly predicted by a set of scalar equations that yield the so called state-evolution (SE). The state-evolution can be used to predict the MSE performance of VAMP at each iteration. The SE was first proven for separable Lipschitz denoisers. In this work, we extend the state-evolution to a larger class of non-separable Lipschitz denoisers. Empirical results show that the SE also accurately

predicts VAMP’s performance in bilinear problems solved that have been converted to linear problems using the “lifting” technique.

In bilinear inverse problems, the forward operator is written as a linear combination of known matrices with unknown weights. Problems such as dictionary learning, CS with matrix uncertainty, self-calibration, are all instances of bilinear inverse problems. We propose a new algorithm called Bilinear Adaptive Vector Approximate Message Passing (BAd-VAMP). Our method is based on the EM-VAMP algorithm, where Expectation Maximization (EM) is combined with VAMP to estimate the unknown parameters in the likelihood and prior. We show that our BAd-VAMP method is efficient and robust to the conditioning of the forward operator, which is a common problem with algorithms based on the earlier AMP algorithm.

Finally, we consider the problem of magnetic resonance imaging (MRI). MRI is a safe, non-invasive method to capture internal images of our body. MRI’s long data-acquisition time is one of its biggest drawbacks. In accelerated MRI, a small number of samples in k-space are acquired and the goal is to accurately reconstruct the image from these few measurements. We develop two algorithms for MRI reconstruction using Approximate Message Passing methods with non-separable denoisers. Our first algorithm is a modification of the existing AMP algorithm, where we propose a scaling in k-space to debias the denoiser input error. Our second algorithm is a modification of VAMP where we propose a robust damping scheme to stabilize VAMP in MRI. We compare our algorithms against conventional plug-and-play algorithms on the fastMRI knee dataset.

This is dedicated to Baba and Ma

Acknowledgments

I am grateful to several people who helped me in completing this dissertation. I would like to first thank my advisor Professor Phil Schniter for his guidance and support throughout my PhD. His creative solutions to difficult problems, work ethic and deep intuition helped me overcome several obstacles in my research. I have learnt to be a better thinker, writer and programmer from him. I am thankful to him for involving me in his research collaborations. This dissertation would not be possible without him.

I would like to thank my committee members, Professor Rizwan Ahmad and Professor Lee Potter, for their helpful feedbacks in my candidacy exam. I thank Tricia Toothman, for organizing social events like the weekly coffee social and monthly meetings with the department chair. My IPS labmates, Nithin Sugavanam, Evan Byrne, Ted Reehorst, Michael Wharton and Saurav Shastri for providing a friendly and supportive community.

I am fortunate to have Bortik Bandyopadhyay and Soumik Mandal as my close friends. I will miss our evening discussions on coding, machine learning and careers. They helped me stay motivated in this unprecedented time, when the entire world is struggling with the pandemic.

Finally, I am indebted to my family for their unconditional love and support. I thank Baba who taught me valuable life lessons and encouraged me to be curious and

think big; Maa, for her love and care; Didi and Kishor da, for always looking after me. I thank my family for trusting my decisions and believing in me. I couldn't have finished my PhD without your love.

Vita

May 2014	B.Tech. Electronics and Communication Engineering, Indian Institute of Technology, Guwahati
May 2016	M.S. Electrical and Computer Engineering, The Ohio State University
2016 - present	Graduate Research Assistant, The Ohio State University

Publications

Research Publications

A. K. Fletcher, P. Pandit, S. Rangan, S. Sarkar, P. Schniter “Plug-in Estimation in High-Dimensional Linear Inverse Problems: A Rigorous Analysis,” *Advances in Neural Information Processing Systems*, pp. 7440-7449, 2018.

S. Sarkar, A. K. Fletcher, S. Rangan, P. Schniter “Bilinear Recovery Using Adaptive Vector-AMP,” *IEEE Transactions on Signal Processing*, vol. 67, no. 13, pp. 3383-3396, 2019.

S. Rangan, P. Schniter, A. K. Fletcher, S. Sarkar “On the Convergence of Approximate Message Passing with Arbitrary Matrices,” *IEEE Transactions on Information Theory*, vol. 65, no. 9, pp. 5339-5351, 2019.

S. Sarkar, A. K. Fletcher, S. Rangan, P. Schniter “Bilinear Recovery Using Adaptive Vector-AMP,” *Proc. Asilomar Conf. on Signals, Systems, and Computers*, (Pacific Grove, CA), Nov. 2019.

K.H. Ngo, M. Guillaud, A. Decurninge, S. Yang, S. Sarkar, P. Schniter “Non-Coherent Multi-User Detection Based on Expectation Propagation,” *Proc. Asilomar Conf. on Signals, Systems, and Computers*, (Pacific Grove, CA), Nov. 2019.

Fields of Study

Major Field: Electrical and Computer Engineering

Studies in:

Machine Learning
Signal Processing

Table of Contents

	Page
Abstract	ii
Dedication	iv
Acknowledgments	v
Vita	vii
List of Tables	xii
List of Figures	xiii
1. Introduction	1
1.1 Standard Linear Model	3
1.2 Standard Bilinear Model	5
1.3 Outline and Contribution	5
1.3.1 State Evolution of VAMP	6
1.3.2 The BAdVAMP Algorithm	6
1.3.3 VAMP for MRI	7
1.4 Notation	8
2. Vector Approximate Message Passing	10
2.1 Introduction	10
2.2 Variational Inference	11
2.3 Expectation Consistent Approximate Inference	12
2.4 State Evolution	17
2.4.1 Large System Limit	18
2.4.2 State Evolution of Non-Separable VAMP	20
2.5 Numerical Experiments	21

2.5.1	Image Recovery via Non-Separable VAMP	21
2.5.2	Bilinear Estimation via Lifting	23
3.	Bilinear Recovery using Adaptive Vector Approximate Message Passing	28
3.1	Introduction	28
3.2	Prior Work	29
3.3	Contributions	33
3.4	Problem Formulation	34
3.5	Expectation Maximization	36
3.6	Bilinear EM-VAMP	38
3.7	Bilinear Adaptive VAMP	40
3.7.1	Estimating the Matrix \mathbf{X}	45
3.7.2	Learning the Parameters $\boldsymbol{\theta}_A$	46
3.7.3	Learning the Noise Precision γ_w	48
3.7.4	Algorithm Enhancements	50
3.8	Relation to Previous Work	51
3.9	Numerical Experiments	53
3.9.1	CS with Matrix Uncertainty	53
3.9.2	Self-Calibration	56
3.9.3	Calibration in Tomography	59
3.9.4	Noiseless Dictionary Learning	64
3.9.5	Noisy, Ill-Conditioned Dictionary Learning	65
4.	Image Reconstruction in Magnetic Resonance Imaging	70
4.1	Introduction	70
4.2	Measurement Model	71
4.3	Signal Denoiser	73
4.4	Plug-and-Play Algorithms	75
4.5	Our Approach	77
4.5.1	AMP for MRI	78
4.5.2	VAMP for MRI	85
4.6	Numerical Experiments	89
4.6.1	Single-coil MRI	90
4.6.2	Multicoil MRI	94
4.6.3	Randomized Single-Coil MRI	98
5.	Conclusions	100

Appendices	102
A. Signal Denoisers	102
A.1 Singular Value Thresholding	102
Bibliography	105

List of Tables

Table	Page
3.1 PSNR (dB) in the tomography experiment	61

List of Figures

Figure	Page
2.1 Image recovery using BM3D-VAMP	22
2.2 SE prediction in VAMP for image recovery	23
2.3 CS with matrix uncertainty: SE prediction in VAMP	24
2.4 CS with matrix uncertainty: NMSE (dB) vs M/P	26
2.5 CS with matrix uncertainty: NMSE (dB) vs $\text{cond}(\mathbf{A})$	26
2.6 Self-Calibration: success rate vs sparsity of lifted VAMP	27
3.1 CS with matrix uncertainty: NMSE (dB) vs sampling ratio	54
3.2 CS with matrix uncertainty: NMSE (dB) vs mean of \mathbf{A}_i	56
3.3 CS with matrix uncertainty: run-time vs sampling ratio	57
3.4 Self-Calibration: comparison of success rates of algorithms	58
3.5 Calibration in tomography: image reconstruction	62
3.6 Calibration in tomography: error images	63
3.7 Dictionary Learning: BAd-VAMP success rate	67
3.8 Noiseless dictionary learning: NMSE (dB) vs dimension and sparsity	68
3.9 Noisy dictionary learning: NMSE (dB) vs condition number	69

4.1	Example sampling masks for acceleration rate $R = 4$	73
4.2	Condition number vs acceleration rate in MRI	74
4.3	NMSE (dB) of AMP in MRI for different type of scaling in k-space	80
4.4	MRI ground truth and sampling masks	89
4.5	Single-coil MRI using Cartesian mask and BM3D denoiser	91
4.6	Single-coil MRI using Cartesian sampling and DnCNN denoiser	92
4.7	Single-coil MRI using variable density sampling and DnCNN denoiser	93
4.8	Coil sensitivities in multicoil MRI	94
4.9	Multicoil MRI using Cartesian sampling and BM3D denoiser	95
4.10	Multicoil MRI using Cartesian sampling and DnCNN denoiser	96
4.11	Multicoil MRI using variable density sampling and DnCNN denoiser	97
4.12	Coded single-coil MRI using Cartesian sampling and DnCNN	99

Chapter 1: Introduction

In many problems of interest in science and engineering we need to reliably recover an unknown high-dimensional vector from noisy measurements. The unknown vector could denote model parameters, sensor data, etc. This type of problem commonly arises in several areas of computational imaging such as tomographic imaging, radar, magnetic resonance imaging (MRI), etc. Several technological advancements in the last two decades has made it possible to solve these complex problems by designing efficient computational algorithms.

In these problems, we have some knowledge of the underlying process that generated the observed measurements from the signal, i.e., we have some knowledge of the likelihood function. Assuming the true signal that we are interested to recover is \mathbf{x} , then the observations $\mathbf{y} \in \mathbb{R}^M$ are generated according to the model

$$\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{w}, \tag{1.1}$$

where \mathcal{A} is the forward operator and \mathbf{w} is measurement noise. The goal in (1.1) is to recover the signal \mathbf{x} from measurements \mathbf{y} . These type of problems are known as “inverse problems” as we aim to invert the measurement process to recover the signal that generated it.

The physics of the measurement device decides \mathcal{A} . In this dissertation, we primarily focus on linear \mathcal{A} . In some applications, we have perfect knowledge of \mathcal{A} and in some cases it could be an approximation of reality. When it is perfectly known, the measurements \mathbf{y} are a linear combination of \mathbf{x} and such problems are also called “linear inverse” problems. We also consider the case when \mathcal{A} is known up to some unknown parameters \mathbf{b} . We represent the parametric forward operator as $\mathcal{A}_{\mathbf{b}}$. The unknown parameters capture the uncertainty in \mathcal{A} , and we aim to jointly recover the signal \mathbf{x} and learn the parameters \mathbf{b} .

Inverse problems are often ill-posed [91] because the forward operator is non-invertible which makes it hard to solve. We need to design recovery algorithms that are robust to noise, model mismatches and ill-conditioning of the problem because in applications such as medical imaging a false diagnosis based on the recovered image can be fatal.

In this dissertation, we develop algorithms for solving linear and bilinear inverse problems using the Approximate Message Passing framework. We develop some theoretical results on the VAMP [73] algorithm. We propose a computationally efficient algorithm called BAdVAMP for solving bilinear inverse problems. In MRI, we modify the AMP [32] and VAMP [73] algorithms by using a robust damping scheme and scaling the k-space measurements and we empirically compare the performance against well-known plug-and-play algorithms such as PnP-ADMM [95] and PnP-FISTA [44].

1.1 Standard Linear Model

In a Standard Linear Model (SLM), \mathcal{A} is a linear operator that can be represented as a matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$. We can write (1.1) as,

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}). \quad (1.2)$$

Compressive Sensing [19, 31] is a well known application where the number of measurements are less than the unknowns, i.e., $M \ll N$. In this case, (1.2) represents an under-determined system of equations and so it is not possible to recover \mathbf{x} unless it has some known structure. In Compressive Sensing, \mathbf{x} is known to be sparse, meaning only a subset of its elements are non-zero. One approach to recover \mathbf{x} is to solve,

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \sqrt{M}\sigma. \quad (1.3)$$

In (1.3), we seek the sparsest \mathbf{x} that is consistent with the measurements by exhaustively searching over all possible solutions. This approach is known to be NP-hard [60] making it intractable even for moderately sized problems. Lasso is a popular alternative approach that recovers \mathbf{x} by solving the ℓ_1 -regularized least squares problem,

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1. \quad (1.4)$$

ISTA [27], AMP [32], VAMP [73], etc., are some algorithms for solving the Lasso optimization problem (1.4). When \mathbf{A} satisfies the Restricted Isometry property (RIP) [34] and \mathbf{x} is sparse enough, there exists theoretical guarantees on the achievable recovery performance of Lasso.

The measurement model in (1.2) is also applicable in several modalities of medical imaging, e.g., Computation Tomography (CT), Magnetic Resonance Imaging (MRI),

etc. In CT, \mathbf{A} is the Radon transform and in MRI \mathbf{A} is a sub-sampled DFT matrix. Images typically admit a sparse representation under a suitable basis, hence a popular approach to solve these imaging problems is using Compressive Sensing. If \mathbf{x} is sparse with respect to basis Φ , then there exists two approaches to recover \mathbf{x} ; i) synthesis, and ii) analysis approach. Formally, the synthesis approach is to solve

$$\hat{\boldsymbol{\alpha}} = \arg \min_{\boldsymbol{\alpha}} \|\mathbf{y} - \mathbf{A}\Phi\boldsymbol{\alpha}\|_2^2 + \lambda\|\boldsymbol{\alpha}\|_1 \quad (1.5a)$$

$$\Rightarrow \hat{\mathbf{x}} = \Phi\hat{\boldsymbol{\alpha}}. \quad (1.5b)$$

Whereas in the analysis approach we solve

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda\|\Phi^H\mathbf{x}\|_1. \quad (1.6)$$

The two approaches are equivalent when Φ is an orthonormal basis. These methods have led to the advancement in computational imaging for several years. A common drawback is that in many practical problems of interest \mathbf{x} may have a more complicated structure than sparsity in a known basis.

To solve (1.4), one usually uses an iterative algorithm where one of the steps in the algorithm involves a proximal operator that enforces prior knowledge of \mathbf{x} . The proximal operator can also be identified as a signal denoiser. It was shown in [79, 95] that we can solve any linear inverse problem if we have a good denoiser for the class of signals we are interested in recovering. In the age of “Big Data”, we may have access to several examples of \mathbf{x} ’s from the signal class that we can use to train a sophisticated denoiser using deep learning [104]. This approach has led to the development of “Plug-and-Play” algorithms achieving superior results.

1.2 Standard Bilinear Model

If the forward operator is not completely known, the inverse problem becomes more challenging. In a bilinear model, the forward operator is a matrix valued linear function,

$$\mathcal{A}_{\mathbf{b}} = \mathbf{A}_0 + \sum_{i=1}^Q b_i \mathbf{A}_i, \quad (1.7)$$

where $\{\mathbf{A}_i\}$ are known matrices and \mathbf{b} are unknown parameters. In some applications, we have tuples of multiple signal and measurements $\{(\mathbf{y}_l, \mathbf{x}_l)\}_{l=1}^L$, also known as the multiple measurement vector (MMV) model. We can stack them column wise to construct the measurement matrix \mathbf{Y} and the signal matrix \mathbf{X}

$$\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_L] \in \mathbb{R}^{M \times L} \quad (1.8a)$$

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_L] \in \mathbb{R}^{N \times L}. \quad (1.8b)$$

The goal is to jointly recover \mathbf{X} and \mathbf{b} from measurements $\mathbf{Y} = \mathcal{A}_{\mathbf{b}}(\mathbf{X}) + \mathbf{W}$, where \mathbf{W} is AWGN noise. It is said to be bilinear because the measurements \mathbf{Y} are linear in \mathbf{X} for fixed \mathbf{b} and linear in \mathbf{b} for fixed \mathbf{X} . This problem has applications in matrix completion [21], robust principle component analysis (RPCA) [20], dictionary learning [81], self-calibration [51], blind deconvolution [47], joint-channel/symbol estimation, compressive sensing with matrix uncertainty [105], and many other tasks.

1.3 Outline and Contribution

In this section, we provide a brief overview of our contributions and the outline of this dissertation.

1.3.1 State Evolution of VAMP

VAMP [73] is a computationally efficient algorithm to solve linear inverse problems. At every iteration t , VAMP generates pseudo-measurements \mathbf{r}_1^t and scalar precision γ_1^t such that for certain large random \mathbf{A} , it ensures the statistical model

$$\mathbf{r}_1^t = \mathbf{x} + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_1^t). \quad (1.9)$$

Moreover, the mean squared error (MSE) performance of VAMP can be exactly predicted by a set of scalar equations known as “state evolution” (SE). These simpler SE equations can be used to study the convergence properties and the accuracy of VAMP in the high dimensional regime for different problems. The state evolution of VAMP was earlier proven for separable Lipschitz denoisers in [73]. In chapter 2, we extend the state evolution of VAMP from separable to a larger class of non-separable Lipschitz denoisers. We also consider bilinear problems such as self-calibration [51] and CS with matrix uncertainty [105], and apply the “lifting” method from [3, 23, 28, 51] to transform it into a linear inverse problem. We then use VAMP to solve the lifted problem and empirically show that the state evolution prediction holds.

1.3.2 The BAdVAMP Algorithm

In chapter 3, we propose the Bilinear Adaptive Vector Approximate Message Passing (BAd-VAMP) algorithm to solve bilinear recovery problems. We combine variational inference with maximum likelihood (ML) estimation to develop our algorithm. We assume the prior density of \mathbf{X} is known up to some parameters $\boldsymbol{\theta}_x$ and the measurement noise is AWGN with unknown noise precision γ_w . Our goal is to estimate the parameters $\boldsymbol{\Theta} \triangleq \{\mathbf{b}, \boldsymbol{\theta}_x, \gamma_w\}$ and recover the unknown matrix $\mathbf{X} \in \mathbb{R}^{N \times L}$ from

measurements $\mathbf{Y} \in \mathbb{R}^{M \times L}$,

$$\mathbf{Y} = \left(\mathbf{A}_0 + \sum_{i=1}^Q b_i \mathbf{A}_i \right) \mathbf{X} + \mathbf{W} \quad (1.10a)$$

$$\mathbf{X} \sim p_{\mathbf{X}}(\cdot; \boldsymbol{\theta}_x) \quad (1.10b)$$

$$\mathbf{W} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1/\gamma_w), \quad (1.10c)$$

where $\{\mathbf{A}_i\}$ are known matrices. We seek to learn the ML estimate of the parameters $\boldsymbol{\Theta}$ and, under that estimate, compute the minimum mean-squared error (MMSE) estimate of \mathbf{X} . We solve several bilinear problems such as dictionary learning, calibration in tomography, CS with matrix uncertainty, etc. We demonstrate numerically that the proposed approach is competitive with other state-of-the-art approaches to bilinear recovery, including lifted VAMP and Bilinear GAMP.

1.3.3 VAMP for MRI

In chapter 4, we consider the problem of image reconstruction in Magnetic Resonance Imaging (MRI). MRI is a non-invasive method to generate high resolution internal images of our body. The MRI scanner produces measurements $\mathbf{y} \in \mathbb{C}^M$ which are noisy sub-sampled Fourier samples of the unknown image \mathbf{x} . The goal is to reconstruct the image $\mathbf{x} \in \mathbb{C}^N$ from $\mathbf{y} \in \mathbb{C}^M$,

$$\mathbf{y} = \mathbf{S}\mathbf{F}\mathbf{x} + \mathbf{w}, \quad (1.11)$$

where \mathbf{F} is the DFT matrix, \mathbf{S} is a sampling matrix and \mathbf{w} is white Gaussian noise. To generate high quality images and reduce the data acquisition time, most modern commercially available MRI scanners use multiple receiver coils in which case the

measurements $\mathbf{y} \in \mathbb{C}^M$ take the form

$$\mathbf{y} = \begin{bmatrix} \mathbf{SFC}_1 \\ \vdots \\ \mathbf{SFC}_K \end{bmatrix} \mathbf{x} + \mathbf{w} \quad (1.12)$$

where $\{\mathbf{C}_k\}_{k=1}^K$ are diagonal matrices representing the coil's sensitivity maps.

We develop algorithms based on the Approximate Message Passing [32] framework using a plug-and-play (PnP) manner. PnP method allows us to use traditional optimization algorithms by calling a powerful image denoiser at every iteration inside the algorithm. The AMP [9] and VAMP [73] algorithms are based on the sensing matrix \mathbf{A} in (1.2) being i.i.d. Gaussian and right rotationally invariant, respectively, which generally doesn't include the sensing matrix in MRI. As a result, these algorithms are known to diverge. In this work, we introduce some adaptive damping schemes to improve their convergence properties in MRI. In VAMP, to solve the LMMSE sub-problem, we use fast iterative solvers such as conjugate gradient [7] or LSQR [65]. In AMP, motivated by the VDAMP [57] algorithm, we use a similar scaling in the k-space which improves the accuracy of the final result. We use the recently released fastMRI [103] dataset for validation.

1.4 Notation

We use boldface uppercase letters to denote matrices (e.g., \mathbf{X}), and we use $\text{tr}\{\mathbf{X}\}$, \mathbf{X}^H , \mathbf{X}^T to denote the trace, conjugate-transpose and transpose of the matrix \mathbf{X} respectively. We denote vectors by boldface lowercase letters (e.g., \mathbf{x}) and scalars by non-bold letters (e.g., x). Given a matrix \mathbf{X} , we use \mathbf{x}_l to denote the l th column and x_{nl} to denote the element in the n th row and l th column. We use $\|\mathbf{X}\|_*$ to denote the nuclear norm of \mathbf{X} , i.e., sum of its singular values. For a square matrix

$\mathbf{X} \in \mathbb{R}^{N \times N}$, we use $\text{diag}(\mathbf{X})$ to denote the vector of elements on the diagonal of \mathbf{X} and we use $\langle \mathbf{X} \rangle \triangleq \frac{1}{N} \sum_{i=1}^N x_{ii}$ to denote the average of the diagonal elements of \mathbf{X} . For a random vector \mathbf{x} with density b , we use $\mathbb{E}[f(\mathbf{x})|b]$ to denote the expectation of $f(\mathbf{x})$, i.e., $\mathbb{E}[f(\mathbf{x})|b] = \int f(\mathbf{x})b(\mathbf{x}) d\mathbf{x}$, and we use $\text{var}[f(\mathbf{x})|b]$ for the corresponding variance. We denote the density of a Gaussian random vector with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ as $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. We use $\text{Diag}(\mathbf{x})$ to denote the diagonal matrix created from vector \mathbf{x} .

Chapter 2: Vector Approximate Message Passing

2.1 Introduction

In this chapter, we will discuss a computationally efficient algorithm to solve the linear inverse problems mentioned in chapter 1. Recall that in linear inverse problems, our goal is to recover the unknown vector $\mathbf{x} \in \mathbb{R}^N$ from measurements $\mathbf{y} \in \mathbb{R}^M$,

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}, \quad (2.1)$$

where $\mathbf{A} \in \mathbb{R}^{M \times N}$ is a known matrix and $\mathbf{w} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \gamma_w^{-1})$ is additive white Gaussian noise of precision γ_w . We assume that the unknown vector \mathbf{x} has a known prior density $p_{\mathbf{x}}(\cdot)$ which encodes all our prior knowledge about \mathbf{x} . We aim to compute the minimum mean-squared error (MMSE) estimate of \mathbf{x} , i.e.,

$$\hat{\mathbf{x}}_{\text{MMSE}} \triangleq \mathbb{E}[\mathbf{x}|\mathbf{y}]. \quad (2.2)$$

The expectation is taken over the posterior density

$$p_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) = \frac{1}{Z(\mathbf{y})} p_{\mathbf{x}}(\mathbf{x}) p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}), \quad (2.3)$$

using the likelihood function

$$p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \gamma_w^{-1}\mathbf{I}) \quad (2.4)$$

and the normalization constant $Z(\mathbf{y})$ given as

$$Z(\mathbf{y}) = p_{\mathbf{y}}(\mathbf{y}) \tag{2.5a}$$

$$= \int p_{\mathbf{x}}(\mathbf{x}) p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) d\mathbf{x}. \tag{2.5b}$$

For high dimensional \mathbf{x} , the integrals in (2.2) and (2.5b) are difficult to compute directly. Thus we need to use alternate methods. Vector Approximate Message Passing (VAMP) [75] is a computationally efficient iterative algorithm to compute the *maximum a posteriori* (MAP) estimate or approximately compute the MMSE estimate of the unknown \mathbf{x} in case of infinitely large right rotationally invariant \mathbf{A} .

2.2 Variational Inference

Variational inference (VI) [98] can be used to bypass the computation of $Z(\mathbf{y})$. For example, notice that the true posterior $p_{\mathbf{x}|\mathbf{y}}$ can be recovered by solving the variational optimization (over densities)

$$\hat{q} = \arg \min_q D_{\text{KL}}(q \| p_{\mathbf{x}|\mathbf{y}}), \tag{2.6}$$

where $D_{\text{KL}}(q \| p)$ denotes the KL divergence from p to q , i.e.,

$$D_{\text{KL}}(q \| p) \triangleq \int q(\mathbf{x}) \ln \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x}. \tag{2.7}$$

Substituting $p_{\mathbf{x}|\mathbf{y}}$ from (2.3) into (2.7) gives us

$$D_{\text{KL}}(q \| p_{\mathbf{x}|\mathbf{y}}) = D_{\text{KL}}(q \| p_{\mathbf{x}}) + D_{\text{KL}}(q \| p_{\mathbf{y}|\mathbf{x}}) + H(q) + \ln Z(\mathbf{y}) \tag{2.8}$$

where $H(q) \triangleq - \int q(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x}$ is the differential entropy of the random vector \mathbf{x} with density q . Thus it follows from (2.6) and (2.8) that

$$\hat{q} = \arg \min_q \{ D_{\text{KL}}(q \| p_{\mathbf{x}}) + D_{\text{KL}}(q \| p_{\mathbf{y}|\mathbf{x}}) + H(q) \}, \tag{2.9}$$

which bypasses $Z(\mathbf{y})$. Still, solving (2.9) is difficult in most cases of interest. The typical response is to constrain q to be in a certain family of distributions, \mathcal{F} , by solving

$$\hat{q} = \arg \min_{q \in \mathcal{F}} \{D_{\text{KL}}(q \| p_{\mathbf{x}}) + D_{\text{KL}}(q \| p_{\mathbf{y}|\mathbf{x}}) + H(q)\}. \quad (2.10)$$

Some popular choices of \mathcal{F} are,

- i. *Mean Field Approximation*: family of separable densities

$$\mathcal{F} = \{q(\cdot) \mid q(\mathbf{x}) = \prod_{i=1}^N q_i(x_i)\} \quad (2.11)$$

- ii. *Exponential Family*: specified by sufficient statistics $\phi(\cdot)$ and parameters η ,

$$\mathcal{F} = \{q(\cdot) \mid q(\mathbf{x}) = \exp\{\eta^\top \phi(\mathbf{x}) - \ln Z(\eta)\}\} \quad (2.12)$$

which includes normal, Dirichlet, gamma, etc.

One should choose \mathcal{F} such that (2.10) can be solved analytically and it should be powerful enough to capture the useful components of the true posterior $p_{\mathbf{x}|\mathbf{y}}$. Since a bad choice of \mathcal{F} can compromise \hat{q} in (2.10) and its mean, we take a different approach.

2.3 Expectation Consistent Approximate Inference

Using the ‘‘Gibbs free energy’’

$$J(q_1, q_2, q_3) \triangleq D_{\text{KL}}(q_1 \| p_{\mathbf{x}}) + D_{\text{KL}}(q_2 \| p_{\mathbf{y}|\mathbf{x}}) + H(q_3), \quad (2.13)$$

one can rewrite (2.9) as

$$\arg \min_{q_1} \min_{q_2} \max_{q_3} J(q_1, q_2, q_3) \quad (2.14a)$$

$$\text{s.t. } q_1 = q_2 = q_3. \quad (2.14b)$$

Algorithm 1 EC algorithm [75]

```

1: initialize:
    $\mathbf{r}_1^0, \gamma_1^0$ 
2: for  $t = 0, \dots, T_{\max}$  do
3:    $\mathbf{x}_1^t = \mathbf{g}_1(\mathbf{r}_1^t, \gamma_1^t)$ 
4:    $1/\eta_1^t = \langle \mathbf{g}'_1(\mathbf{r}_1^t, \gamma_1^t) \rangle / \gamma_1^t$ 
5:    $\gamma_2^t = \eta_1^t - \gamma_1^t$ 
6:    $\mathbf{r}_2^t = (\eta_1^t \mathbf{x}_1^t - \gamma_1^t \mathbf{r}_1^t) / \gamma_2^t$ 
7:    $\mathbf{x}_2^t = \mathbf{g}_2(\mathbf{r}_2^t, \gamma_2^t)$ 
8:    $1/\eta_2^t = \langle \mathbf{g}'_2(\mathbf{r}_2^t, \gamma_2^t) \rangle / \gamma_2^t$ 
9:    $\gamma_1^{t+1} = \eta_2^t - \gamma_2^t$ 
10:   $\mathbf{r}_1^{t+1} = (\eta_2^t \mathbf{x}_2^t - \gamma_2^t \mathbf{r}_2^t) / (\eta_2^t - \gamma_2^t)$ 
11: end for

```

We minimize over q_1 and q_2 because $D_{\text{KL}}(q_1 \| p_{\mathbf{x}})$ and $D_{\text{KL}}(q_2 \| p_{\mathbf{y}|\mathbf{x}})$ are convex, while we maximize over q_3 because $H(q_3)$ is concave. But, as discussed earlier, (2.14) is difficult to solve.

In the *expectation consistent approximate inference* (EC) scheme proposed by Oppé and Winther in [63], the density constraint (2.14b) is relaxed to the first and second moment matching constraints, i.e.,

$$\arg \min_{q_1} \min_{q_2} \max_{q_3} J(q_1, q_2, q_3) \quad (2.15a)$$

$$\text{s.t. } \mathbb{E}[\mathbf{x}|q_1] = \mathbb{E}[\mathbf{x}|q_2] = \mathbb{E}[\mathbf{x}|q_3] \quad (2.15b)$$

$$\text{tr}\{\text{Cov}[\mathbf{x}|q_1]\} = \text{tr}\{\text{Cov}[\mathbf{x}|q_2]\} = \text{tr}\{\text{Cov}[\mathbf{x}|q_3]\}, \quad (2.15c)$$

where $\mathbb{E}[\mathbf{x}|q_i]$ denotes $\mathbb{E}[\mathbf{x}]$ under $\mathbf{x} \sim q_i$. This yields stationary points of the form

$$q_1(\mathbf{x}) \propto p_{\mathbf{x}}(\mathbf{x}) \exp\left(-\frac{\gamma_1}{2} \|\mathbf{x} - \mathbf{r}_1\|_2^2\right) \quad (2.16a)$$

$$q_2(\mathbf{x}) \propto p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \exp\left(-\frac{\gamma_2}{2} \|\mathbf{x} - \mathbf{r}_2\|_2^2\right) \quad (2.16b)$$

$$q_3(\mathbf{x}) \propto \exp\left(-\frac{\eta}{2} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2\right), \quad (2.16c)$$

for $\{\mathbf{r}_1, \gamma_1, \mathbf{r}_2, \gamma_2, \hat{\mathbf{x}}, \eta\}$ that lead to satisfaction of (2.15b)-(2.15c). There are several algorithms whose fixed points are the stationary point of EC (e.g., ADATAP [64], S-AMP [17]).

Various approaches can be used to solve for $\{\mathbf{r}_1, \gamma_1, \mathbf{r}_2, \gamma_2, \hat{\mathbf{x}}, \eta\}$. One is to alternate the update of $\{(\mathbf{r}_1, \gamma_1), (\hat{\mathbf{x}}, \eta)\}$ and $\{(\mathbf{r}_2, \gamma_2), (\hat{\mathbf{x}}, \eta)\}$ such that, at each iteration, the moments of q_3 are consistent with either q_1 or q_2 . This approach is summarized in Alg. 1 using

$$\mathbf{g}_1(\mathbf{r}_1, \gamma_1) \triangleq \frac{\int \mathbf{x} p_{\mathbf{x}}(\mathbf{x}) \mathcal{N}(\mathbf{x}; \mathbf{r}_1, \mathbf{I}/\gamma_1) d\mathbf{x}}{\int p_{\mathbf{x}}(\mathbf{x}) \mathcal{N}(\mathbf{x}; \mathbf{r}_1, \mathbf{I}/\gamma_1) d\mathbf{x}} \quad (2.17)$$

$$\mathbf{g}_2(\mathbf{r}_2, \gamma_2) \triangleq \frac{\int \mathbf{x} p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \mathcal{N}(\mathbf{x}; \mathbf{r}_2, \mathbf{I}/\gamma_2) d\mathbf{x}}{\int p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \mathcal{N}(\mathbf{x}; \mathbf{r}_2, \mathbf{I}/\gamma_2) d\mathbf{x}}, \quad (2.18)$$

which, under these definitions of \mathbf{g}_1 and \mathbf{g}_2 , can be recognized as an instance of *expectation propagation* (EP) [58, Sec. 3.2], [41, 86]. In lines 4 and 8 of Alg. 1, $\mathbf{g}'_i(\mathbf{r}_i, \gamma_i) \in \mathbb{R}^N$ denotes the diagonal of the Jacobian matrix of $\mathbf{g}_i(\cdot, \gamma_i)$ at \mathbf{r}_i , i.e.,

$$\mathbf{g}'_i(\mathbf{r}_i, \gamma_i) \triangleq \text{diag} \left(\frac{\partial \mathbf{g}_i(\mathbf{r}_i, \gamma_i)}{\partial \mathbf{r}_i} \right), \quad (2.19)$$

and $\langle \mathbf{x} \rangle$ denotes the average coefficient value, i.e., $\langle \mathbf{x} \rangle \triangleq \frac{1}{N} \sum_{i=1}^N x_i$ for $\mathbf{x} \in \mathbb{R}^N$.

Recall the likelihood function for the standard linear model,

$$p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \gamma_w^{-1} \mathbf{I}). \quad (2.20)$$

Using the definition of $\mathbf{g}_2(\cdot)$ in (2.18) and the likelihood function in (2.20), we see that \mathbf{g}_2 is performing an LMMSE estimation whose closed form solution is

$$\mathbf{g}_2(\mathbf{r}_2, \gamma_2) = (\gamma_2 \mathbf{I} + \gamma_w \mathbf{A}^\top \mathbf{A})^{-1} (\gamma_2 \mathbf{r}_2 + \gamma_w \mathbf{A}^\top \mathbf{y}) \quad (2.21)$$

$$\langle \mathbf{g}'_2(\mathbf{r}_2, \gamma_2) \rangle = \gamma_2 \text{tr} \{ (\gamma_2 \mathbf{I} + \gamma_w \mathbf{A}^\top \mathbf{A})^{-1} \} / N. \quad (2.22)$$

Meanwhile, the form of \mathbf{g}_1 depends on $p_{\mathbf{x}}$ through (2.17).

Algorithm 2 EP algorithm (SVD Form) [75]

```

1: define:
    Economy SVD:  $\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^\top \in \mathbb{R}^{M \times N}$ 
     $R = \text{rank}(\mathbf{A})$ 
     $\tilde{\mathbf{y}} = \text{Diag}(\mathbf{s})^{-1} \mathbf{U}^\top \mathbf{y}$ 
2: initialize:
     $\mathbf{r}^0, \gamma^0$ 
3: for  $t = 0, \dots, T_{\max}$  do
4:    $\mathbf{x}^t = \mathbf{g}_1(\mathbf{r}^t, \gamma^t)$ 
5:    $\alpha^t = \langle \mathbf{g}'_1(\mathbf{r}^t, \gamma^t) \rangle$ 
6:    $\tilde{\mathbf{r}}^t = (\mathbf{x}^t - \alpha^t \mathbf{r}^t) / (1 - \alpha^t)$ 
7:    $\tilde{\gamma}^t = \gamma^t (1 - \alpha^t) / \alpha^t$ 
8:    $\mathbf{d}^t = \gamma_w \text{Diag}(\gamma_w \mathbf{s}^2 + \tilde{\gamma}^t \mathbf{1})^{-1} \mathbf{s}^2$ 
9:    $\gamma^{t+1} = \tilde{\gamma}^t R \langle \mathbf{d}^t \rangle / (N - R \langle \mathbf{d}^t \rangle)$ 
10:   $\mathbf{r}^{t+1} = \tilde{\mathbf{r}}^t + \frac{N}{R} \mathbf{V} \text{Diag}(\mathbf{d}^t / \langle \mathbf{d}^t \rangle) (\tilde{\mathbf{y}} - \mathbf{V}^\top \tilde{\mathbf{r}}^t)$ 
11: end for

```

Using the definition of \mathbf{g}_2 in (2.21), we see that in line 7 of Alg. 1 a matrix inversion of complexity $O(N^3)$ is needed at every iteration, which can significantly slow down the algorithm. However, if we were to pre-compute the (economy) SVD $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$, then the per iteration cost of VAMP is dominated by $O(RN)$, where R is the rank of \mathbf{A} . The SVD version of Alg. 1 is summarized in Alg. 2.

Based on the description above, one might wonder whether the EC stationary point $\hat{\mathbf{x}} = \mathbb{E}[\mathbf{x}|q_1] = \mathbb{E}[\mathbf{x}|q_2] = \mathbb{E}[\mathbf{x}|q_3]$ is a good approximation of the true conditional mean $\mathbb{E}[\mathbf{x}|\mathbf{y}]$, and additionally one might question whether Algorithm 1 converges to this $\hat{\mathbf{x}}$. Both of these concerns were resolved in the VAMP paper [75]. In particular, [75] showed that, when \mathbf{A} is right rotationally invariant and asymptotically large, the per-iteration behavior of Alg. 1 with \mathbf{g}_2 from (2.21) and Lipschitz \mathbf{g}_1 is exactly predicted by a scalar state evolution. Furthermore, in the case where \mathbf{g}_1 is matched to $p_{\mathbf{x}}$, the prior of \mathbf{x} , as in (2.17), and where \mathbf{g}_2 uses the true AWGN precision $\gamma_w < \infty$

as in (2.18), the MSE of the fixed point $\hat{\mathbf{x}}$ of Alg. 1 was shown in [75] to match the MMSE predicted by the replica method [93]. This replica prediction is conjectured to be correct [77], in which case the $\hat{\mathbf{x}}$ generated by Alg. 1 under (2.17) and (2.21) will be MMSE for infinitely large, right rotationally invariant \mathbf{A} when the state evolution has unique fixed points. Note that, for infinitely large i.i.d. \mathbf{A} , the replica prediction has been proven to be correct [6, 78].

In the sequel, we will refer to Alg. 1 with generic Lipschitz \mathbf{g}_1 as the VAMP algorithm, noting that it coincides with EP in the special case of Bayesian \mathbf{g}_1 from (2.17). VAMP is more general than EP because it can be used with denoisers \mathbf{g}_1 that are MAP [36], or that have no probabilistic interpretation, and still lead to precisely predictable behavior under infinitely large right rotationally invariant \mathbf{A} [35, 75].

In Algorithm 1, lines 6 and 10 perform what is referred to as *Onsager correction* in statistical physics [90]. It ensures that, under a certain large random \mathbf{A} , the VAMP quantity $(\mathbf{r}_1^t, \gamma_1^t)$ obeys the following AWGN statistical model at every iteration:

$$\mathbf{r}_1^t = \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, \gamma_1^t \mathbf{I}), \quad (2.23)$$

where \mathbf{x}_0 is the true vector we are interested in recovering in (2.1). This AWGN-like property is one of the main reasons for VAMP's fast convergence properties. There are several ways to interpret Algorithm 1; one can interpret it as an instance of expectation propagation (EP) [63, 89], or as an instance of belief-propagation (BP) on a non-loopy factor graph with vector-valued nodes [73], or as a variant of alternating direction method of multipliers (ADMM) [14] involving two dual updates with adaptive step size [36]. Another intuitive interpretation of Algorithm 1 is the following:

(i) In line 3, the function \mathbf{g}_1 estimates the signal assuming from AWGN-corrupted measurements $\mathbf{r}_1^t = \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_1^t)$.

(ii) Lines 5 & 6 compute $\mathbf{r}_2^t, \gamma_2^t$ to yield the model $\mathbf{x}_0 = \mathbf{r}_2^t + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_2^t)$.

(iii) In line 7, \mathbf{g}_2 performs LMMSE estimation of \mathbf{x} from \mathbf{y} under the model

$$\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w), \quad \mathbf{x}_0 = \mathbf{r}_2^t + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_2^t).$$

(iv) Lines 9 & 10 compute $\mathbf{r}_1^{t+1}, \gamma_1^{t+1}$ such that the AWGN like model in step (i) holds.

The above discussion suggests that \mathbf{g}_1 is performing a denoising operation. In the next section, we provide a rigorous analysis of the AWGN-like property of VAMP.

2.4 State Evolution

Under some conditions, the noise precision in the AWGN model in (2.23) can be exactly predicted by a set of scalar recursive equations known as the “state-evolution” (SE). This state-evolution formalism allows us to predict the MSE performance of VAMP at every iteration and study its asymptotic convergence properties. In this section, we introduce the SE formalism in VAMP and the conditions under which it is guaranteed to be true. The SE of VAMP was first proven for separable \mathbf{g}_1 in [73]. The following results are from our work in [35], where we extended the SE analysis to the non-separable case. We first briefly describe the “large system limit” model and then state the main results.

2.4.1 Large System Limit

We consider a sequence of problems defined by the problem size, N . For each N , we assume $\mathbf{x}_0 \in \mathbb{R}^N$ is a “true” vector that we are interested in recovering from measurements $\mathbf{y} \in \mathbb{R}^M$ generated according to the linear model,

$$\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \gamma_{w0}^{-1}\mathbf{I}). \quad (2.24)$$

The “true” noise precision is γ_{w0} , which could be different from the assumed noise precision used in (2.21) and (2.22). Here, we assume $M = N$ without loss of generality.

We can easily extend the following results to the case $M \neq N$, by zero padding \mathbf{s} [73].

We assume the SVD of \mathbf{A} takes the form,

$$\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^\top, \quad \mathbf{s} = (s_1, \dots, s_N), \quad (2.25)$$

where $\mathbf{U} \in \mathbb{R}^{N \times N}$ is any arbitrary orthogonal matrix, \mathbf{s} is i.i.d. with components $s_i \in [0, s_{\max}]$ almost surely, and $\mathbf{V} \in \mathbb{R}^{N \times N}$ is Haar distributed. The latter means that \mathbf{V} is uniformly distributed on the set of $N \times N$ orthogonal matrices, which implies \mathbf{A} is *right rotationally invariant*, meaning if we right multiply \mathbf{A} by any deterministic orthogonal matrix \mathbf{V}_0 then the distribution of $\mathbf{A}\mathbf{V}_0$ is the same as \mathbf{A} , i.e., $\mathbf{A}\mathbf{V}_0 \stackrel{d}{=} \mathbf{A}$. We also assume that the quantities $\{\mathbf{x}_0, \mathbf{w}, \mathbf{s}, \mathbf{V}\}$ are all independent. In the large system limit, we let $N \rightarrow \infty$ and study the performance of VAMP averaged over the random quantities $\{\mathbf{w}, \mathbf{s}, \mathbf{V}\}$. We next introduce two important definitions the denoiser \mathbf{g}_1 needs to satisfy.

Definition 1. A denoiser $\mathbf{g}(\cdot, \cdot)$ is uniformly Lipschitz continuous if there exists constants A, B and $C > 0$, such that $\forall \mathbf{r}_1, \mathbf{r}_2 \in \mathbb{R}^N, \gamma_1, \gamma_2 > 0$ and N , we have,

$$\|\mathbf{g}(\mathbf{r}_2, \gamma_2) - \mathbf{g}(\mathbf{r}_1, \gamma_1)\| \leq (A + B|\gamma_2 - \gamma_1|)\|\mathbf{r}_2 - \mathbf{r}_1\| + C\sqrt{N}|\gamma_2 - \gamma_1|. \quad (2.26)$$

Definition 2. A denoiser $\mathbf{g}(\cdot, \cdot)$ and vector $\mathbf{u} \in \mathbb{R}^N$ are said to be convergent under Gaussian noise if the following condition holds: Let $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^N$ be two i.i.d. sequences with $\begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{S} \otimes \mathbf{I})$ for some positive definite matrix $\mathbf{S} \in \mathbb{R}^{2 \times 2}$. Then, all the following limits must exist almost surely $\forall \gamma_1, \gamma_2 > 0, \mathbf{S}$:

$$\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{g}(\mathbf{u} + \mathbf{z}_1, \gamma_1)^T \mathbf{g}(\mathbf{u} + \mathbf{z}_2, \gamma_2), \quad \lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{g}(\mathbf{u} + \mathbf{z}_1, \gamma_1)^T \mathbf{u}, \quad (2.27a)$$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{u}^T \mathbf{z}_1, \quad \lim_{N \rightarrow \infty} \frac{1}{N} \|\mathbf{u}\|^2 \quad (2.27b)$$

$$\lim_{N \rightarrow \infty} \langle \nabla \mathbf{g}(\mathbf{u} + \mathbf{z}_1, \gamma_1) \rangle = \frac{1}{N S_{12}} \mathbf{g}(\mathbf{u} + \mathbf{z}_1, \gamma_1)^T \mathbf{z}_2. \quad (2.27c)$$

Moreover, the limiting values are continuous in \mathbf{S} , γ_1 and γ_2 .

Several denoisers used in practice satisfy the above definitions. In Appendix A.1, we show that the singular value thresholding denoiser, commonly used in low rank matrix recovery [16], satisfy (2.26) and (2.27a).

Given a denoiser satisfying definitions 1 and 2, and an assumed noise precision γ_w , we run Algorithm 1 with the initialization,

$$\mathbf{r}_1^0 = \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, \tau_1^0 \mathbf{I}), \quad (2.28)$$

for some initial error variance τ_1^0 . In addition, we assume

$$\lim_{N \rightarrow \infty} \tau_1^0 = \bar{\tau}_1^0, \quad (2.29)$$

almost surely for some $\bar{\tau}_1^0 \geq 0$. To measure the performance of VAMP, we define two error functions, $\mathcal{E}_1(\cdot)$ and $\mathcal{E}_2(\cdot)$. The error function $\mathcal{E}_1(\cdot)$,

$$\mathcal{E}_1(\gamma_1, \tau_1) \triangleq \lim_{N \rightarrow \infty} \frac{1}{N} \|\mathbf{g}_1(\mathbf{x}_0 + \mathbf{z}_1, \gamma_1) - \mathbf{x}_0\|^2 \quad (2.30a)$$

$$\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \tau_1 \mathbf{I}) \quad (2.30b)$$

characterizes the MSE of the denoiser $\mathbf{g}_1(\cdot, \cdot)$ whereas $\mathcal{E}_2(\cdot)$,

$$\mathcal{E}_2(\gamma_2, \tau_2) \triangleq \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E} \|\mathbf{g}_2(\mathbf{x}_0 + \mathbf{z}_2, \gamma_2) - \mathbf{x}_0\|^2 \quad (2.31a)$$

$$\mathbf{z}_2 \sim \mathcal{N}(\mathbf{0}, \tau_2 \mathbf{I}) \quad (2.31b)$$

$$\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathcal{N}(\mathbf{0}, \gamma_{w_0}^{-1} \mathbf{I}). \quad (2.31c)$$

characterizes the MSE of the linear estimator $\mathbf{g}_2(\cdot, \cdot)$. The limit (2.30b) exists almost surely due to the assumption of $\mathbf{g}_1(\cdot, \cdot)$ being convergent under Gaussian noise. We also define the *sensitivity functions* as

$$\mathcal{A}_i(\gamma_i, \tau_i) := \lim_{N \rightarrow \infty} \langle \nabla \mathbf{g}_i(\mathbf{x}_0 + \mathbf{z}_i, \gamma_i) \rangle, \quad \mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \tau_i \mathbf{I}). \quad (2.32)$$

2.4.2 State Evolution of Non-Separable VAMP

Consider the following error vectors in VAMP,

$$\mathbf{p}^t \triangleq \mathbf{r}_1^t - \mathbf{x}_0 \quad (2.33a)$$

$$\mathbf{q}^t \triangleq \mathbf{V}^\top (\mathbf{r}_2^t - \mathbf{x}_0). \quad (2.33b)$$

At iteration t , the error in the input to the denoiser $\mathbf{g}_1(\cdot, \cdot)$ is \mathbf{p}^t and the transformed error in the input to the LMMSE estimator $\mathbf{g}_2(\cdot, \cdot)$ is \mathbf{q}^t . In the state evolution analysis, we show that asymptotically these errors are Gaussian. Moreover, we can also exactly predict the error precisions (2.33) via the following recursive scalar *state evolution* equations initialized with $t = 0$, τ_1^0 in (2.28) and $\bar{\gamma}_1^0$ in (2.29):

$$\bar{\alpha}_1^t = \mathcal{A}_1(\bar{\gamma}_1^t, \tau_1^t), \quad \bar{\eta}_1^t = \frac{\bar{\gamma}_1^t}{\bar{\alpha}_1^t}, \quad \bar{\gamma}_2^t = \bar{\eta}_1^t - \bar{\gamma}_1^t \quad (2.34a)$$

$$\tau_2^t = \frac{1}{(1 - \bar{\alpha}_1^t)^2} [\mathcal{E}_1(\bar{\gamma}_1^t, \tau_1^t) - (\bar{\alpha}_1^t)^2 \tau_1^t], \quad (2.34b)$$

$$\bar{\alpha}_2^t = \mathcal{A}_2(\bar{\gamma}_2^t, \tau_2^t), \quad \bar{\eta}_2^t = \frac{\bar{\gamma}_2^t}{\bar{\alpha}_2^t}, \quad \bar{\gamma}_1^{t+1} = \bar{\eta}_2^t - \bar{\gamma}_2^t \quad (2.34c)$$

$$\tau_1^{t+1} = \frac{1}{(1 - \bar{\alpha}_2^t)^2} [\mathcal{E}_2(\bar{\gamma}_2^t, \tau_2^t) - (\bar{\alpha}_2^t)^2 \tau_2^t] \quad (2.34d)$$

We can now state our main result.

Theorem 1. *Assume that the true vector \mathbf{x}_0 and the denoiser $\mathbf{g}_2(\cdot, \cdot)$ satisfies definitions 1 and 2. Also assume that for all iterations t , the solution $\bar{\alpha}_1^t$ from the SE equations (2.34) satisfies $\bar{\alpha}_1^t \in (0, 1)$ and $\bar{\gamma}_{ik} > 0$. Then,*

(a) *The error vectors \mathbf{p}^t and \mathbf{q}^t defined in (2.33) satisfy,*

$$\lim_{N \rightarrow \infty} \frac{1}{N} \|\mathbf{p}^t - \tilde{\mathbf{p}}^t\|^2 = 0 \quad (2.35)$$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \|\mathbf{q}^t - \tilde{\mathbf{q}}^t\|^2 = 0 \quad (2.36)$$

where, $\tilde{\mathbf{p}}^t \sim \mathcal{N}(\mathbf{0}, \tau_1^t \mathbf{I})$, $\tilde{\mathbf{q}}^t \sim \mathcal{N}(\mathbf{0}, \tau_2^t \mathbf{I})$ and $\tilde{\mathbf{p}}^t, \tilde{\mathbf{q}}^t$ are independent.

(b) *For any iteration t , and $i = 1, 2$, we have, almost surely*

$$\lim_{N \rightarrow \infty} \frac{1}{N} \|\mathbf{x}_i^t - \mathbf{x}_0\|^2 = \frac{1}{\bar{\eta}_i^t}, \quad \lim_{N \rightarrow \infty} (\alpha_i^t, \eta_i^t, \gamma_i^t) = (\bar{\alpha}_i^t, \bar{\eta}_i^t, \bar{\gamma}_i^t). \quad (2.37)$$

Proof. See [35]. □

Theorem 1 thus shows that the error vectors \mathbf{p}^t and \mathbf{q}^t in (2.33) are approximately i.i.d. Gaussian in a certain high dimensional limit. This result is an extension of the main result in [73] from separable denoisers to non-separable denoisers. The error variances and the mean squared error (MSE) of the VAMP estimates $(\mathbf{x}_1^t, \mathbf{x}_2^t)$ can be exactly predicted by (2.34), which are the same as the SE equations in [73].

2.5 Numerical Experiments

2.5.1 Image Recovery via Non-Separable VAMP

We consider the problem of recovering an image $\mathbf{x}_0 \in \mathbb{R}^N$ from measurements $\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{w} \in \mathbb{R}^M$ when $M \ll N$. This problem is commonly seen in computed



Figure 2.1: Image reconstruction using BM3D-VAMP at undersampling ratio of 0.5.

tomography, magnetic resonance imaging, etc., where \mathbf{A} depends on the application. In this experiment, we provide empirical evidence of the state evolution in VAMP. We randomly generated $\mathbf{A} \in \mathbb{R}^{M \times N}$ from i.i.d. $\mathcal{N}(0, 1)$ for $M/N = 0.5$ and set the noise precision γ_w that achieved 40 dB signal to noise ratio (SNR),

$$\text{SNR} \triangleq \frac{\mathbb{E} \|\mathbf{A}\mathbf{x}_0\|^2}{\mathbb{E} \|\mathbf{w}\|^2}. \quad (2.38)$$

We used the *Barbara* image of size 128×128 (i.e., $N = 128^2$) as the true vector \mathbf{x}_0 and BM3D [26] as the denoiser $\mathbf{g}_1(\cdot, \cdot)$ in VAMP. Figure 2.2 shows the PSNR¹ predicted by the state evolution and the PSNR in VAMP at every iteration averaged over 10 random draws of \mathbf{A} and \mathbf{w} . The state evolution accurately predicts the PSNR in VAMP at every iteration.

¹We define the PSNR as $\frac{\max_i |x_{0,i}|}{\|\hat{\mathbf{x}} - \mathbf{x}_0\|}$.

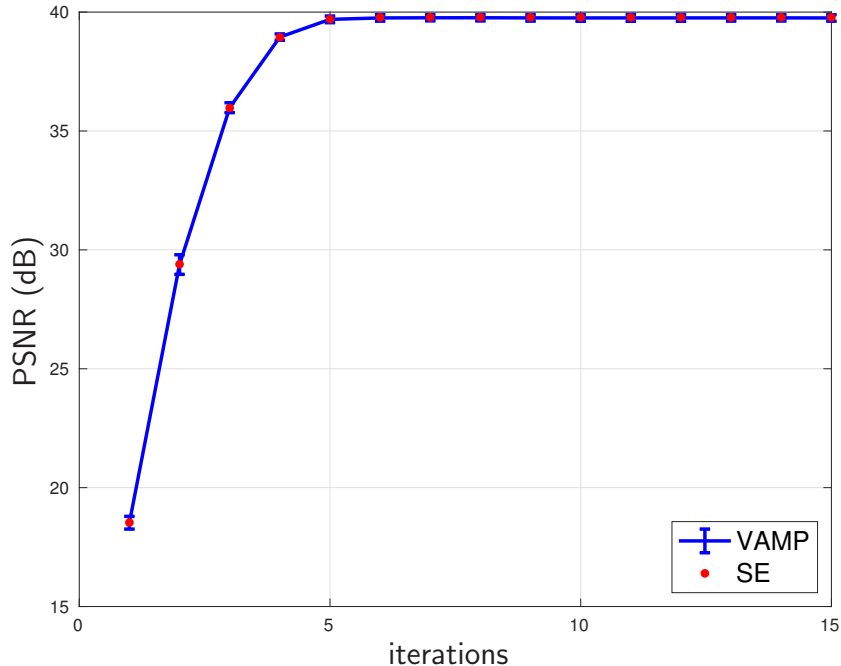


Figure 2.2: SE prediction in VAMP for image recovery

2.5.2 Bilinear Estimation via Lifting

We now use the standard linear inverse model (2.1) to tackle problems in *bilinear* estimation through a technique known as “lifting” [3, 23, 28, 51]. In doing so, we are motivated by applications like blind deconvolution [11], self-calibration [51], compressed sensing (CS) with matrix uncertainty [105] and joint channel-symbol estimation [88]. All cases yield measurements \mathbf{y} of the form

$$\mathbf{y} = \left(\sum_{l=1}^L b_l \mathbf{A}_l \right) \mathbf{c} + \mathbf{w} \in \mathbb{R}^M, \quad (2.39)$$

where $\{\mathbf{A}_l\}_{l=1}^L$ are known, $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w)$, and the objective is to recover both $\mathbf{b} \triangleq [b_1, \dots, b_L]^\top$ and $\mathbf{c} \in \mathbb{R}^P$. This bilinear problem can be “lifted” into a linear

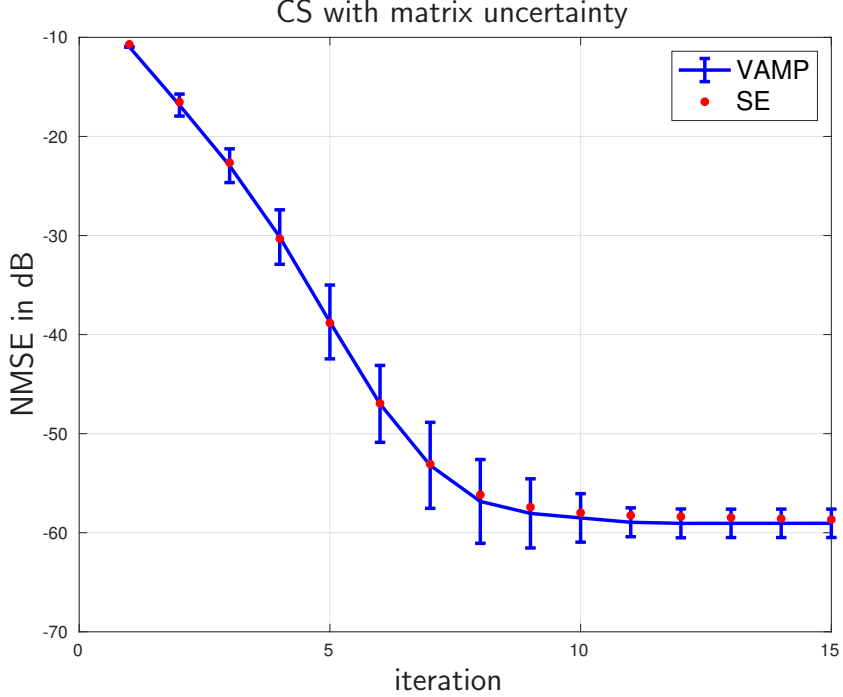


Figure 2.3: SE prediction in VAMP for CS with matrix uncertainty

problem of the form (2.1) by setting

$$\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_L] \in \mathbb{R}^{M \times LP} \quad (2.40a)$$

$$\mathbf{x} = \text{vec}(\mathbf{c}\mathbf{b}^\top) \in \mathbb{R}^{LP}, \quad (2.40b)$$

where $\text{vec}(\mathbf{X})$ vectorizes \mathbf{X} by concatenating its columns. When \mathbf{b} and \mathbf{c} are i.i.d. with known priors, the MMSE denoiser $\mathbf{g}_1(\mathbf{r}_1, \gamma_1) = \mathbb{E}(\mathbf{x}|\mathbf{r}_1 = \mathbf{x} + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_1))$ can be implemented near-optimally by the rank-one AMP algorithm from [72] (see also [30, 48, 55]), with divergence estimated as in [56].

We first consider *CS with matrix uncertainty* [105], where b_1 is known. For these experiments, we generated the unknown $\{b_i\}_{i=2}^L$ as i.i.d. $\mathcal{N}(0, 1)$ and the unknown $\mathbf{c} \in \mathbb{R}^P$ as K -sparse with $\mathcal{N}(0, 1)$ nonzero entries. For $b_1 = \sqrt{20}$, $L = 11$, $P = 256$, $K = 10$, i.i.d. $\mathcal{N}(0, 1)$ matrices $\{\mathbf{A}_i\}$, and $\text{SNR} = 40$ dB, Fig. 2.3 shows that the

NMSE on \mathbf{x} of lifted VAMP is very close to its SE prediction. We then compared lifted VAMP to PBiGAMP from [66], which applies AMP directly to the (non-lifted) bilinear problem, and to WSS-TLS from [105], which uses non-convex optimization. We also compared to MMSE estimation of \mathbf{b} under oracle knowledge of \mathbf{c} , and MMSE estimation of \mathbf{c} under oracle knowledge of $\text{support}(\mathbf{c})$ and \mathbf{b} . Fig. 2.4 shows the normalized MSE on \mathbf{b} and \mathbf{c} versus sampling ratio M/P ,

$$\text{NMSE}(\hat{\mathbf{b}}) \triangleq \frac{\mathbb{E} \|\hat{\mathbf{b}} - \mathbf{b}\|^2}{\mathbb{E} \|\mathbf{b}\|^2}. \quad (2.41)$$

This figure demonstrates that lifted VAMP and PBiGAMP perform close to the oracles and much better than WSS-TLS.

Although lifted VAMP performs similarly to PBiGAMP in Fig. 2.4, its advantage over PBiGAMP becomes apparent with non-i.i.d. \mathbf{A} . For illustration, we repeated the previous experiment, but with \mathbf{A} constructed using the SVD $\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^T$ with Haar distributed \mathbf{U} and \mathbf{V} and geometrically spaced \mathbf{s} . Also, to make the problem more difficult, we set $b_1 = 1$. Figure 2.5 shows the normalized MSE on \mathbf{b} and \mathbf{c} versus $\text{cond}(\mathbf{A})$ at $M/P = 0.6$. There it can be seen that lifted VAMP is much more robust than PBiGAMP to the conditioning of \mathbf{A} .

We next consider the *self-calibration* problem [51], where the measurements take the form

$$\mathbf{y} = \text{Diag}(\mathbf{H}\mathbf{b})\Psi\mathbf{c} + \mathbf{w} \in \mathbb{R}^M. \quad (2.42)$$

Here the matrices $\mathbf{H} \in \mathbb{R}^{M \times L}$ and $\Psi \in \mathbb{R}^{M \times P}$ are known and the objective is to recover the unknown vectors \mathbf{b} and \mathbf{c} . Physically, the vector $\mathbf{H}\mathbf{b}$ represents unknown calibration gains that lie in a known subspace, specified by \mathbf{H} . Note that (2.42) is

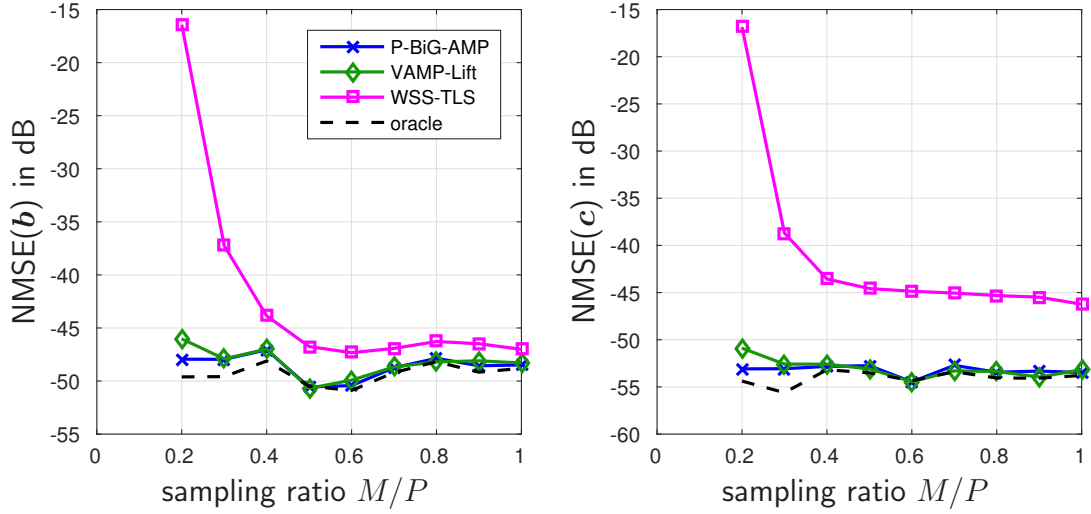


Figure 2.4: NMSE vs. M/P with i.i.d. $\mathcal{N}(0, 1)$ \mathbf{A} .

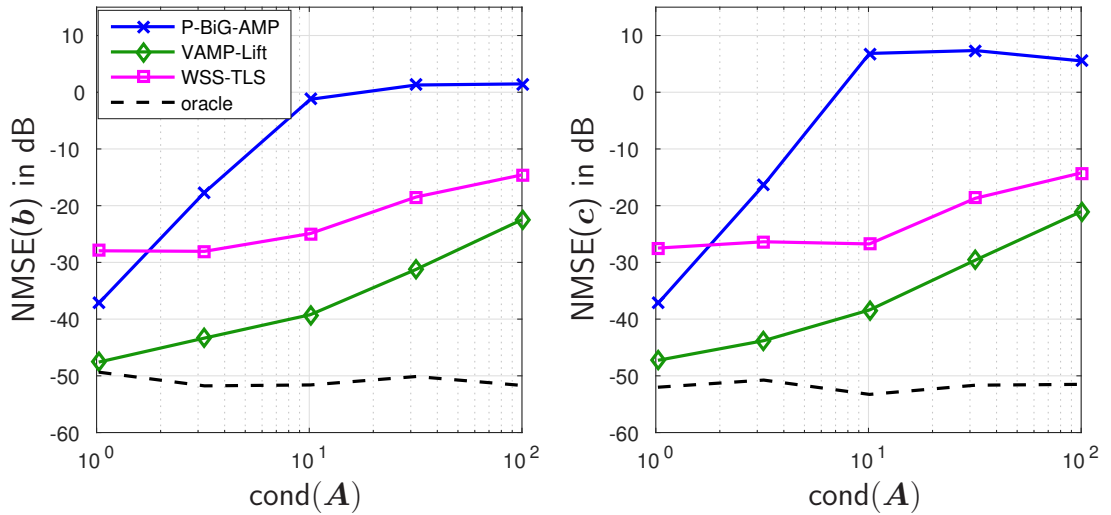


Figure 2.5: NMSE vs. $\text{cond}(\mathbf{A})$ at $M/P = 0.6$.

an instance of (2.39) with $\mathbf{A}_l = \text{Diag}(\mathbf{h}_l)\Psi$, where \mathbf{h}_l denotes the l th column of \mathbf{H} . Different from “CS with matrix uncertainty”, all elements in \mathbf{b} are now unknown,

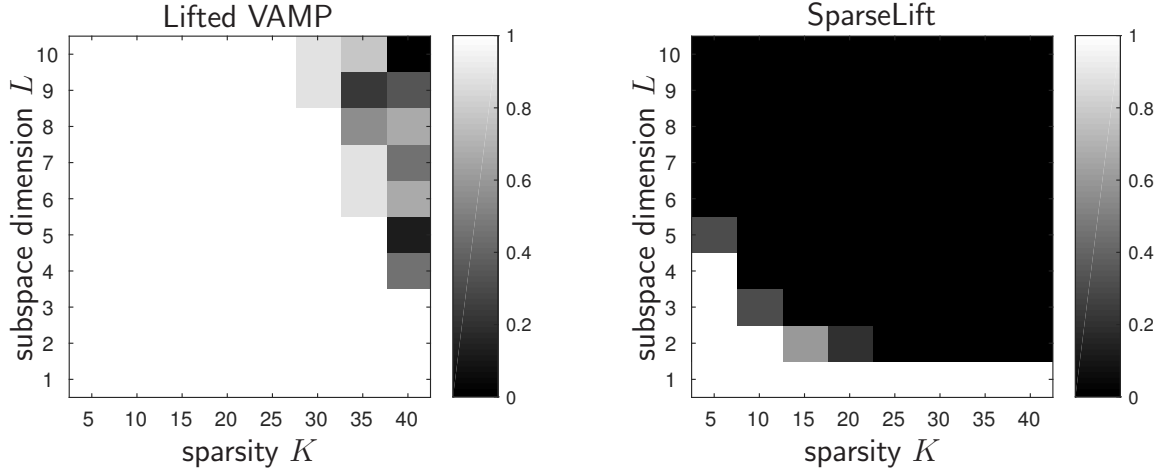


Figure 2.6: Success rate vs. sparsity K and subspace dimension L in self-calibration

and so WSS-TLS [105] cannot be applied. Instead, we compare lifted VAMP to the SparseLift approach from [51], which is based on convex relaxation and has provable guarantees.

For our experiment, we generated Ψ and $\mathbf{b} \in \mathbb{R}^L$ as i.i.d. $\mathcal{N}(0, 1)$; \mathbf{c} as K -sparse with $\mathcal{N}(0, 1)$ nonzero entries; \mathbf{H} as randomly chosen columns of a Hadamard matrix; and $\mathbf{w} = \mathbf{0}$. Figure 2.6 plots the success rate versus L and K , where “success” is defined as $\mathbb{E} \|\widehat{\mathbf{c}}\mathbf{b}^\top - \mathbf{c}\mathbf{b}^\top\|_F^2 / \mathbb{E} \|\mathbf{c}\mathbf{b}^\top\|_F^2 < -60$ dB. The figure shows that, relative to SparseLift, lifted VAMP gives successful recoveries for a wider range of L and K .

Chapter 3: Bilinear Recovery using Adaptive Vector Approximate Message Passing

3.1 Introduction

In this chapter, we are interested in problems that can be formulated as estimation of an unknown structured matrix \mathbf{Z} from noisy or incomplete measurements. The type of structure in \mathbf{Z} determines the specific subproblem to be solved.

For example, when \mathbf{Z} has a low-rank structure and only a subset of its entries are observed, the problem is known as *matrix completion* [21]. When $\mathbf{Z} = \mathbf{L} + \mathbf{S}$ for low-rank \mathbf{L} and sparse \mathbf{S} , the problem of estimating \mathbf{L} and \mathbf{S} is known as *robust principle components analysis* (RPCA) [20]. When $\mathbf{Z} = \mathbf{BC}$ with sparse \mathbf{C} , the problem of estimating \mathbf{B} and \mathbf{C} is known as *dictionary learning* [81]. When $\mathbf{Z} = \mathbf{BC}$ with nonnegative \mathbf{B} and \mathbf{C} , the problem is known as *nonnegative matrix factorization* (NMF) [47].

Sometimes \mathbf{Z} has a more complicated structure. For example, the problems of *self-calibration* and *blind (circular) deconvolution* [51] can be formulated using $\mathbf{Z} = \text{Diag}(\mathbf{H}\mathbf{b})\mathbf{\Psi}\mathbf{C}$, where \mathbf{H} and $\mathbf{\Psi}$ are known and \mathbf{b} and \mathbf{C} are to be estimated. The problem of *compressive sensing (CS) with matrix uncertainty* [105] can be formulated using $\mathbf{z} = \sum_i b_i \mathbf{A}_i \mathbf{c}$, where $\{\mathbf{A}_i\}$ are known and where \mathbf{b} and sparse \mathbf{c} are to be

estimated. The latter covers the problem of *joint channel-symbol estimation* [43], in which case b_i are the data symbols, \mathbf{c} contains (possibly sparse) channel coefficients, and the known $\{\mathbf{A}_i\}$ are determined by the modulation scheme. The more general problem of *matrix CS* [22, 100] results from

$$z_m = \text{tr}\{\mathbf{A}_m^\top(\mathbf{L} + \mathbf{S})\} \text{ for } m = 1, \dots, M, \quad (3.1)$$

where $\{\mathbf{A}_m\}$ are known and the goal is to estimate low-rank \mathbf{L} and sparse \mathbf{S} .

3.2 Prior Work

Many algorithms have been developed to solve the above problems. Some solve a convex relaxation of the original problem, while others attack non-convex formulations via alternating methods, greedy methods, variational methods, message-passing methods, and other techniques.

For matrix completion, well-known approaches include the nuclear-norm-based convex optimization method IALM [50], the non-convex successive over-relaxation approach LMAFit [101], the Grassmanian gradient-descent approach GROUSE [5], the greedy hard-thresholding approach Matrix-ALPS [46], and the variational-Bayes method VSBL [4]. For RPCA, there are also versions of IALM [50], LMAFit [101], and VSBL [4], as well as a robust cousin of GROUSE, called GRASTA [39]. For dictionary learning, there is the greedy K-SVD algorithm [2], the online SPAMS approach [54], and the ER-SpUD approach from [87]. A unified approach to matrix completion, RPCA, and dictionary learning was proposed in [42, 67, 68] using an extension of the approximate message-passing (AMP) methodology from [32, 71]. The resulting “bilinear generalized AMP” (BiGAMP) algorithm was compared to the aforementioned methods in [68] and found (empirically) to be competitive, if not

superior, in phase transition and runtime. A related approach known as LowRAMP was proposed [55] and analyzed in [49, 59].

For self-calibration and blind deconvolution, well-known approaches include the convex relaxations from [3, 13, 51] and the alternating method from [40]. For CS with matrix uncertainty, there is the award-winning non-convex method [105]. For matrix CS, the well-known papers [1, 22, 102] proposed convex approaches and [46, 100] proposed greedy approaches. See the recent overview [28] for many other works. An AMP-based approach to self-calibration, blind deconvolution, CS with matrix uncertainty, and matrix CS was proposed in [66] and analyzed in [85]. This “parametric BiGAMP” (PBiGAMP) was compared to the above works in [66] and found to yield improved empirical phase transitions.

More recently, AMP methods for bilinear inference were proposed using the “lifting” approach (see, e.g., [3, 23, 28, 51] for seminal papers on lifting). We used the lifting method in chapter 2 to solve some bilinear problems by framing it as a linear inverse problem. We illustrate the main idea behind lifting here. Suppose the measurement vector $\mathbf{y} \in \mathbb{R}^M$ takes the form

$$\mathbf{y} = \sum_{i=1}^Q \sum_{j=1}^N b_i \mathbf{a}_{i,j} c_j + \mathbf{w}, \quad (3.2)$$

where $\mathbf{a}_{i,j} \in \mathbb{R}^M$ is known for all i, j and the goal is to recover $\mathbf{b} = [b_1, \dots, b_Q]^\top$ and $\mathbf{c} = [c_1, \dots, c_N]^\top$ in the presence of white noise \mathbf{w} . Rewriting the measurements as

$$\mathbf{y} = \sum_{i=1}^Q b_i \underbrace{[\mathbf{a}_{i,1}, \dots, \mathbf{a}_{i,N}]}_{\triangleq \mathbf{A}_i} \mathbf{c} + \mathbf{w} \quad (3.3)$$

$$= \underbrace{[\mathbf{A}_1 \cdots \mathbf{A}_Q]}_{\triangleq \mathbf{A}} \begin{bmatrix} b_1 \mathbf{c} \\ \vdots \\ b_Q \mathbf{c} \end{bmatrix} + \mathbf{w} \quad (3.4)$$

$$= \mathbf{A} \mathbf{x} + \mathbf{w} \quad \text{for } \mathbf{x} = \mathbf{b} \otimes \mathbf{c} = \text{vec}(\mathbf{c} \mathbf{b}^\top), \quad (3.5)$$

we see that the noisy bilinear recovery problem (3.2) can be rewritten as the noisy linear recovery problem (3.5) with a rank-one structure on (the matrix form of) \mathbf{x} . Thus, if this low-rank signal structure can be exploited by a linear inference algorithm, then bilinear inference can be accomplished. This is precisely what was proposed in [80], building on the non-separable-denoising version of the AMP algorithm from [56]. A rigorous analysis of “lifted AMP” was presented in [12].

The trouble with AMP is that its behavior is understood only in the case of large [83] or infinitely large, i.i.d. (sub) Gaussian \mathbf{A} [8,9]. Even small deviations from this scenario (e.g., mildly ill-conditioned and/or non-zero-mean \mathbf{A}) can cause AMP to diverge [18, 74, 96]. The Vector AMP (VAMP) algorithm discussed in chapter 2 addresses this issue. In [75], it was first established that, if \mathbf{A} is an infinitely large right-rotationally invariant² random matrix and the denoising function used by VAMP is separable and Lipschitz, then VAMP’s performance can be exactly predicted by a scalar state-evolution that also provides testable conditions for optimality. Since the class of right-rotationally invariant matrices is much larger than the class of

²If \mathbf{A} is right-rotationally invariant then its singular value decomposition $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$ has Haar distributed \mathbf{V} , i.e., \mathbf{V} is uniformly distributed over the group of orthogonal matrices.

i.i.d. Gaussian matrices, VAMP is much more robust than AMP with regards to the construction of \mathbf{A} . For example, VAMP has no problem with ill-conditioned or mean-shifted matrices [75].

In chapter 2 Sec. 2.4, we performed a rigorous analysis of VAMP under *non*-separable Lipschitz denoisers, showing that—here too—VAMP’s behavior is exactly predicted by a scalar state-evolution when \mathbf{A} is infinitely large and right-rotationally invariant. Moreover, we also demonstrated the success of lifted VAMP on bilinear problems such as self-calibration and CS with matrix uncertainty. In addition, we gave evidence that, like AMP, the PBiGAMP algorithm is sensitive to deviations from the i.i.d. assumptions used in its derivation [66] and analysis [85]. For this reason, lifted VAMP significantly outperformed PBiGAMP in some cases.

In our numerical experiments (Sec. 2.5), we showed that lifted VAMP has good performance and a rigorous analyses under infinitely large right-rotationally invariant random \mathbf{A} . However, it suffers from computational issues brought on by the lifting itself: The $N + Q$ unknowns $[\mathbf{b}, \mathbf{c}]$ in the bilinear problem (3.2) manifest as NQ unknowns \mathbf{x} after lifting to (3.5). This is a serious problem when N and Q are both large. As a concrete example, consider the application of lifting to (square) dictionary learning, where the goal is to recover $\mathbf{B} \in \mathbb{R}^{N \times N}$ and sparse $\mathbf{C} \in \mathbb{R}^{N \times L}$ from noisy measurements $\mathbf{Y} = \mathbf{BC} + \mathbf{W}$. This bilinear relationship can be lifted via

$$\mathbf{Y} = \sum_{ij} b_{i,j} \mathbf{A}_{i,j} \mathbf{C} + \mathbf{W} \tag{3.6}$$

$$= \underbrace{[\mathbf{A}_{1,1} \cdots \mathbf{A}_{N,N}]}_{\triangleq \mathbf{A} \in \mathbb{R}^{N \times N^3}} \underbrace{(\mathbf{b} \otimes \mathbf{C})}_{\triangleq \mathbf{X} \in \mathbb{R}^{N^3 \times L}} + \mathbf{W}, \tag{3.7}$$

where $\mathbf{A}_{i,j} \in \mathbb{R}^{N \times N}$ is constructed with a 1 in the (i, j) th position and zeros elsewhere, and where $\mathbf{b} = [b_{1,1}, \dots, b_{N,N}]^T \in \mathbb{R}^{N^2}$. Even at the relatively small patch size of 8×8

(i.e., $N = 64$), the matrix \mathbf{A} has dimension $64 \times 262\,144$, and the unknown matrix \mathbf{X} has dimension $262\,144 \times L$. The rule-of-thumb $L = 5N \ln N$ [87] then gives $L = 1331$, in which case \mathbf{X} contains 3.5×10^8 entries, which leads to difficulties with computation and memory.

3.3 Contributions

In this chapter, we present a novel VAMP-based approach to bilinear recovery. With the aim of computational efficiency, we avoid lifting and instead build on the recently proposed Adaptive VAMP framework from [37]. However, different from [37], which focused on noisy *linear* recovery, we focus on noisy *bilinear* recovery.

In particular, we focus on recovering the unknown parameters $\{b_i\}_{i=1}^Q$ and the unknown random matrix $\mathbf{C} \in \mathbb{R}^{N \times L}$ from noisy measurements $\mathbf{Y} \in \mathbb{R}^{M \times L}$ of the form

$$\mathbf{Y} = \sum_{i=1}^Q b_i \mathbf{A}_i \mathbf{C} + \mathbf{W}, \quad (3.8)$$

where $\{\mathbf{A}_i\}$ are known and \mathbf{W} is additive white Gaussian noise (AWGN). Note that (3.8) is a multiple-measurement vector (MMV) extension of (3.3), and that it covers all of the motivating problems discussed in Sec. 3.1. For example, in self-calibration, where we estimate \mathbf{b} and \mathbf{C} from $\mathbf{Y} = \text{Diag}(\mathbf{H}\mathbf{b})\mathbf{\Psi}\mathbf{C} + \mathbf{W}$, we can set $\mathbf{A}_i = \text{Diag}(\mathbf{h}_i)\mathbf{\Psi}$, where \mathbf{h}_i is the i th column of \mathbf{H} . Or, in dictionary learning, where we estimate \mathbf{B} and \mathbf{C} from $\mathbf{Y} = \mathbf{B}\mathbf{C} + \mathbf{W}$, we can write $\mathbf{B} = \sum_{i=1}^{MN} b_i \mathbf{A}_i$ for $\mathbf{A}_i = \mathbf{e}_{\langle i-1 \rangle_M} \mathbf{e}_{\lfloor (i-1)/M \rfloor}^\top$, where $\langle i \rangle_M$ denotes i -modulo- M , $\lfloor \cdot \rfloor$ denotes floor, and $\{\mathbf{e}_i\}$ is the standard basis.

When deriving³ the proposed method, we treat $\{b_i\}$ as deterministic unknowns and the entries of \mathbf{C} as random variables. The prior distribution on \mathbf{C} is assumed to be known up to some (possibly) unknown hyperparameters, which are learned jointly with $\{b_i\}$ and \mathbf{C} . Also, \mathbf{W} is treated as additive white Gaussian noise (AWGN) with an unknown variance that is also learned. More details are provided in the sequel.

We show (empirically) that the proposed Bilinear Adaptive VAMP (BAd-VAMP) method performs as well as the EM-PBiGAMP algorithm from [66], with regard to accuracy and computational complexity, when the underlying matrices are i.i.d., as assumed for the derivation of PBiGAMP. However, we will show that BAd-VAMP outperforms EM-PBiGAMP when the underlying matrices become ill-conditioned. In the ill-conditioned case, we show that BAd-VAMP performs as well as, and sometimes significantly better than, lifted VAMP. However, BAd-VAMP is much more computationally efficient due to its avoidance of lifting. In this sense, the proposed BAd-VAMP is shown to be accurate, robust, and computationally efficient.

3.4 Problem Formulation

In an effort to make our algorithmic development more consistent with our discussion on VAMP in chapter 2, we now make the following notational changes relative to (3.8):

- (a) We will use the notation $\mathbf{A}(\mathbf{b}) \triangleq \sum_i b_i \mathbf{A}_i$ to be concise.
- (b) The quantities b_i and \mathbf{C} in (3.8) will be changed to $\theta_{A,i}$ and \mathbf{X} respectively.

³Although the derivation treats the entries of \mathbf{C} as random variables and the associated denoiser as Bayesian, the final algorithm is more general in that it only requires the denoiser to be Lipschitz.

We can thus state our problem of interest as follows: estimate the matrix $\mathbf{X} \in \mathbb{R}^{N \times L}$ and learn the parameters $\Theta \triangleq \{\theta_A, \theta_x, \gamma_w\}$ in the statistical model

$$\mathbf{Y} = \mathbf{A}(\theta_A)\mathbf{X} + \mathbf{W} \quad (3.9a)$$

$$\mathbf{x}_l \sim p_{\mathbf{x}}(\cdot; \theta_x) \quad \forall l \in \{1, \dots, L\} \quad (3.9b)$$

$$w_{ml} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \gamma_w^{-1}), \quad (3.9c)$$

where $\mathbf{A}(\cdot)$ is a known matrix-valued linear function, and $p_{\mathbf{x}}(\cdot; \theta_x)$ is a prior density on \mathbf{X} parameterized by the unknown vector θ_x . Here, γ_w is the noise precision, i.e., the inverse noise variance.

More precisely, we aim to compute the maximum-likelihood (ML) estimate of Θ and, under that estimate, compute the minimum mean-squared error (MMSE) estimate of \mathbf{X} , i.e.,

$$\hat{\Theta}_{\text{ML}} = \arg \max_{\Theta} p_{\mathbf{Y}}(\mathbf{Y}; \Theta) \quad (3.10)$$

$$\hat{\mathbf{X}}_{\text{MMSE}} = \mathbb{E}[\mathbf{X} | \mathbf{Y}; \hat{\Theta}_{\text{ML}}]. \quad (3.11)$$

In (3.10), $p_{\mathbf{Y}}(\mathbf{Y}; \Theta)$ is the likelihood function of Θ , which can be written as

$$p_{\mathbf{Y}}(\mathbf{Y}; \Theta) = \int p_{\mathbf{X}}(\mathbf{X}; \Theta) p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\mathbf{X}; \Theta) d\mathbf{X}. \quad (3.12)$$

In (3.11), the expectation is taken over the posterior density

$$p_{\mathbf{X}|\mathbf{Y}}(\mathbf{X}|\mathbf{Y}; \hat{\Theta}_{\text{ML}}) = \frac{p_{\mathbf{X}}(\mathbf{X}; \hat{\Theta}_{\text{ML}}) p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\mathbf{X}; \hat{\Theta}_{\text{ML}})}{p_{\mathbf{Y}}(\mathbf{Y}; \hat{\Theta}_{\text{ML}})}. \quad (3.13)$$

The statistical model (3.9) implies that⁴

$$p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\mathbf{X}; \Theta) = \prod_{l=1}^L p_{y_l|\mathbf{x}}(y_l|\mathbf{x}_l; \Theta) \quad (3.14)$$

⁴In (3.14)-(3.16), to promote notational simplicity, the left side of the equation is written using Θ even though the right side depends on a subset of Θ .

where each factor becomes

$$p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}; \Theta) = \mathcal{N}(\mathbf{y}; \mathbf{A}(\boldsymbol{\theta}_A)\mathbf{x}, \mathbf{I}/\gamma_w). \quad (3.15)$$

Since the columns $\{\mathbf{x}_l\}_{l=1}^L$ are independently drawn from density $p_{\mathbf{x}}(\cdot; \boldsymbol{\theta}_x)$, we have

$$p_{\mathbf{X}}(\mathbf{X}; \Theta) = \prod_{l=1}^L p_{\mathbf{x}}(\mathbf{x}_l; \boldsymbol{\theta}_x) \quad (3.16)$$

for some density $p_{\mathbf{x}}(\cdot; \boldsymbol{\theta}_x)$ parameterized by $\boldsymbol{\theta}_x$. In this case, the posterior density $p_{\mathbf{X}|\mathbf{Y}}$ decouples across the columns of \mathbf{X} and \mathbf{Y} as

$$p_{\mathbf{X}|\mathbf{Y}}(\mathbf{X}|\mathbf{Y}; \Theta) \propto \prod_{l=1}^L p_{\mathbf{x}}(\mathbf{x}_l; \Theta) p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}_l|\mathbf{x}_l; \Theta). \quad (3.17)$$

3.5 Expectation Maximization

We can use maximum likelihood (ML) estimation to learn the unknown parameters Θ in (3.9). From (3.10) and (3.12), we have

$$\hat{\Theta}_{\text{ML}} = \arg \min_{\Theta} - \ln \int p_{\mathbf{X}}(\mathbf{X}; \Theta) p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\mathbf{X}; \Theta) d\mathbf{X}, \quad (3.18)$$

but (3.18) is impractical to optimize directly due to the high dimensional integral.

Expectation-maximization (EM) [29] is a well known iterative approach to ML that alternates between the following two steps:

- (a) Minimizing an upper-bound of the negative log-likelihood.
- (b) Tightening the upper-bound.

The EM algorithm is usually written as

$$Q(\Theta; \hat{\Theta}^t) \triangleq -\mathbb{E} [\ln p_{\mathbf{X},\mathbf{Y}}(\mathbf{X}, \mathbf{Y}; \Theta) | \mathbf{Y}; \hat{\Theta}^t] \quad (3.19a)$$

$$\hat{\Theta}^{t+1} = \arg \min_{\Theta} Q(\Theta; \hat{\Theta}^t). \quad (3.19b)$$

Letting $q^t = p_{\mathbf{X}|\mathbf{Y}}(\cdot|\mathbf{Y}; \widehat{\Theta}^t)$, we can write

$$\begin{aligned} Q(\Theta; \widehat{\Theta}^t) &= -\mathbb{E} [\ln p_{\mathbf{X}}(\mathbf{X}; \Theta) | \mathbf{Y}; \widehat{\Theta}^t] \\ &\quad - \mathbb{E} [\ln p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\mathbf{X}; \Theta) | \mathbf{Y}; \widehat{\Theta}^t] \end{aligned} \quad (3.20a)$$

$$\begin{aligned} &= -\mathbb{E} [\ln p_{\mathbf{X}}(\mathbf{X}; \Theta) | q^t] \\ &\quad - \mathbb{E} [\ln p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\mathbf{X}; \Theta) | q^t] \end{aligned} \quad (3.20b)$$

$$= J(q^t, q^t, q^t; \Theta) + \text{const.} \quad (3.20c)$$

where J is the Gibbs free energy from (2.13) where there were no unknown parameters.

In the presence of unknown parameters Θ , it takes the form

$$J(q_1, q_2, q_3; \Theta) \triangleq D_{\text{KL}}(q_1 \| p_{\mathbf{X}}(\cdot, \Theta)) + D_{\text{KL}}(q_2 \| p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\cdot; \Theta)) + H(q_3). \quad (3.21)$$

Thus, (3.19) can also be written as the following two steps [61]:

$$q^t = p_{\mathbf{X}|\mathbf{Y}}(\cdot|\mathbf{Y}; \widehat{\Theta}^t) \quad (3.22a)$$

$$\widehat{\Theta}^{t+1} = \arg \min_{\Theta} J(q^t, q^t, q^t; \Theta). \quad (3.22b)$$

EM is an type of Majorization-Minimization algorithm as it minimizes an upper bound to the actual cost function that we want to minimize. We can show that $J(q^t, q^t, q^t; \Theta)$ is an upper bound on $-\ln p_{\mathbf{Y}}(\mathbf{Y}; \Theta)$ for *any* q^t since

$$J(q^t, q^t, q^t; \Theta) = -\ln p_{\mathbf{Y}}(\mathbf{Y}; \Theta) + D_{\text{KL}}(q^t \| p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\cdot; \Theta)) \quad (3.23)$$

$$\stackrel{(a)}{\geq} -\ln p_{\mathbf{Y}}(\mathbf{Y}; \Theta), \quad (3.24)$$

where (a) follows from $D_{\text{KL}} \geq 0$ by construction. Thus, while the specific choice of q^t in (3.22a) yields a tight upper bound in that

$$J(q^t, q^t, q^t; \widehat{\Theta}^t) = -\ln p_{\mathbf{Y}}(\mathbf{Y}; \widehat{\Theta}^t), \quad (3.25)$$

Algorithm 3 Bilinear EM-VAMP

```

1: initialize:
    $\forall l : \mathbf{r}_{1,l}^0, \gamma_{1,l}^0, \boldsymbol{\theta}_x^0, \boldsymbol{\theta}_A^0, \gamma_w^0$ 
2: for  $t = 0, \dots, T_{\max}$  do
3:    $\forall l : \mathbf{x}_{1,l}^t = \mathbf{g}_1(\mathbf{r}_{1,l}^t, \gamma_{1,l}^t; \boldsymbol{\theta}_x^t)$ 
4:    $\forall l : 1/\eta_{1,l}^t = \langle \mathbf{g}'_1(\mathbf{r}_{1,l}^t, \gamma_{1,l}^t; \boldsymbol{\theta}_x^t) \rangle / \gamma_{1,l}^t$ 
5:    $q_1^t(\mathbf{X}) \propto \prod_{l=1}^L p_{\mathbf{x}}(\mathbf{x}_l; \boldsymbol{\theta}_x^t) e^{-\frac{1}{2}\gamma_{1,l}^t \|\mathbf{x}_l - \mathbf{r}_{1,l}^t\|^2}$ 
6:    $\boldsymbol{\theta}_x^{t+1} = \arg \max_{\boldsymbol{\theta}_x} \mathbb{E}[\ln p_{\mathbf{x}}(\mathbf{X}; \boldsymbol{\theta}_x) | q_1^t]$ 
7:    $\forall l : \gamma_{2,l}^t = \eta_{1,l}^t - \gamma_{1,l}^t$ 
8:    $\forall l : \mathbf{r}_{2,l}^t = (\eta_{1,l}^t \mathbf{x}_{1,l}^t - \gamma_{1,l}^t \mathbf{r}_{1,l}^t) / \gamma_{2,l}^t$ 
9:    $\forall l : \mathbf{x}_{2,l}^t = \mathbf{g}_{2,l}(\mathbf{r}_{2,l}^t, \gamma_{2,l}^t; \boldsymbol{\theta}_A^t, \gamma_w^t)$ 
10:   $\forall l : 1/\eta_{2,l}^t = \langle \mathbf{g}'_{2,l}(\mathbf{r}_{2,l}^t, \gamma_{2,l}^t; \boldsymbol{\theta}_A^t, \gamma_w^t) \rangle / \gamma_{2,l}^t$ 
11:   $q_2^t(\mathbf{X}) \propto \prod_{l=1}^L p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}_l | \mathbf{x}_l; \boldsymbol{\theta}_A^t, \gamma_w^t) e^{-\frac{1}{2}\gamma_{2,l}^t \|\mathbf{x}_l - \mathbf{r}_{2,l}^t\|^2}$ 
12:   $\boldsymbol{\theta}_A^{t+1} = \arg \max_{\boldsymbol{\theta}_A} \mathbb{E}[\ln p_{\mathbf{y}|\mathbf{x}}(\mathbf{Y} | \mathbf{X}; \boldsymbol{\theta}_A, \gamma_w^t) | \mathbf{Y}, q_2^t]$ 
13:   $\gamma_w^{t+1} = \arg \max_{\gamma_w} \mathbb{E}[\ln p_{\mathbf{y}|\mathbf{x}}(\mathbf{Y} | \mathbf{X}; \boldsymbol{\theta}_A^{t+1}, \gamma_w) | \mathbf{Y}, q_2^t]$ 
14:   $\forall l : \gamma_{1,l}^{t+1} = \eta_{2,l}^t - \gamma_{2,l}^t$ 
15:   $\forall l : \mathbf{r}_{1,l}^{t+1} = (\eta_{2,l}^t \mathbf{x}_{2,l}^t - \gamma_{2,l}^t \mathbf{r}_{2,l}^t) / \gamma_{1,l}^{t+1}$ 
16: end for

```

other choices of bounding q^t can also be used in EM [61]. EM is used in many parameter estimation problems in signal processing most notable being training hidden Markov Models in speech [69]. In the next section, we will use the variational interpretation of EM to develop an algorithm to solve the bilinear recovery problem in (3.9).

3.6 Bilinear EM-VAMP

From the descriptions of VAMP in chapter 2 and EM algorithm in Sec. 3.5, we see that both the algorithms minimize the same Gibbs free energy cost $J(q_1, q_2, q_3; \boldsymbol{\Theta})$ from (3.21), but w.r.t. different variables. VAMP minimizes J w.r.t. the beliefs $\{q_1, q_2, q_3\}$ under the first and second moment-matching constraints for a given value of $\boldsymbol{\Theta}$. Meanwhile, EM minimizes J w.r.t. the parameters $\boldsymbol{\Theta}$ for a given beliefs

$\{q_1, q_2, q_3\}$. As a result, the two approaches can be straightforwardly merged for *joint* estimation of $\{q_1, q_2, q_3\}$ and Θ . In doing so, the goal is to solve the following constrained optimization problem

$$\arg \min_{\Theta, q_1} \min_{q_2} \max_{q_3} J(q_1, q_2, q_3; \Theta) \quad (3.26a)$$

$$\text{s.t.} \quad \mathbb{E}[\mathbf{x}|q_1] = \mathbb{E}[\mathbf{x}|q_2] = \mathbb{E}[\mathbf{x}|q_3] \quad (3.26b)$$

$$\text{tr}\{\text{Cov}[\mathbf{x}|q_1]\} = \text{tr}\{\text{Cov}[\mathbf{x}|q_2]\} = \text{tr}\{\text{Cov}[\mathbf{x}|q_3]\}, \quad (3.26c)$$

and the proposed methodology is to “interleave” the VAMP and EM algorithms, as specified in Alg. 3. In lines 3-4 of Alg. 3, the estimation function \mathbf{g}_1 is similar to that in VAMP (2.17) but with parameters θ_x ,

$$\mathbf{g}_1(\mathbf{r}_{1,l}, \gamma_{1,l}; \theta_x) \triangleq \frac{\int \mathbf{x} p_{\mathbf{x}}(\mathbf{x}; \theta_x) \mathcal{N}(\mathbf{x}; \mathbf{r}_{1,l}, \mathbf{I}/\gamma_{1,l}) d\mathbf{x}}{\int p_{\mathbf{x}}(\mathbf{x}; \theta_x) \mathcal{N}(\mathbf{x}; \mathbf{r}_{1,l}, \mathbf{I}/\gamma_{1,l}) d\mathbf{x}}. \quad (3.27)$$

And in lines 9-10, the estimation function $\mathbf{g}_{2,l}$ is the MMV version of the linear estimator \mathbf{g}_2 from (2.18) and it is defined as

$$\mathbf{g}_{2,l}(\mathbf{r}_{2,l}, \gamma_{2,l}; \theta_A, \gamma_w) \triangleq \frac{\int \mathbf{x} p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}_l|\mathbf{x}; \theta_A, \gamma_w) \mathcal{N}(\mathbf{x}; \mathbf{r}_{2,l}, \mathbf{I}/\gamma_{2,l}) d\mathbf{x}}{\int p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}_l|\mathbf{x}; \theta_A, \gamma_w) \mathcal{N}(\mathbf{x}; \mathbf{r}_{2,l}, \mathbf{I}/\gamma_{2,l}) d\mathbf{x}}. \quad (3.28)$$

It is the LMMSE estimator applied on each column \mathbf{x}_l in the model

$$\mathbf{y}_l = \mathbf{A}(\theta_A)\mathbf{x}_l + \mathbf{w}_l \quad (3.29a)$$

$$\mathbf{x}_l \sim \mathcal{N}(\mathbf{r}_{2,l}, \mathbf{I}/\gamma_{2,l}) \quad (3.29b)$$

$$\mathbf{w}_l \sim \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w). \quad (3.29c)$$

The other lines in Alg. 3 will be detailed in sections from 3.7.1 to 3.7.3.

3.7 Bilinear Adaptive VAMP

Assume at every iteration t , the following conditions are satisfied in BAdVAMP

- i) The matrix $\mathbf{A}(\boldsymbol{\theta}_A^t)$ is infinitely large and right-rotationally invariant.
- ii) The estimation functions \mathbf{g}_1 and $\{\mathbf{g}_{2,l}\}_{l=1}^L$ are “matched” (i.e., MMSE) for the statistical model generating (\mathbf{X}, \mathbf{Y}) .

Then, the VAMP state-evolution from (2.34) shows that at every iteration t , the VAMP quantities $\{(\mathbf{r}_{1,l}^t, \gamma_{1,l}^t, \mathbf{r}_{2,l}^t, \gamma_{2,l}^t)\}_{l=1}^L$ obey the statistical model

$$\mathbf{r}_{1,l}^t = \mathbf{x}_l + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_{1,l}^t) \quad \forall l \quad (3.30a)$$

$$\mathbf{x}_l = \mathbf{r}_{2,l}^t + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_{2,l}^t) \quad \forall l, \quad (3.30b)$$

where \mathbf{x}_l is the l th column of the true signal realization \mathbf{X} that we aim to recover. That is, $\mathbf{r}_{1,l}^t$ is an AWGN-corrupted version of the true signal \mathbf{x}_l with *known* AWGN precision $\gamma_{1,l}^t$, and the true signal \mathbf{x}_l is an AWGN-corrupted version of $\mathbf{r}_{2,l}^t$ with *known* AWGN precision $\gamma_{2,l}^t$. In the context of EM-VAMP under (3.9), this “matched” condition requires that $\boldsymbol{\theta}_A^t$, $\boldsymbol{\theta}_x^t$, and γ_w^t are all perfect estimates. When $\boldsymbol{\theta}_A^t$, $\boldsymbol{\theta}_x^t$, or γ_w^t are *not* perfect, so that \mathbf{g}_1 and $\mathbf{g}_{2,l}$ are mismatched, the VAMP state-evolution shows that $\mathbf{r}_{1,l}^t$ is still an AWGN corrupted version of \mathbf{x}_l , but with an AWGN precision *different* than $\gamma_{1,l}^t$.

In the mismatched case, the impact on EM-VAMP is the following. While the algorithm is trying to learn $\Theta = \{\boldsymbol{\theta}_A, \boldsymbol{\theta}_x, \gamma_w\}$, the value of $\gamma_{i,l}^t$ does *not* correctly characterize the noise precision in $\mathbf{r}_{i,l}^t$. As a result, the beliefs q_1^t and q_2^t in lines 5 and 11 of Alg. 3 become mismatched, which compromises the EM updates of Θ^t .

To remedy this situation, it was proposed in [45] (in the context of EM-GAMP [97]) to explicitly estimate the precision of the AWGN corruption on $\mathbf{r}_{1,l}^t$ and $\mathbf{r}_{2,l}^t$ and use it in place of the AMP-supplied estimates $\gamma_{1,l}^t$ and $\gamma_{2,l}^t$. This approach was coined “Adaptive” GAMP in [45] and later extended to (linear) Adaptive VAMP in [37].

For Bilinear Adaptive VAMP, the first goal is to replace the estimation of $\boldsymbol{\theta}_x$ in line 6 of Alg. 3 with the joint ML estimation

$$(\boldsymbol{\theta}_x^t, \boldsymbol{\gamma}_1^t) = \arg \max_{\boldsymbol{\theta}_x, \boldsymbol{\gamma}_1} p(\mathbf{R}_1^t; \boldsymbol{\gamma}_1, \boldsymbol{\theta}_x) \quad (3.31)$$

under the following statistical model

$$\mathbf{r}_{1,l}^t = \mathbf{x}_l + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_{1,l}) \quad \forall l \quad (3.32)$$

$$\mathbf{x}_l \sim p_{\mathbf{x}}(\cdot; \boldsymbol{\theta}_x) \quad \forall l, \quad (3.33)$$

with independence across $l = 1, \dots, L$. For this subproblem, we propose to use (inner) EM iterations indexed by τ , i.e.,

$$(\boldsymbol{\theta}_x^{\tau+1}, \boldsymbol{\gamma}_1^{\tau+1}) = \arg \max_{\boldsymbol{\theta}_x, \boldsymbol{\gamma}_1} \mathbb{E} \left[\ln p(\mathbf{X}, \mathbf{R}_1^t; \boldsymbol{\gamma}_1, \boldsymbol{\theta}_x) \mid \mathbf{R}_1^t; \boldsymbol{\gamma}_1^\tau, \boldsymbol{\theta}_x^\tau \right] \quad (3.34)$$

$$\begin{aligned} &= \arg \max_{\boldsymbol{\theta}_x, \boldsymbol{\gamma}_1} \left\{ \mathbb{E} \left[\ln p(\mathbf{X}; \boldsymbol{\theta}_x) \mid \mathbf{R}_1^t; \boldsymbol{\gamma}_1^\tau, \boldsymbol{\theta}_x^\tau \right] \right. \\ &\quad \left. + \mathbb{E} \left[\ln p(\mathbf{R}_1^t \mid \mathbf{X}; \boldsymbol{\gamma}_1) \mid \mathbf{R}_1^t; \boldsymbol{\gamma}_1^\tau, \boldsymbol{\theta}_x^\tau \right] \right\}. \end{aligned} \quad (3.35)$$

From (3.32), we can write the following

$$p(\mathbf{R}_1^t \mid \mathbf{X}; \boldsymbol{\gamma}_1) = \prod_{l=1}^L p(\mathbf{r}_{1,l}^t \mid \mathbf{x}_l; \gamma_{1,l}) \quad (3.36)$$

$$= \prod_{l=1}^L \mathcal{N}(\mathbf{r}_{1,l}^t; \mathbf{x}_l, \mathbf{I}/\gamma_{1,l}). \quad (3.37)$$

The optimization problem in (3.34) decouples into

$$\boldsymbol{\theta}_x^{\tau+1} = \arg \max_{\boldsymbol{\theta}_x} \mathbb{E} \left[\ln p(\mathbf{X}; \boldsymbol{\theta}_x) \mid \mathbf{R}_1^t; \boldsymbol{\gamma}_1^\tau, \boldsymbol{\theta}_x^\tau \right] \quad (3.38)$$

and

$$\boldsymbol{\gamma}_1^{\tau+1} = \arg \max_{\boldsymbol{\gamma}_1} \mathbb{E} [\ln p(\mathbf{R}_1^t | \mathbf{X}; \boldsymbol{\gamma}_1) | \mathbf{R}_1^t; \boldsymbol{\gamma}_1^\tau, \boldsymbol{\theta}_x^\tau] \quad (3.39)$$

$$= \arg \max_{\boldsymbol{\gamma}_1} \sum_{l=1}^L \mathbb{E} [\ln p(\mathbf{r}_{1,l}^t | \mathbf{x}_l; \boldsymbol{\gamma}_{1,l}) | \mathbf{r}_{1,l}^t; \boldsymbol{\gamma}_{1,l}^\tau, \boldsymbol{\theta}_x^\tau], \quad (3.40)$$

where the latter optimization decouples further into

$$\boldsymbol{\gamma}_{1,l}^{\tau+1} = \arg \max_{\boldsymbol{\gamma}_{1,l}} \left\{ \frac{N}{2} \ln \boldsymbol{\gamma}_{1,l} - \frac{\boldsymbol{\gamma}_{1,l}}{2} \mathbb{E} [\|\mathbf{x}_l - \mathbf{r}_{1,l}^t\|_2^2 | \mathbf{r}_{1,l}^t; \boldsymbol{\gamma}_{1,l}^\tau, \boldsymbol{\theta}_x^\tau] \right\} \quad (3.41)$$

$$= N \left\{ \mathbb{E} [\|\mathbf{x}_l - \mathbf{r}_{1,l}^t\|_2^2 | \mathbf{r}_{1,l}^t; \boldsymbol{\gamma}_{1,l}^\tau, \boldsymbol{\theta}_x^\tau] \right\}^{-1} \quad (3.42)$$

$$= \left\{ \frac{1}{N} \sum_{n=1}^N \mathbb{E} [(x_{nl} - r_{1,nl}^t)^2 | \mathbf{r}_{1,l}^t; \boldsymbol{\gamma}_{1,l}^\tau, \boldsymbol{\theta}_x^\tau] \right\}^{-1} \quad (3.43)$$

$$= \left\{ \frac{1}{N} \|\boldsymbol{\mathbf{x}}_{1,l}^\tau - \mathbf{r}_{1,l}^t\|^2 + \frac{1}{\eta_{1,l}^\tau} \right\}^{-1}, \quad (3.44)$$

for $l = 1, \dots, L$ and

$$\boldsymbol{\mathbf{x}}_{1,l}^\tau \triangleq \mathbb{E} [\mathbf{x}_l | \mathbf{r}_{1,l}^t; \boldsymbol{\gamma}_{1,l}^\tau, \boldsymbol{\theta}_x^\tau] \quad (3.45)$$

$$= \mathbf{g}_1(\mathbf{r}_{1,l}^t, \boldsymbol{\gamma}_{1,l}^\tau; \boldsymbol{\theta}_x^\tau) \quad (3.46)$$

$$1/\eta_{1,l}^\tau \triangleq \text{tr} \{ \text{Cov} [\mathbf{x}_l | \mathbf{r}_{1,l}^t; \boldsymbol{\gamma}_{1,l}^\tau, \boldsymbol{\theta}_x^\tau] \} / N \quad (3.47)$$

$$= \langle \mathbf{g}'_1(\mathbf{r}_{1,l}^t, \boldsymbol{\gamma}_{1,l}^\tau; \boldsymbol{\theta}_x^\tau) \rangle / \boldsymbol{\gamma}_{1,l}^\tau. \quad (3.48)$$

Above, we detailed the re-estimation of $\boldsymbol{\gamma}_1^t$. A similar procedure can be used for re-estimation of the precisions $\boldsymbol{\gamma}_2^t$. The resulting Bilinear Adaptive VAMP (BAD-VAMP) is summarized in Alg. 4 using $\tau_{1,\max}$ EM iterations for the first inner loop and $\tau_{2,\max}$ EM iterations for the second inner loop. To avoid the complications of a dual-index notation (i.e., t and τ), we use only the single index t in Alg. 4 and over-write the quantities in each inner loop. Note that, when $\tau_{1,\max} = \tau_{2,\max} = 0$, BAD-VAMP (i.e., Alg. 4) reduces to bilinear EM-VAMP (i.e., Alg. 3). In practice, we

noticed that re-estimating the precision γ_1 is important especially in problems like dictionary learning and often a single EM iteration for γ_1 (i.e., $\tau_{1,\max} = 1$) suffices. Re-estimating the precisions γ_2 in VAMP can be very computationally expensive as it involves a matrix inversion, however, in many applications we noticed that γ_2 is often accurate and we can skip it.

Algorithm 4 Bilinear Adaptive VAMP

```

1: initialize:
    $\forall l : \mathbf{r}_{1,l}^0, \gamma_{1,l}^0, \boldsymbol{\theta}_x^0, \boldsymbol{\theta}_A^0, \gamma_w^0$ 
2: for  $t = 0, \dots, T_{\max}$  do
3:   for  $\tau = 0, \dots, \tau_{1,\max}$  do
4:      $\forall l : \mathbf{x}_{1,l}^t \leftarrow \mathbf{g}_1(\mathbf{r}_{1,l}^t, \gamma_{1,l}^t; \boldsymbol{\theta}_x^t)$ 
5:      $\forall l : 1/\eta_{1,l}^t \leftarrow \langle \mathbf{g}'_1(\mathbf{r}_{1,l}^t, \gamma_{1,l}^t; \boldsymbol{\theta}_x^t) \rangle / \gamma_{1,l}^t$ 
6:      $\forall l : 1/\gamma_{1,l}^t \leftarrow \frac{1}{N} \|\mathbf{x}_{1,l}^t - \mathbf{r}_{1,l}^t\|^2 + 1/\eta_{1,l}^t$ 
7:      $q_1^t(\mathbf{X}) \propto \prod_{l=1}^L p_{\mathbf{x}}(\mathbf{x}_l; \boldsymbol{\theta}_x^t) e^{-\frac{1}{2}\gamma_{1,l}^t \|\mathbf{x}_l - \mathbf{r}_{1,l}^t\|^2}$ 
8:      $\boldsymbol{\theta}_x^t \leftarrow \arg \max_{\boldsymbol{\theta}_x} \mathbb{E}[\ln p_{\mathbf{x}}(\mathbf{X}; \boldsymbol{\theta}_x) | q_1^t]$ 
9:   end for
10:   $\boldsymbol{\theta}_x^{t+1} = \boldsymbol{\theta}_x^t$ 
11:   $\forall l : \gamma_{2,l}^t = \eta_{1,l}^t - \gamma_{1,l}^t$ 
12:   $\forall l : \mathbf{r}_{2,l}^t = (\eta_{1,l}^t \mathbf{x}_{1,l}^t - \gamma_{1,l}^t \mathbf{r}_{1,l}^t) / \gamma_{2,l}^t$ 
13:  for  $\tau = 0, \dots, \tau_{2,\max}$  do
14:     $\forall l : \mathbf{x}_{2,l}^t \leftarrow \mathbf{g}_{2,l}(\mathbf{r}_{2,l}^t, \gamma_{2,l}^t; \boldsymbol{\theta}_A^t, \gamma_w^t)$ 
15:     $\forall l : 1/\eta_{2,l}^t \leftarrow \langle \mathbf{g}'_{2,l}(\mathbf{r}_{2,l}^t, \gamma_{2,l}^t; \boldsymbol{\theta}_A^t, \gamma_w^t) \rangle / \gamma_{2,l}^t$ 
16:     $\forall l : 1/\gamma_{2,l}^t \leftarrow \frac{1}{N} \|\mathbf{x}_{2,l}^t - \mathbf{r}_{2,l}^t\|^2 + 1/\eta_{2,l}^t$ 
17:     $q_2^t(\mathbf{X}) \propto \prod_l p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}_l | \mathbf{x}_l; \boldsymbol{\theta}_A^t, \gamma_w^t) e^{-\frac{1}{2}\gamma_{2,l}^t \|\mathbf{x}_l - \mathbf{r}_{2,l}^t\|^2}$ 
18:     $\boldsymbol{\theta}_A^t \leftarrow \arg \max_{\boldsymbol{\theta}_A} \mathbb{E}[\ln p_{\mathbf{y}|\mathbf{x}}(\mathbf{Y} | \mathbf{X}; \boldsymbol{\theta}_A, \gamma_w^t) | \mathbf{Y}, q_2^t]$ 
19:     $\gamma_w^t \leftarrow \arg \max_{\gamma_w} \mathbb{E}[\ln p_{\mathbf{y}|\mathbf{x}}(\mathbf{Y} | \mathbf{X}; \boldsymbol{\theta}_A^t, \gamma_w) | \mathbf{Y}, q_2^t]$ 
20:  end for
21:   $\boldsymbol{\theta}_A^{t+1} = \boldsymbol{\theta}_A^t$ 
22:   $\gamma_w^{t+1} = \gamma_w^t$ 
23:   $\forall l : \gamma_{1,l}^{t+1} = \eta_{2,l}^t - \gamma_{2,l}^t$ 
24:   $\forall l : \mathbf{r}_{1,l}^{t+1} = (\eta_{2,l}^t \mathbf{x}_{2,l}^t - \gamma_{2,l}^t \mathbf{r}_{2,l}^t) / \gamma_{1,l}^{t+1}$ 
25: end for

```

We now provide additional details on the steps in Alg. 4.

3.7.1 Estimating the Matrix \mathbf{X}

Recalling the definition of $\mathbf{g}_{2,l}$ in (3.28), the form of $p_{\mathbf{y}|\mathbf{x}}$ in (3.15) implies that

$$\mathbf{g}_{2,l}(\mathbf{r}_{2,l}^t, \gamma_{2,l}^t; \boldsymbol{\theta}_A^t, \gamma_w^t) = \mathbf{C}_l^t (\gamma_{2,l}^t \mathbf{r}_{2,l}^t + \gamma_w^t \mathbf{A}(\boldsymbol{\theta}_A^t)^\top \mathbf{y}_l) \quad (3.49)$$

$$\langle \mathbf{g}'_{2,l}(\mathbf{r}_{2,l}^t, \gamma_{2,l}^t; \boldsymbol{\theta}_A^t, \gamma_w^t) \rangle = \gamma_{2,l}^t \text{tr} \{ \mathbf{C}_l^t \} / N \quad (3.50)$$

for the covariance matrix

$$\mathbf{C}_l^t \triangleq (\gamma_{2,l}^t \mathbf{I}_N + \gamma_w^t \mathbf{A}(\boldsymbol{\theta}_A^t)^\top \mathbf{A}(\boldsymbol{\theta}_A^t))^{-1}. \quad (3.51)$$

To avoid computing a separate matrix inverse (3.51) for each column $l = 1, \dots, L$, one could instead compute the eigenvalue decomposition

$$\mathbf{A}(\boldsymbol{\theta}_A^t)^\top \mathbf{A}(\boldsymbol{\theta}_A^t) = \mathbf{U}^t \text{Diag}(\mathbf{s}^t) \mathbf{U}^{t\top}, \quad (3.52)$$

and then leverage the fact that

$$\mathbf{C}_l^t = (\gamma_{2,l}^t \mathbf{I} + \gamma_w^t \mathbf{U}^t \text{Diag}(\mathbf{s}^t) \mathbf{U}^{t\top})^{-1} \quad (3.53)$$

$$= (\mathbf{U}^t (\gamma_{2,l}^t \mathbf{I} + \gamma_w^t \text{Diag}(\mathbf{s}^t)) \mathbf{U}^{t\top})^{-1} \quad (3.54)$$

$$= \mathbf{U}^t \text{Diag}(\gamma_{2,l}^t \mathbf{1} + \gamma_w^t \mathbf{s}^t)^{-1} \mathbf{U}^{t\top}, \quad (3.55)$$

which reduces to the inversion of a diagonal matrix for each $l = 1, \dots, L$.

3.7.2 Learning the Parameters $\boldsymbol{\theta}_A$

We now provide details on the update of $\boldsymbol{\theta}_A$ line 18 of Alg. 4. Given the form of $p_{\mathbf{Y}|\mathbf{X}}$ in (3.14)-(3.15), we have that

$$\ln p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\mathbf{X}; \boldsymbol{\theta}_A, \gamma_w) = \frac{ML}{2} \ln \gamma_w - \frac{\gamma_w}{2} \|\mathbf{Y} - \mathbf{A}(\boldsymbol{\theta}_A)\mathbf{X}\|_F^2 + \text{const} \quad (3.56)$$

$$\begin{aligned} &= \frac{ML}{2} \ln \gamma_w - \frac{\gamma_w}{2} \left(\text{tr}\{\mathbf{Y}\mathbf{Y}^\top\} - 2 \text{tr}\{\mathbf{A}(\boldsymbol{\theta}_A)\mathbf{X}\mathbf{Y}^\top\} \right. \\ &\quad \left. + \text{tr}\{\mathbf{A}(\boldsymbol{\theta}_A)\mathbf{X}\mathbf{X}^\top\mathbf{A}(\boldsymbol{\theta}_A)^\top\} \right) + \text{const}. \end{aligned} \quad (3.57)$$

Since the first and second moments of belief q_2^t are the following

$$\mathbb{E}[\mathbf{X}|q_2^t] = \mathbf{X}_2^t \quad (3.58)$$

$$\mathbb{E}[\mathbf{X}\mathbf{X}^\top|q_2^t] = \sum_{l=1}^L \mathbb{E}[\mathbf{x}_l\mathbf{x}_l^\top|q_{2,l}^t] = \mathbf{X}_2^t\mathbf{X}_2^{t^\top} + \underbrace{\sum_{l=1}^L \mathbf{C}_l^t}_{\triangleq \mathbf{C}^t}, \quad (3.59)$$

we have that

$$\begin{aligned} &\mathbb{E}[\ln p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\mathbf{X}; \boldsymbol{\theta}_A, \gamma_w) | \mathbf{Y}, q_2^t] \\ &= \frac{ML}{2} \ln \gamma_w - \frac{\gamma_w}{2} \left(\text{tr}\{\mathbf{Y}\mathbf{Y}^\top\} - 2 \text{tr}\{\mathbf{A}(\boldsymbol{\theta}_A)\mathbf{X}_2^t\mathbf{Y}^\top\} \right. \\ &\quad \left. + \text{tr}\{\mathbf{A}(\boldsymbol{\theta}_A)\mathbf{X}_2^t\mathbf{X}_2^{t^\top}\mathbf{A}(\boldsymbol{\theta}_A)^\top\} + \text{tr}\{\mathbf{A}(\boldsymbol{\theta}_A)\mathbf{C}^t\mathbf{A}(\boldsymbol{\theta}_A)^\top\} \right) \end{aligned} \quad (3.60)$$

$$\begin{aligned} &= \frac{ML}{2} \ln \gamma_w - \frac{\gamma_w}{2} \left(\|\mathbf{Y} - \mathbf{A}(\boldsymbol{\theta}_A)\mathbf{X}_2^t\|_F^2 \right. \\ &\quad \left. + \text{tr}\{\mathbf{A}(\boldsymbol{\theta}_A)\mathbf{C}^t\mathbf{A}(\boldsymbol{\theta}_A)^\top\} \right) + \text{const}. \end{aligned} \quad (3.61)$$

To maximize (3.61) over $\boldsymbol{\theta}_A = [\theta_{A,1}, \dots, \theta_{A,Q}]$ with fixed γ_w , we consider the affine-linear model

$$\mathbf{A}(\boldsymbol{\theta}_A) = \mathbf{A}_0 + \sum_{i=1}^Q \theta_{A,i}\mathbf{A}_i, \quad (3.62)$$

noting that non-linear models could be handled using similar techniques. Plugging (3.62) into (3.60), we get

$$\begin{aligned} & \mathbb{E} [\ln p_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y}|\mathbf{X}; \boldsymbol{\theta}_A, \gamma_w) | \mathbf{Y}, q_2^t] \\ &= \text{const} - \frac{\gamma_w}{2} \sum_{i=1}^Q \sum_{j=1}^Q \theta_{A,i} \text{tr} \{ \mathbf{A}_i (\mathbf{C}^t + \mathbf{X}_2^t \mathbf{X}_2^{t\top}) \mathbf{A}_j^\top \} \theta_{A,j} \\ & \quad - \gamma_w \sum_{i=1}^Q \theta_{A,i} (\text{tr} \{ \mathbf{A}_i (\mathbf{C}^t + \mathbf{X}_2^t \mathbf{X}_2^{t\top}) \mathbf{A}_0^\top \} - \text{tr} \{ \mathbf{A}_i \mathbf{X}_2^t \mathbf{Y}^\top \}) \end{aligned} \quad (3.63)$$

$$= -\frac{\gamma_w}{2} (\boldsymbol{\theta}_A^\top \mathbf{H}^t \boldsymbol{\theta}_A - 2\boldsymbol{\theta}_A^\top \boldsymbol{\beta}^t) + \text{const} \quad (3.64)$$

for the quantities \mathbf{H}^t and $\boldsymbol{\beta}^t$ defined as

$$[\mathbf{H}^t]_{ij} \triangleq \text{tr} \{ \mathbf{A}_i (\mathbf{C}^t + \mathbf{X}_2^t \mathbf{X}_2^{t\top}) \mathbf{A}_j^\top \} \quad (3.65)$$

$$= \text{tr} \{ \mathbf{A}_j^\top \mathbf{A}_i (\mathbf{C}^t + \mathbf{X}_2^t \mathbf{X}_2^{t\top}) \} \quad (3.66)$$

$$[\boldsymbol{\beta}^t]_i \triangleq \text{tr} \{ \mathbf{A}_i \mathbf{X}_2^t \mathbf{Y}^\top \} - \text{tr} \{ \mathbf{A}_i (\mathbf{C}^t + \mathbf{X}_2^t \mathbf{X}_2^{t\top}) \mathbf{A}_0^\top \} \quad (3.67)$$

$$= \text{tr} \{ \mathbf{Y}^\top \mathbf{A}_i \mathbf{X}_2^t \} - \text{tr} \{ \mathbf{A}_0^\top \mathbf{A}_i (\mathbf{C}^t + \mathbf{X}_2^t \mathbf{X}_2^{t\top}) \}, \quad (3.68)$$

where $\mathbf{A}_j^\top \mathbf{A}_i$ and $\mathbf{Y}^\top \mathbf{A}_i$ can be pre-computed to speed up the algorithm. We can find the maximizer of (3.64) by computing its gradient w.r.t. $\boldsymbol{\theta}_A$ and setting it to zero. We then get the update $\boldsymbol{\theta}_A^{t+1}$ as

$$\mathbf{H}^t \boldsymbol{\theta}_A^{t+1} = \boldsymbol{\beta}^t \quad (3.69)$$

$$\Rightarrow \boldsymbol{\theta}_A^{t+1} = (\mathbf{H}^t)^{-1} \boldsymbol{\beta}^t. \quad (3.70)$$

A special case of (3.62) is where $\mathbf{A}(\cdot)$ has no structure, i.e.,

$$\mathbf{A}(\boldsymbol{\theta}_A) = \sum_{m=1}^M \sum_{n=1}^N \theta_{A,m,n} \mathbf{e}_m \mathbf{e}_n^\top. \quad (3.71)$$

where \mathbf{e}_m denotes the m th standard basis vector. In this case, it can be shown that

$$\mathbf{A}(\boldsymbol{\theta}_A^{t+1}) = \mathbf{Y} \mathbf{X}_2^{t\top} (\mathbf{C}^t + \mathbf{X}_2^t \mathbf{X}_2^{t\top})^{-1}. \quad (3.72)$$

3.7.3 Learning the Noise Precision γ_w

We now provide details on the update of the noise precision γ_w in line 19 of Alg. 4. To maximize (3.61) over γ_w with fixed $\boldsymbol{\theta}_A = \boldsymbol{\theta}_A^{t+1}$, we search for the values of γ_w that zero the derivative of (3.61). The unique solution is straightforwardly shown to be

$$1/\gamma_w^{t+1} = \frac{1}{ML} \left(\|\mathbf{Y} - \mathbf{A}(\boldsymbol{\theta}_A^{t+1})\mathbf{X}_2^t\|_F^2 + \text{tr} \{ \mathbf{A}(\boldsymbol{\theta}_A^{t+1})\mathbf{C}^t \mathbf{A}(\boldsymbol{\theta}_A^{t+1})^\top \} \right). \quad (3.73)$$

In Alg. 5, BAd-VAMP is rewritten with detailed expressions for the updates of $\mathbf{x}_{2,l}^t$, $\eta_{2,l}^t$, $\boldsymbol{\theta}_A^t$, and γ_w^t .

Algorithm 5 Bilinear Adaptive VAMP (Detailed)

```

1: initialize:
    $\forall l : \mathbf{r}_{1,l}^0, \gamma_{1,l}^0, \boldsymbol{\theta}_x^0, \boldsymbol{\theta}_A^0, \gamma_w^0$ 
2: for  $t = 0, \dots, T_{\max}$  do
3:   for  $\tau = 0, \dots, \tau_{1,\max}$  do
4:      $\forall l : \mathbf{x}_{1,l}^t \leftarrow \mathbf{g}_1(\mathbf{r}_{1,l}^t, \gamma_{1,l}^t; \boldsymbol{\theta}_x^t)$ 
5:      $\forall l : 1/\eta_{1,l}^t \leftarrow \langle \mathbf{g}'_1(\mathbf{r}_{1,l}^t, \gamma_{1,l}^t; \boldsymbol{\theta}_x^t) \rangle / \gamma_{1,l}^t$ 
6:      $\forall l : 1/\gamma_{1,l}^t \leftarrow \frac{1}{N} \|\mathbf{x}_{1,l}^t - \mathbf{r}_{1,l}^t\|^2 + 1/\eta_{1,l}^t$ 
7:      $q_1^t(\mathbf{X}) \propto \prod_{l=1}^L p_{\mathbf{x}}(\mathbf{x}_l; \boldsymbol{\theta}_x^t) e^{-\frac{1}{2} \gamma_{1,l}^t \|\mathbf{x}_l - \mathbf{r}_{1,l}^t\|^2}$ 
8:      $\boldsymbol{\theta}_x^t \leftarrow \arg \max_{\boldsymbol{\theta}_x} \mathbb{E}[\ln p_{\mathbf{x}}(\mathbf{X}; \boldsymbol{\theta}_x) | q_1^t]$ 
9:   end for
10:   $\boldsymbol{\theta}_x^{t+1} = \boldsymbol{\theta}_x^t$ 
11:   $\forall l : \gamma_{2,l}^t = \eta_{1,l}^t - \gamma_{1,l}^t$ 
12:   $\forall l : \mathbf{r}_{2,l}^t = (\eta_{1,l}^t \mathbf{x}_{1,l}^t - \gamma_{1,l}^t \mathbf{r}_{1,l}^t) / \gamma_{2,l}^t$ 
13:  for  $\tau = 0, \dots, \tau_{2,\max}$  do
14:     $\forall l : \mathbf{C}_l^t \leftarrow (\gamma_{2,l}^t \mathbf{I}_N + \gamma_w^t \mathbf{A}(\boldsymbol{\theta}_A^t)^\top \mathbf{A}(\boldsymbol{\theta}_A^t))^{-1}$ 
15:     $\forall l : \mathbf{x}_{2,l}^t \leftarrow \mathbf{C}_l^t (\gamma_{2,l}^t \mathbf{r}_{2,l}^t + \gamma_w^t \mathbf{A}(\boldsymbol{\theta}_A^t)^\top \mathbf{y}_l)$ 
16:     $\forall l : 1/\eta_{2,l}^t \leftarrow \text{tr}\{\mathbf{C}_l^t\} / N$ 
17:     $\mathbf{C}^t \leftarrow \sum_{l=1}^L \mathbf{C}_l^t$ 
18:     $\forall i, j : [\mathbf{H}^t]_{ij} \leftarrow \text{tr}\{\mathbf{A}_j^\top \mathbf{A}_i (\mathbf{C}^t + \mathbf{X}_2^t \mathbf{X}_2^{t\top})\}$ 
19:     $\forall i : [\boldsymbol{\beta}^t]_i \leftarrow \text{tr}\{\mathbf{Y}^\top \mathbf{A}_i \mathbf{X}_2^t\}$ 
20:     $\phantom{\forall i :} - \text{tr}\{\mathbf{A}_0^\top \mathbf{A}_i (\mathbf{C}^t + \mathbf{X}_2^t \mathbf{X}_2^{t\top})\}$ 
21:     $\boldsymbol{\theta}_A^t \leftarrow (\mathbf{H}^t)^{-1} \boldsymbol{\beta}^t$ 
22:     $1/\gamma_w^t \leftarrow \frac{1}{ML} (\|\mathbf{Y} - \mathbf{A}(\boldsymbol{\theta}_A^t) \mathbf{X}_2^t\|_F^2$ 
23:     $\phantom{1/\gamma_w^t \leftarrow} + \text{tr}\{\mathbf{A}(\boldsymbol{\theta}_A^t) \mathbf{C}^t \mathbf{A}(\boldsymbol{\theta}_A^t)^\top\})$ 
24:  end for
25:   $\boldsymbol{\theta}_A^{t+1} = \boldsymbol{\theta}_A^t$ 
26:   $\gamma_w^{t+1} = \gamma_w^t$ 
27:   $\forall l : \gamma_{1,l}^{t+1} = \eta_{2,l}^t - \gamma_{2,l}^t$ 
28:   $\forall l : \mathbf{r}_{1,l}^{t+1} = (\eta_{2,l}^t \mathbf{x}_{2,l}^t - \gamma_{2,l}^t \mathbf{r}_{2,l}^t) / \gamma_{1,l}^{t+1}$ 
29: end for

```

3.7.4 Algorithm Enhancements

We now propose several enhancements to the BAd-VAMP algorithm presented in Alg. 4 and detailed in Alg. 5.

Damping

For fixed Θ^t and infinitely large right-rotationally invariant $\mathbf{A}(\theta_A^t)$, the state-evolution of VAMP guarantees its convergence. But when $\mathbf{A}(\theta_A^t)$ deviates from this assumption, damping the VAMP iterations can help maintain convergence [75]. With damping, lines 25-26 of Alg. 5 (or lines 23-24 of Alg. 4) would be replaced by

$$\gamma_{1,l}^{t+1} = (1 - \zeta)\gamma_{1,l}^t + \zeta(\eta_{2,l}^t - \gamma_{2,l}^t) \quad (3.74)$$

$$\mathbf{r}_{1,l}^{t+1} = (1 - \zeta)\mathbf{r}_{1,l}^t + \zeta(\eta_{2,l}^t \mathbf{x}_{2,l}^t - \gamma_{2,l}^t \mathbf{r}_{2,l}^t) / (\eta_{2,l}^t - \gamma_{2,l}^t) \quad (3.75)$$

for some $\zeta \in (0, 1)$. We can perform a similar damping on $\{(\mathbf{r}_{2,l}^t, \gamma_{2,l}^t)\}_{l=1}^L$ and replace lines 11-12 of Alg. 5 (or lines 11-12 of Alg. 4) by the following,

$$\gamma_{2,l}^t = (1 - \zeta)\gamma_{1,l}^{t-1} + \zeta(\eta_{1,l}^t - \gamma_{1,l}^t) \quad (3.76)$$

$$\mathbf{r}_{1,l}^t = (1 - \zeta)\mathbf{r}_{1,l}^{t-1} + \zeta(\eta_{1,l}^t \mathbf{x}_{1,l}^t - \gamma_{1,l}^t \mathbf{r}_{1,l}^t) / (\eta_{1,l}^t - \gamma_{1,l}^t). \quad (3.77)$$

The case $\zeta = 1$ corresponds to no damping.

Negative precisions

Sometimes the precisions $\{\gamma_{1,l}, \gamma_{2,l}\}_{l=1}^L$ can be negative. To avoid the precisions from being negative, we suggest to restrict them to the interval $[\gamma_{\min}, \infty)$, for very small $\gamma_{\min} > 0$, in lines 11 and 25 of Alg. 5 (or lines 11 and 23 of Alg. 4).

Restarts

Due to the non-convex nature of the bilinear inference problem, the algorithm may get stuck at local minima or slowed by saddle points. To mitigate these issues, it sometimes helps to restart the algorithm. For each restart, we suggest to initialize $\boldsymbol{\theta}_A^0$ at the final estimate of $\boldsymbol{\theta}_A$ returned by the previous run.

3.8 Relation to Previous Work

The proposed Bilinear Adaptive VAMP algorithm extends the (linear) Adaptive VAMP algorithm of [37] from the case where $\mathbf{A}(\boldsymbol{\theta}_A)$ is known to the case where $\mathbf{A}(\boldsymbol{\theta}_A)$ is unknown. In the known- $\mathbf{A}(\boldsymbol{\theta}_A)$ setting, where $\mathbf{A}(\boldsymbol{\theta}_A)$ is infinitely large and right-rotationally invariant, it was rigorously established in [37] that Adaptive VAMP obeys a state-evolution similar to that of VAMP, and that its estimates of $\{\boldsymbol{\theta}_x, \gamma_w\}$ are asymptotically consistent if certain identifiability conditions are met, i.e., they converge to the true values as $t \rightarrow \infty$. As future work, it would be interesting to understand whether Bilinear Adaptive VAMP also obeys a state evolution for certain classes of $\mathbf{A}(\cdot)$.

The proposed BAd-VAMP algorithm targets the same class of bilinear recovery problems as the EM-PBiGAMP algorithm from [66], and both leverage EM for automated hyperparameter tuning. The BiGAMP algorithm [67]- [68] is a special case of the PBiGAMP algorithm [66]. BiGAMP applies to the recovery of \mathbf{A} and \mathbf{X} from measurements of the form $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{W}$, whereas PBiGAMP applies to the recovery of \mathbf{b} and \mathbf{X} from $\mathbf{Y} = \mathbf{A}(\mathbf{b})\mathbf{X} + \mathbf{W}$ under known affine linear $\mathbf{A}(\cdot)$. Both can be combined with EM for hyperparameter learning. Note that the “AMP” aspects of these algorithms are fundamentally different. PBiGAMP treats the vectors $\{\mathbf{a}_{i,j}\}$

in (3.2) as i.i.d. Gaussian for its derivation, whereas BAd-VAMP treats the matrix $\mathbf{A}(\mathbf{b}) = \sum_{i=1}^Q b_i \mathbf{A}_i$ as right rotationally-invariant for its derivation. The latter allows more freedom in the singular values of $\mathbf{A}(\mathbf{b})$, which leads to increased robustness in practice, as demonstrated by the numerical experiments in Sec. 3.9.

BAd-VAMP and lifted VAMP both leverage the VAMP approach from [75] to solve bilinear inference problems. However, they do so in very different ways. As discussed in Sec. 3.1, lifted VAMP “lifts” the bilinear problem into a higher-dimensional linear problem, and then uses non-separable denoising to jointly estimate \mathbf{b} and \mathbf{c} in (3.5). An unfortunate consequence of lifting is a possibly significant increase in computational complexity and memory. In contrast, BAd-VAMP avoids lifting, and it employs EM to estimate \mathbf{b} and VAMP to estimate \mathbf{c} . Interestingly, Sec. 3.9 shows BAd-VAMP performing equal or better to lifted VAMP in all experiments.

3.9 Numerical Experiments

In this section, we present the results of numerical experiments that study the behavior of the BAd-VAMP algorithm from Alg. 5, in comparison to other state-of-the-art algorithms, on several bilinear recovery problems. In all cases, we ran BAd-VAMP with $\tau_{1,\max} = 1$ and $\tau_{2,\max} = 0$ inner EM iterations and we assumed that the signal prior $p_{\mathbf{x}}$ is fully known (i.e., $\boldsymbol{\theta}_x$ is known). We nominally used a damping coefficient of $\zeta = 0.8$ and minimum precision of $\gamma_{\min} = 10^{-6}$. We initialized BAd-VAMP by $\gamma_w^0 = 0.1$, $\gamma_{1,l}^0 = 10^{-3} \forall l$, and we set $\mathbf{r}_{1,l}^0$ and $\boldsymbol{\theta}_A^0$ to random vectors drawn i.i.d. from $\mathcal{N}(0, 10)$ and $\mathcal{N}(0, 1)$ respectively. Unless otherwise noted, no restarts were used in that experiment.

3.9.1 CS with Matrix Uncertainty

We consider the compressive sensing (CS) with matrix uncertainty [105] problem that we discussed in Sec. 2.5.2. We describe the problem formulation here for convenience. The goal is to recover the K -sparse signal $\mathbf{c} \in \mathbb{R}^N$ and the uncertainty parameters $\mathbf{b} \in \mathbb{R}^Q$ from length- M measurements \mathbf{y} ,

$$\mathbf{y} = \mathbf{A}(\mathbf{b})\mathbf{c} + \mathbf{w} \quad (3.78a)$$

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w) \quad (3.78b)$$

$$\mathbf{A}(\mathbf{b}) = \mathbf{A}_0 + \sum_{i=1}^Q b_i \mathbf{A}_i, \quad (3.78c)$$

where $\{\mathbf{A}_i\}_{i=0}^Q$ are known matrices. For our experiments, we used $Q = 10$, $N = 256$ and $K = 10$. We selected the noise precision γ_w so that the SNR as defined is 40 dB,

$$\text{SNR} \triangleq \frac{\mathbb{E}[\|\mathbf{A}(\mathbf{b})\mathbf{c}\|^2]}{\mathbb{E}[\|\mathbf{w}\|^2]} \quad (3.79)$$

$$= \frac{\gamma_w}{M} \mathbb{E}[\|\mathbf{A}(\mathbf{b})\mathbf{c}\|^2]. \quad (3.80)$$

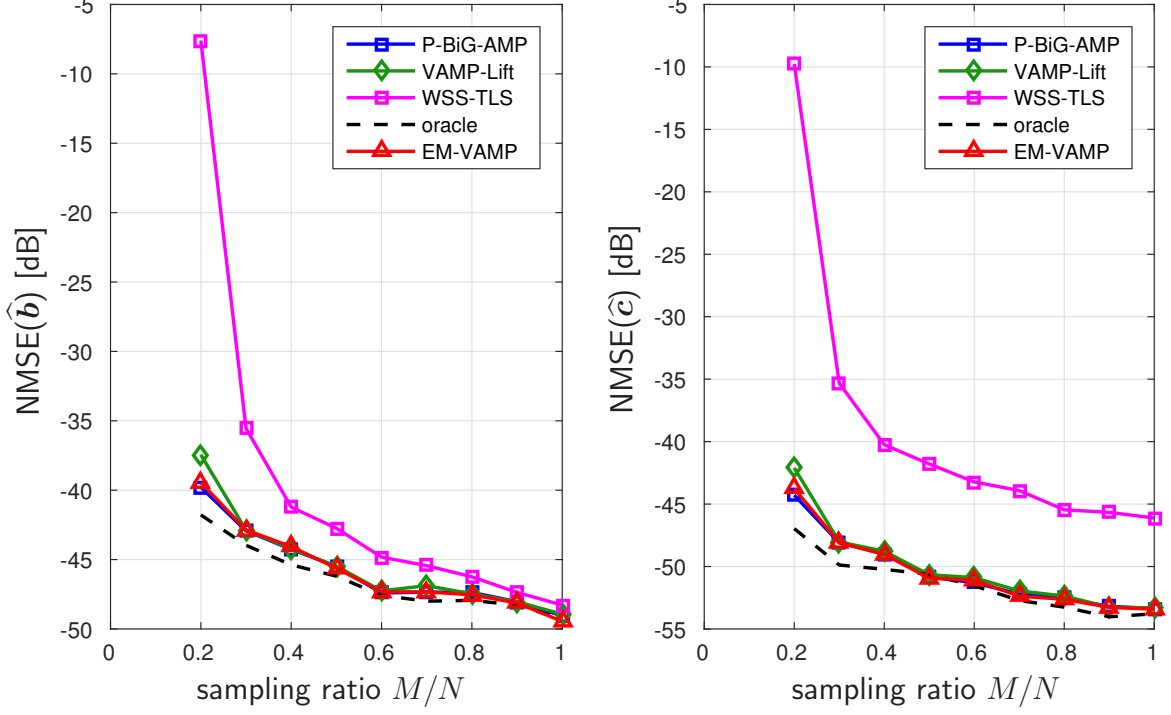


Figure 3.1: CS with matrix uncertainty: Median NMSE (over 50 trials) on signal \mathbf{c} and uncertainty parameters \mathbf{b} versus sampling ratio M/N .

Also, the uncertainty parameters \mathbf{b} were drawn $\mathcal{N}(\mathbf{0}, \mathbf{I})$, and \mathbf{c} was drawn with uniformly random support and with K non-zero elements from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. We measured performance of our algorithm using the following metrics

$$\text{NMSE}(\hat{\mathbf{b}}) \triangleq \frac{\|\hat{\mathbf{b}} - \mathbf{b}\|^2}{\|\mathbf{b}\|^2} \quad (3.81)$$

$$\text{NMSE}(\hat{\mathbf{c}}) \triangleq \frac{\|\hat{\mathbf{c}} - \mathbf{c}\|^2}{\|\mathbf{c}\|^2}. \quad (3.82)$$

As a reference, we considered two oracle estimators: the MMSE estimator for \mathbf{b} assuming \mathbf{c} is known, and the MMSE estimator for \mathbf{c} assuming \mathbf{b} and the support of \mathbf{c} are known.

For our first experiment, the elements of $\{\mathbf{A}_i\}_{i=1}^Q$ were drawn i.i.d. $\mathcal{N}(0, 1)$ and the elements of \mathbf{A}_0 were drawn $\mathcal{N}(0, 20)$. BAd-VAMP was run for a maximum of 200 iterations with a maximum of 2 restarts and damping $\zeta = 0.86$.

Figure 3.1 shows that the AMP-based algorithms gave near-oracle performance for the tested range of M/N , although lifted VAMP performed slightly worse than the others when $M/N = 0.2$. In contrast, the performance of WSS-TLS from the award-winning paper [105] was significantly worse than the AMP approaches. WSS-TLS aims to solve the non-convex optimization problem

$$(\hat{\mathbf{b}}, \hat{\mathbf{c}}) = \arg \min_{\mathbf{b}, \mathbf{c}} \left\| \left(\mathbf{A}_0 + \sum_{i=1}^Q b_i \mathbf{A}_i \right) \mathbf{c} - \mathbf{y} \right\|^2 + \|\mathbf{b}\|^2 / \gamma_w + \lambda \|\mathbf{c}\|_1 \quad (3.83)$$

using alternating minimization. For WSS-TLS, we used oracle knowledge of γ_w , oracle tuning of the regularization parameter λ , and code from the authors' website.

For our second experiment, we tested algorithm's robustness to non-zero mean in $\mathbf{A}(\mathbf{b})$ ⁵, since this is a known issue with many AMP algorithms [18, 74, 96]. For this, we fixed the sampling ratio at $M/N = 0.6$, drew the elements of $\{\mathbf{A}_i\}_{i=1}^Q$ from i.i.d. $\mathcal{N}(\mu, 1)$, and drew the elements of \mathbf{A}_0 from i.i.d. $\mathcal{N}(\mu, 20)$. Figure 3.2 reports the median NMSE versus mean μ , and shows that BAd-VAMP is much more robust to $\mu > 0$ than the other tested AMP algorithms as well as WSS-TLS.

Figure 3.3 shows the runtime of the algorithms in Figure 3.1. Our implementation used MATLAB (R2015b) on an RHEL workstation with an 8-core Intel i7 processor. Although, for WSS-TLS, we used a grid-search to optimize λ in (3.83), Figure 3.3

⁵For the simpler case where \mathbf{b} is known and the objective is to recover \mathbf{c} from $\mathbf{y} = \mathbf{A}\mathbf{c} + \mathbf{w}$, modifications of AMP that temporarily remove the mean from \mathbf{A} have been proposed [96]. However, it is not clear how to extend this approach to the bilinear problem of recovering \mathbf{b} and \mathbf{c} from $\mathbf{y} = \mathbf{A}(\mathbf{b})\mathbf{c} + \mathbf{w}$.

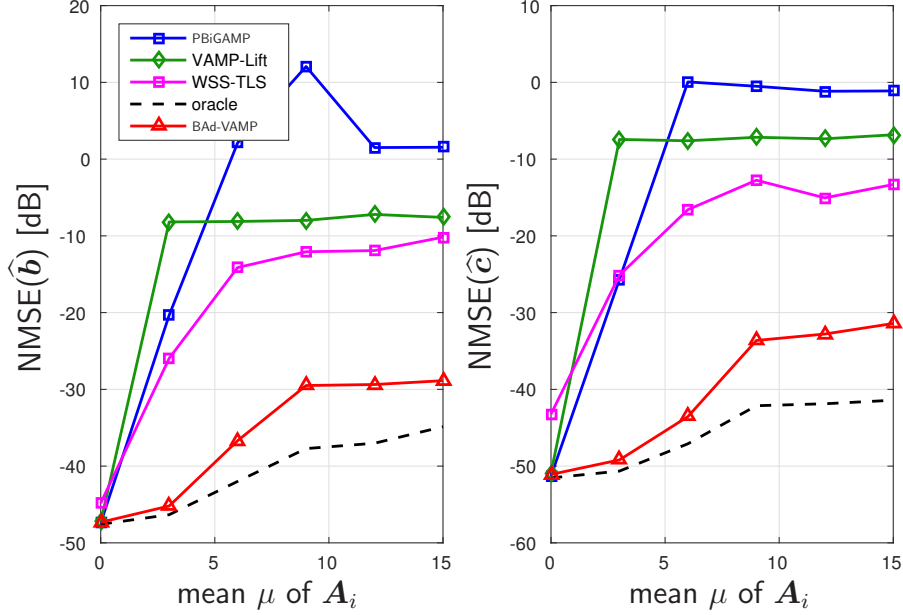


Figure 3.2: CS with matrix uncertainty: Median NMSE (over 50 trials) on signal \mathbf{c} and uncertainty parameters \mathbf{b} versus mean of matrices \mathbf{A}_i at $M/N = 0.6$.

only shows the runtime of WSS-TLS after λ was chosen. Figure 3.3 shows BAd-VAMP running much faster than lifted VAMP and WSS-TLS, and slightly slower than PBiGAMP.

3.9.2 Self-Calibration

We consider the self-calibration problem [51] that we also discussed earlier in Sec. 2.5.2. We present the problem formulation here for convenience. The goal is to recover the K -sparse signal vector \mathbf{c} and the calibration parameters \mathbf{b} from measurements of the form $\mathbf{y} = \text{Diag}(\mathbf{H}\mathbf{b})\Psi\mathbf{c}$ with known $\mathbf{H} \in \mathbb{R}^{M \times Q}$ and $\Psi \in \mathbb{R}^{M \times N}$. Here, $\mathbf{H}\mathbf{b}$ represents an unknown vector of gains on the measurements, where the gain vector is believed to lie in the Q -dimensional subspace spanned by the columns of \mathbf{H} . For our experiment, $M = 128$, $N = 256$, Ψ and \mathbf{b} were drawn i.i.d. $\mathcal{N}(0, 1)$, \mathbf{H} was

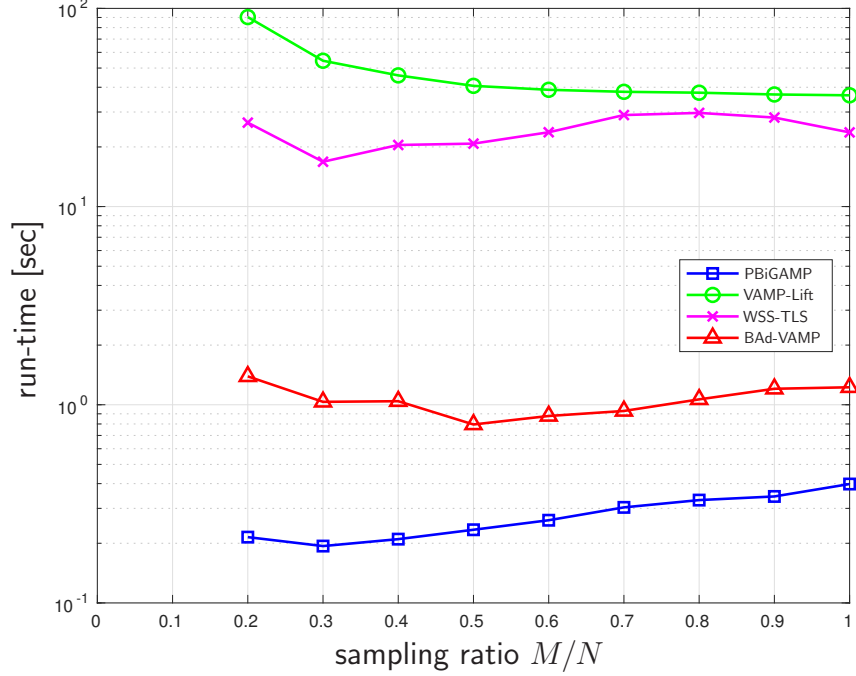


Figure 3.3: CS with matrix uncertainty: Median run-time in seconds (over 10 trials) versus sampling ratio M/N .

constructed using Q randomly selected columns of the Hadamard matrix, and \mathbf{c} was drawn with uniformly random support and with K non-zero elements from $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

Figure 3.4 shows the rate of successful recovery versus subspace dimension Q and sparsity K for several algorithms. A recovery $(\hat{\mathbf{b}}, \hat{\mathbf{c}})$ was considered “successful” when $\|\hat{\mathbf{b}}\hat{\mathbf{c}}^T - \mathbf{b}\mathbf{c}^T\|_F^2 / \|\mathbf{b}\mathbf{c}^T\|_F^2 \leq -50$ dB. From the figure, we see that the performance of BAd-VAMP is similar to that of EM-PBiGAMP, and even slightly better when Q is small and K is large. Meanwhile, BAd-VAMP appears significantly better than both lifted VAMP and SparseLift from [51]. SparseLift is a convex relaxation with provable guarantees [51]. For computational reasons (recall the discussion in Sec. 3.2), it was difficult to simulate lifted VAMP for $Q \geq 10$.

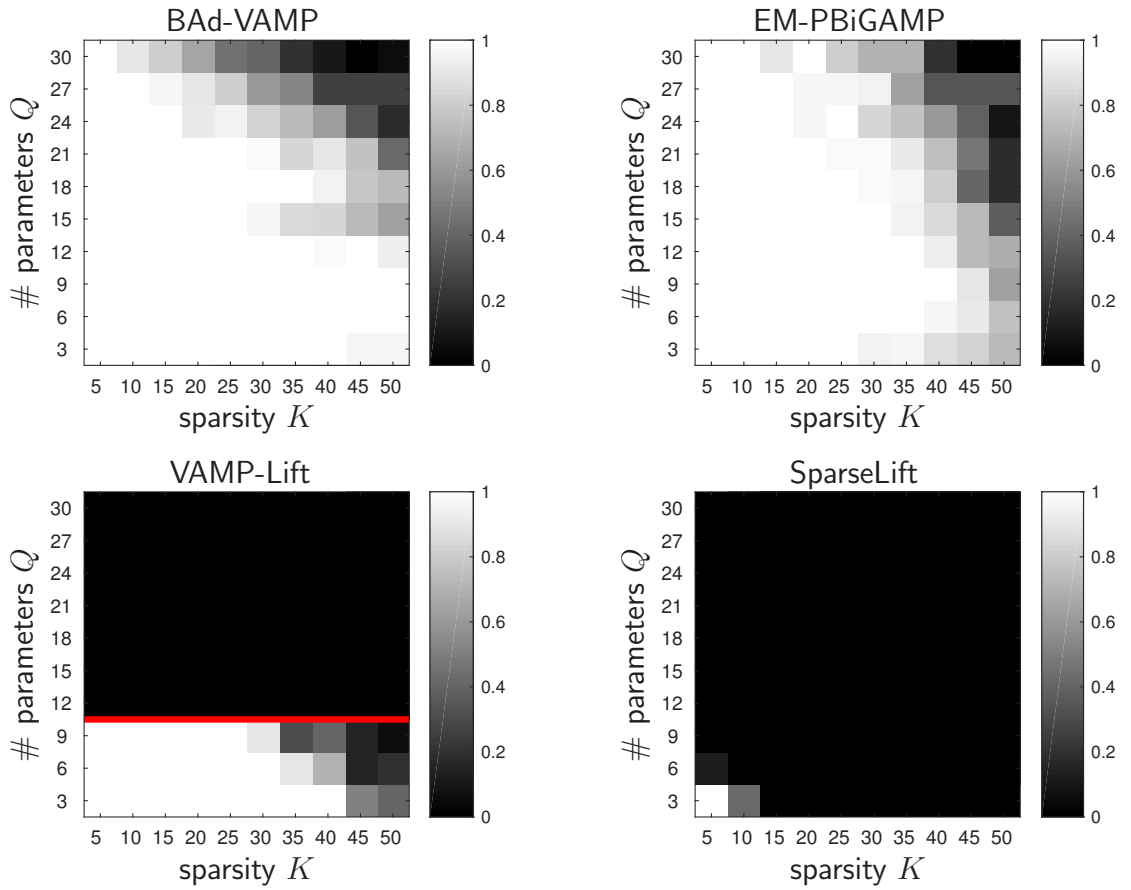


Figure 3.4: Self-calibration: Empirical success rate (over 50 trials) for several algorithms versus number of calibration parameters Q and sparsity K . For computational reasons, VAMP-Lift was simulated only for $Q < 10$.

3.9.3 Calibration in Tomography

We consider the problem of reconstructing an image from a sequence of tomographic projections, where the projections along each direction are scaled by an unknown calibration gain. In particular, let Ψ_ω be the tomographic projection matrix corresponding to angle $\omega \in [0, \pi]$. We used the matrix form of the Radon transform instead of the operator form (i.e., Radon transform) to avoid numerical error when implementing the adjoint. Our goal is to reconstruct the image \mathbf{x} from measurements

$$\mathbf{y} = \begin{bmatrix} b_1 \Psi_{\omega_1} \\ \vdots \\ b_K \Psi_{\omega_K} \end{bmatrix} \mathbf{x} + \mathbf{w} \quad (3.84a)$$

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w), \quad (3.84b)$$

where $b_k \sim \mathcal{N}(1, \sigma_b^2)$ are the unknown gains unknown. Note that, by defining

$$\mathbf{A}_k = \begin{bmatrix} \mathbf{0} \\ \vdots \\ \Psi_{\omega_k} \\ \vdots \\ \mathbf{0} \end{bmatrix} \quad \text{for } k = 1, \dots, K \quad (3.85)$$

we can write $\mathbf{y} = \mathbf{A}(\mathbf{b})\mathbf{x} + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w)$ for $\mathbf{A}(\mathbf{b}) = \sum_{k=1}^K b_k \mathbf{A}_k$, which matches (3.8).

In an attempt to solve the above problem, we used BAd-VAMP to recover the image \mathbf{x} while simultaneously learning the calibration gains \mathbf{b} . For this, we used BAd-VAMP in “plug-and-play” mode briefly described in Sec. 1.1, where the BM3D image denoiser [26] was used to implement the $\mathbf{g}_1(\cdot)$ function in Alg. 1. There is an inherent scaling ambiguity of the problem, i.e., if $(\hat{\mathbf{x}}, \hat{\mathbf{b}})$ is a solution then so is $(\alpha \hat{\mathbf{x}}, \alpha^{-1} \hat{\mathbf{b}})$ for any $\alpha > 0$. To avoid the scaling problem, we scaled the image estimate $\hat{\mathbf{x}}$ by the α that minimized $\|\mathbf{x} - \alpha \hat{\mathbf{x}}\|$ before computing the PSNR.

As baselines, we also tested the VAMP [75], total variation (TV) [10, 82] and regularization-by-denoising (RED) [76, 79] approaches (see descriptions below). These

approaches all assume a noisy *linear* data model of the form $\mathbf{y} = \widehat{\mathbf{A}}\mathbf{x} + \mathcal{N}(0, \gamma_w^{-1}\mathbf{I})$ with known $\widehat{\mathbf{A}}$. To apply them to (3.84b), we considered two cases: the *genie-calibrated* (GC) case, where a genie supplies the true gains \mathbf{b} and the algorithm uses $\widehat{\mathbf{A}} = \mathbf{A}(\mathbf{b})$, and the *un-calibrated* (UC) case, where \mathbf{b} is unknown and the algorithms assume $\widehat{\mathbf{A}} = \mathbf{A}(\mathbf{1})$. In the latter case, the $\widehat{\mathbf{A}}$ -based model is mismatched to the data-generation model (3.84b). For fair comparison, we scaled the GC and UC image estimates $\widehat{\mathbf{x}}$ by the α that minimized $\|\mathbf{x} - \alpha\widehat{\mathbf{x}}\|_2$ before computing the PSNR.

We now provide additional details on the experimental setup. For \mathbf{x} , we used the modified Shepp-Logan phantom of size 64×64 , shown in the top-left panel of Fig. 3.5. For $\mathbf{A}(\cdot)$, we used $K = 25$ projections spaced uniformly in $\omega \in [0, \pi]$. The calibration gains $\mathbf{b} \stackrel{\text{iid}}{\sim} \mathcal{N}(1, \sigma_b^2)$ were generated using $\sigma_b = 0.06$, and the noise precision γ_w was set to achieve an SNR of $\mathbb{E}[\|\mathbf{A}(\mathbf{b})\mathbf{x}\|^2] / \mathbb{E}[\|\mathbf{w}\|^2] = 40$ dB. The TV method [82] computes

$$\widehat{\mathbf{x}} = \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{y} - \widehat{\mathbf{A}}\mathbf{x}\|_2^2 + \lambda_t \|\nabla \mathbf{x}\|_{2,1} \right\} \quad (3.86)$$

for the isotropic TV operator

$$\|\nabla \mathbf{x}\|_{2,1} = \sum_{i,j} \sqrt{(x_{i,j} - x_{i,j-1})^2 + (x_{i,j} - x_{i-1,j})^2}. \quad (3.87)$$

We solved (3.86) using FASTA [38], and tuned λ_t to maximize the PSNR. RED [76, 79] solves the fixed-point equation

$$\widehat{\mathbf{A}}^H(\widehat{\mathbf{A}}\widehat{\mathbf{x}} - \mathbf{y}) + \lambda_r(\widehat{\mathbf{x}} - \boldsymbol{\rho}(\widehat{\mathbf{x}}, \tau)) = \mathbf{0} \quad (3.88)$$

for $\widehat{\mathbf{x}}$, where $\boldsymbol{\rho}(\cdot, \tau)$ is an image denoising algorithm with noise-variance τ . For our experiment, we used BM3D for $\boldsymbol{\rho}(\cdot, \tau)$, solved (3.88) using the ADMM method from [79] with 200 iterations, and tuned both λ_r and τ to maximize PSNR. For BAd-VAMP, we initialized \mathbf{b} to $\mathbf{1}$, used damping $\zeta = 0.1$, assumed known noise precision

Table 3.1: PSNR (dB) in the tomography experiment

measurements	BAd-VAMP	VAMP	RED	TV
genie calibrated (GC)	—	39.57	36.56	33.27
un-calibrated (UC)	38.27	31.62	31.24	26.79

γ_w , and used at most 100 iterations. For VAMP, we initialized $\gamma_1^0 = 10^{-4}$ and used at most 100 iterations.

Table 3.1 reports the median PSNR achieved by each algorithm across 10 random draws of \mathbf{b} and \mathbf{w} , Fig. 3.5 shows example image recoveries, and Fig. 3.6 shows the corresponding error images. From Table 3.1, we see that the PSNR performance of BAd-VAMP (which does not know \mathbf{b}) is nearly as good as genie-calibrated VAMP, and 1.7 dB better than genie-calibrated RED. Furthermore, the PSNR performance of BAd-VAMP is more than 6.6 dB better than un-calibrated VAMP and 7 dB better than un-calibrated RED. The uncalibrated VAMP, RED, and TV recoveries in Fig. 3.5 are plagued by either streaking artifacts and/or loss of detail (e.g., note the disappearance of the small white dots in uncalibrated TV). But the BAdVAMP image recovery in Fig. 3.5 shows no streaking artifacts and a high level of detail. Likewise, Fig. 3.6 shows that TV has trouble correctly recovering the white outer ellipse, RED has trouble in the interior region, but BAd-VAMP does well throughout.

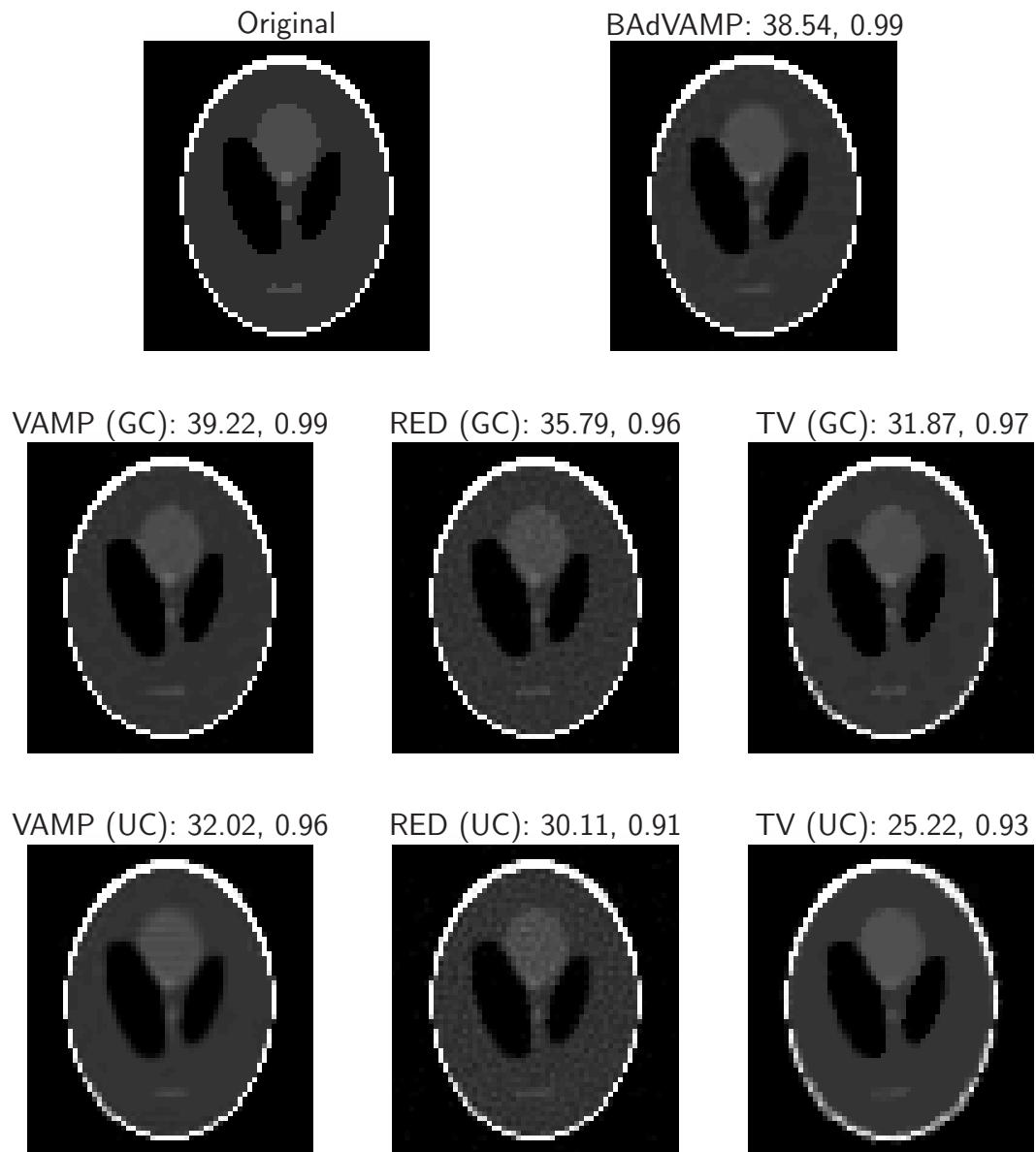


Figure 3.5: Calibration in tomography: Reconstruction PSNR in dB and SSIM of 64×64 Shepp-Logan phantom from 25 equally spaced tomographic projections. In the genie-calibrated (GC) case, $\hat{\mathbf{A}} = \mathbf{A}(\mathbf{b})$, while in the un-calibrated (UC) case, $\hat{\mathbf{A}} = \mathbf{A}(\mathbf{1})$.

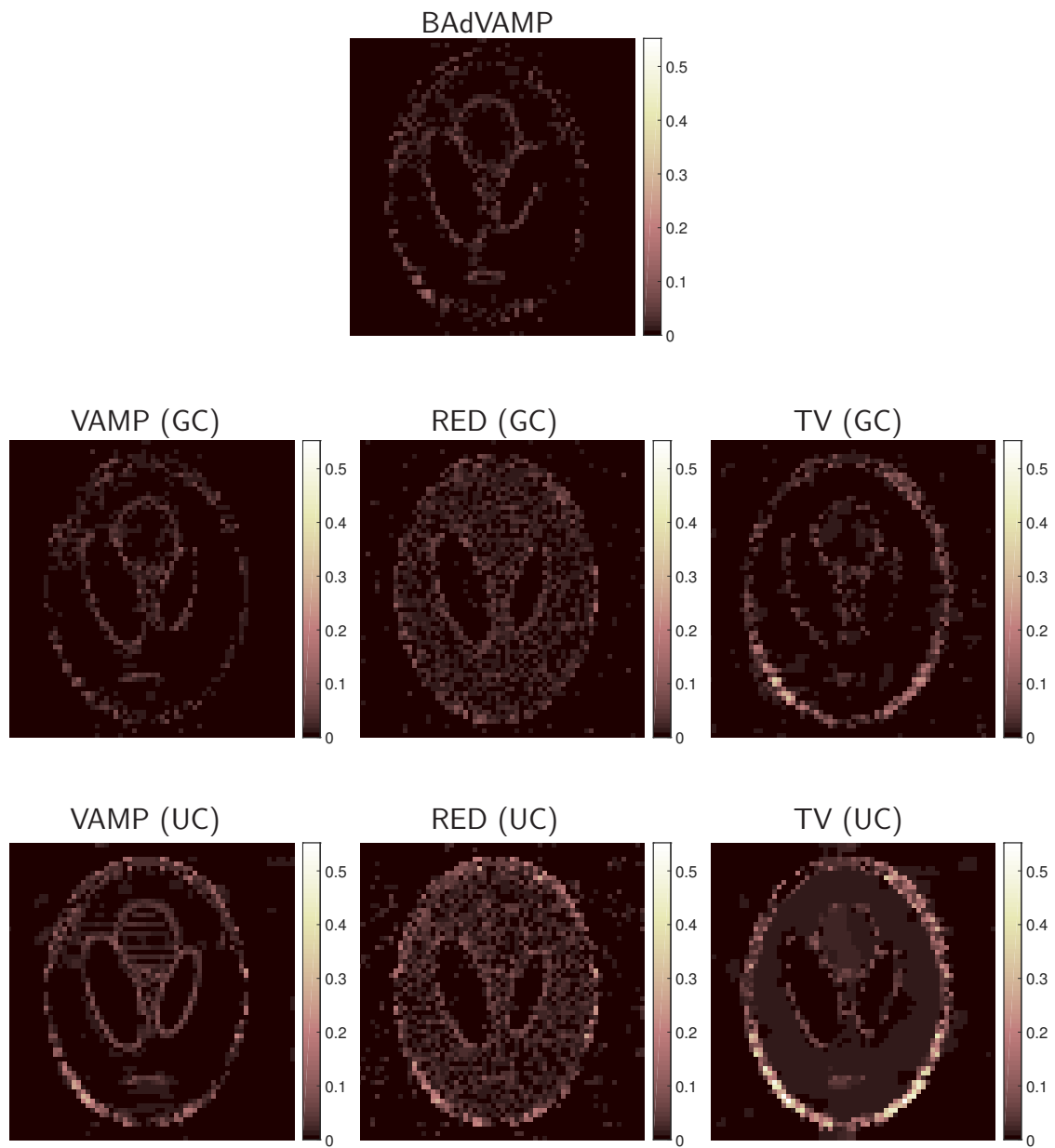


Figure 3.6: Calibration in tomography: Error images for the reconstructions shown in Fig. 3.5.

3.9.4 Noiseless Dictionary Learning

In dictionary learning (DL) [81], the goal is to find a dictionary matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ and a sparse matrix $\mathbf{X} \in \mathbb{R}^{N \times L}$ such that a given matrix $\mathbf{Y} \in \mathbb{R}^{M \times L}$ can be approximately factored as $\mathbf{Y} \approx \mathbf{A}\mathbf{X}$. In this section, we test the proposed BAd-VAMP algorithm for DL by generating $\mathbf{Y} = \mathbf{A}\mathbf{X}$ such that \mathbf{X} has K -sparse columns, and measuring the NMSE on the resulting estimates of \mathbf{A} and \mathbf{X} .

We consider two cases: i) where the true \mathbf{A} is structured as $\mathbf{A} = \sum_{i=1}^Q b_i \mathbf{A}_i$ with known $\{\mathbf{A}_i\}_{i=1}^Q$ (recall (3.62)), and ii) where the true \mathbf{A} is unstructured (recall (3.71)). In either case, the pair (\mathbf{A}, \mathbf{X}) is recoverable only up to an ambiguity: a scalar ambiguity in the structured case and a generalized permutation ambiguity in the unstructured case. Thus, when measuring reconstruction quality, we consider

$$\text{NMSE}(\hat{\mathbf{A}}) \triangleq \min_{\lambda \in \mathbb{R}} \frac{\|\mathbf{A} - \lambda \hat{\mathbf{A}}\|_F^2}{\|\mathbf{A}\|_F^2} \quad (3.89)$$

in the structured case and

$$\text{NMSE}(\hat{\mathbf{A}}) \triangleq \min_{\mathbf{P} \in \mathcal{P}} \frac{\|\mathbf{A} - \hat{\mathbf{A}}\mathbf{P}\|_F^2}{\|\mathbf{A}\|_F^2} \quad (3.90)$$

in the unstructured case, where \mathcal{P} denotes the set of generalized permutation matrices.⁶ For our experiments, we drew the coefficients of $\{\mathbf{A}_i\}_{i=1}^Q$ and \mathbf{b} as i.i.d. $\mathcal{N}(0, 1)$ with $Q = N$ in the structured case, and we drew the coefficients of \mathbf{A} as i.i.d. $\mathcal{N}(0, 1)$ in the unstructured case.

In our first experiment, we fixed the sparsity rate at $K/N = 0.2$ and we varied both the dictionary dimension N and the training length L . The top-right panel of Fig. 3.7 suggests that, as the dimension N grows, a *fixed* training length L is sufficient

⁶If \mathbf{P} is a generalized permutation matrix then $\mathbf{P} = \mathbf{\Pi}\mathbf{D}$, where $\mathbf{\Pi}$ is a permutation matrix and \mathbf{D} is a diagonal matrix.

to successfully recover \mathbf{A} in the structured case with $Q = N$. By “successfully recover,” we mean that $\text{NMSE}(\widehat{\mathbf{A}}) \leq -50$ dB. Note that this latter prescription for L is consistent with the theoretical analysis in [87]. In the unstructured case, the bottom panel of Figure 3.7 shows the median $\text{NMSE}(\widehat{\mathbf{A}})$ versus N when $L = 6N \ln N$. Together, the top-left and bottom panels of Fig. 3.7 suggest that a training length of $L = O(N \ln N)$ suffices to successfully recover \mathbf{A} .

In our second experiment, we focused on the unstructured case, fixed the training length at $L = 5N \ln N$, and varied both the dictionary dimension N and the sparsity K in the columns of \mathbf{X} . Figure 3.8 shows that BAd-VAMP performed similarly to EM-BiGAMP [68] for all but very small N , and much better than K-SVD [2] and SPAMS [54]. The advantage of BAd-VAMP over EM-BiGAMP for DL will be illustrated in the sequel.

3.9.5 Noisy, Ill-Conditioned Dictionary Learning

In this section, we show the robustness of BAd-VAMP over EM-BiGAMP [67] when learning ill-conditioned dictionaries from noisy measurements. To do so, we generated the measurements as $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{W}$ and tested the algorithms in recovering \mathbf{A} and \mathbf{X} (up to appropriate ambiguities). The elements of \mathbf{W} were drawn i.i.d. $\mathcal{N}(0, 1/\gamma_w)$ with γ_w chosen to achieve $\text{SNR} \triangleq \mathbb{E}[\|\mathbf{A}\mathbf{X}\|_F^2] / \mathbb{E}[\|\mathbf{W}\|_F^2] = 40$ dB. The true dictionary was generated as $\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^T$, where \mathbf{U} and \mathbf{V} were drawn uniformly over the group of orthogonal matrices, and where the singular values in \mathbf{s} were chosen so that $s_i/s_{i-1} = \rho \forall i$. The values of s_0 and ρ were selected to obtain a desired condition number $\kappa(\mathbf{A})$ while also ensuring $\|\mathbf{A}\|_F^2 = N$.

Figure 3.9 reports median $\text{NMSE}(\widehat{\mathbf{A}})$ and $\text{NMSE}(\widehat{\mathbf{X}})$ versus condition number $\kappa(\mathbf{A})$ for the recovery of $\mathbf{A} \in \mathbb{R}^{N \times N}$ and K -sparse $\mathbf{X} \in \mathbb{R}^{N \times L}$ from noisy measurements \mathbf{Y} . For this figure, we used $K = 13$, $N = 64$, $L = 5N \ln N$, the unstructured definition of $\text{NMSE}(\widehat{\mathbf{A}})$ from (3.90), and a similar definition for $\text{NMSE}(\widehat{\mathbf{X}})$. In addition to showing the performance of BAd-VAMP and EM-PBiGAMP, the figure shows the performance of the known- \mathbf{X} oracle for the estimation of \mathbf{A} , as well as the known- \mathbf{A} and known-support oracle for the estimation of \mathbf{X} . Figure 3.9 shows that EM-BiGAMP gave near-oracle NMSE for $\kappa(\mathbf{A}) \leq 40$, but its performance degraded significantly for larger $\kappa(\mathbf{A})$. In contrast, BAd-VAMP gave near-oracle NMSE for $\kappa(\mathbf{A}) \leq 110$, which suggests increased robustness to ill-conditioned dictionaries \mathbf{A} .

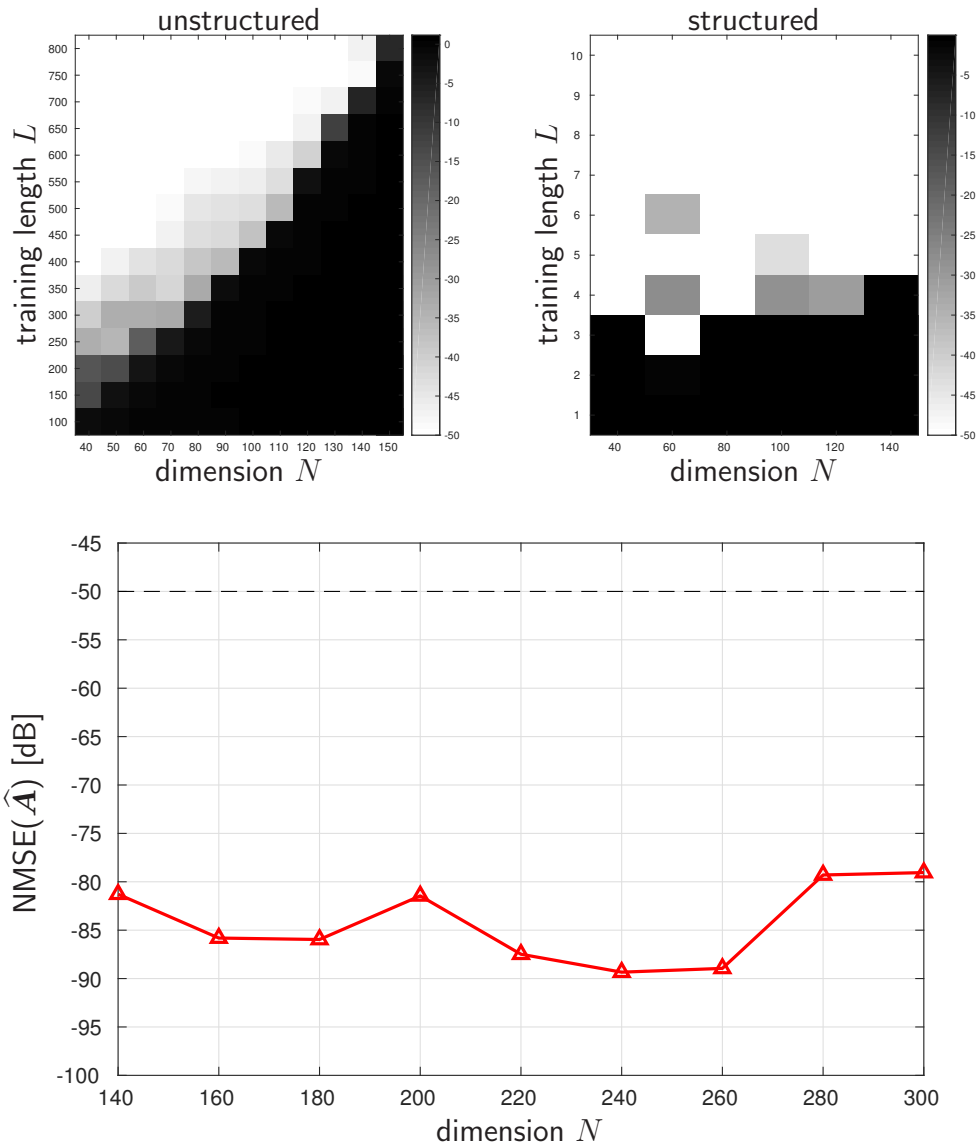


Figure 3.7: BAd-VAMP dictionary learning with $N \times N$ dictionary \mathbf{A} and $N \times L$ code matrix \mathbf{X} with sparsity rate $K/N = 0.2$. Top left: success-rate for unstructured dictionary versus N and L . Top right: success-rate for structured dictionary with $Q = N$ free parameters. Bottom: Median $\text{NMSE}(\hat{\mathbf{A}})$ over 20 trials versus N for an unstructured dictionary with $L = 6N \ln N$.

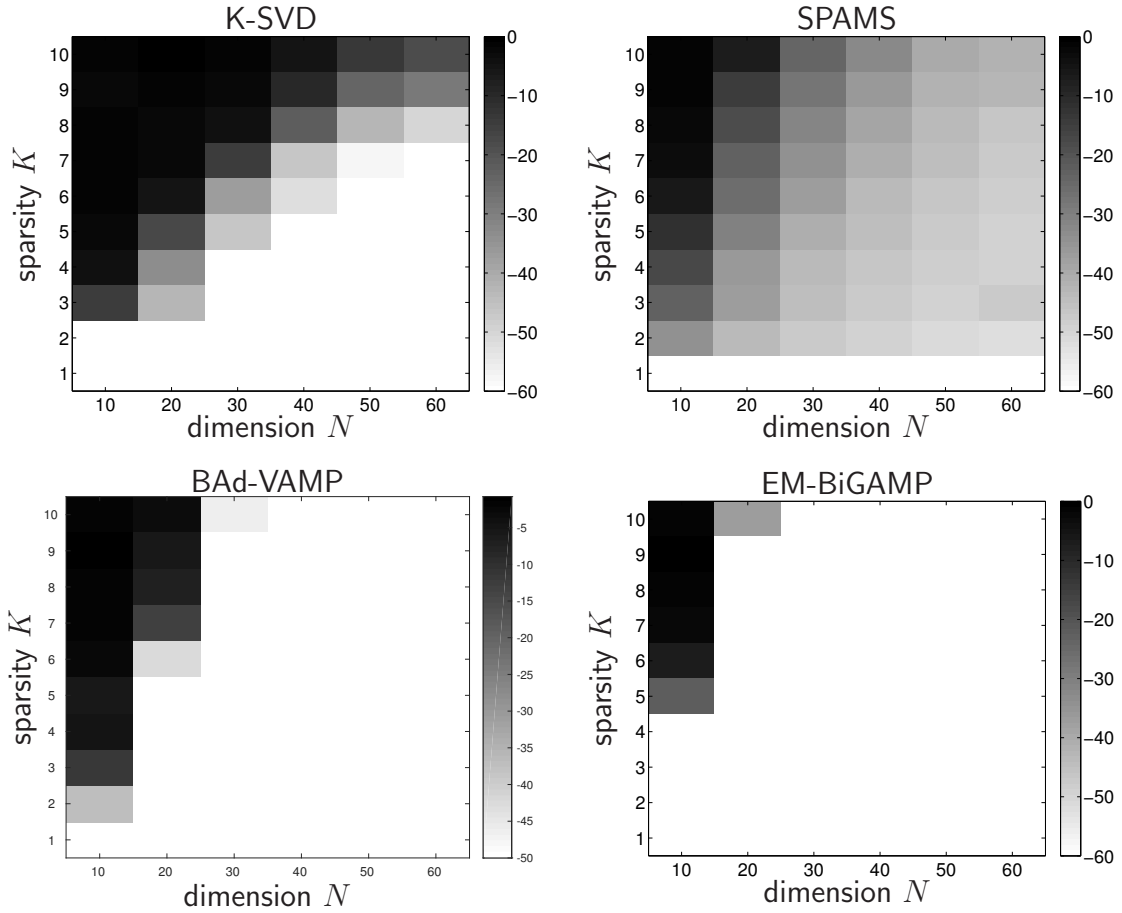


Figure 3.8: Noiseless dictionary learning: Median $\text{NMSE}(\hat{\mathbf{A}})$ in dB (over 20 trials) versus dictionary dimension N and sparsity K with unstructured dictionary $\mathbf{A} \in \mathbb{R}^{N \times N}$ and training length $L = 5N \ln N$.

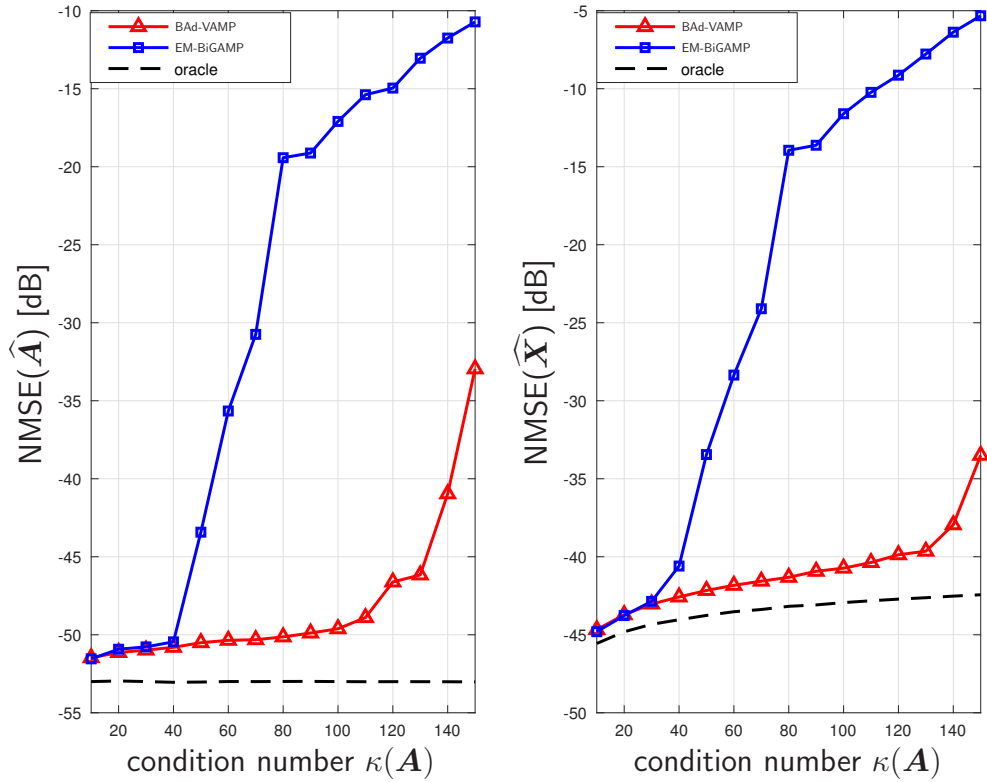


Figure 3.9: Noisy dictionary learning: Median NMSE in dB (over 50 trials) versus condition number $\kappa(\mathbf{A})$ for unstructured $\mathbf{A} \in \mathbb{R}^{N \times N}$ with $N = 64$, sparsity $K = 13$, training length $L = 5N \ln N$, and SNR = 40 dB.

Chapter 4: Image Reconstruction in Magnetic Resonance Imaging

4.1 Introduction

Magnetic Resonance Imaging (MRI) is a non-invasive technique used in radiology to generate internal images of the body. It uses a combination of strong static magnetic field, gradient magnetic field, and a sequence of RF pulses to capture measurements. The underlying physics of MRI is a complex quantum mechanical phenomenon. In classical physics, a static magnetic field magnetizes the atomic nucleus which causes all the nuclei to spin coherently at an angular frequency proportional to the external magnetic field. When an RF pulse of the same frequency is applied, the nuclei get excited to a higher energy state and upon relaxation emit signals measured by receiver coils in the MRI scanner. The gradient magnetic fields are used to spatially encode the signal. Since our body has an abundance of Hydrogen atoms, MRI mostly maps the location of water molecules and fat tissues in the body.

Unlike CT, which uses high energy ionizing X-rays, MRI is considered to be much safer. MRI provides better contrast between fat, water, muscle and other soft tissue structures, making it the most preferred imaging method for knee, brain, etc.

There are several challenges in MRI. Data acquisition is a very slow process, e.g., an MRI scan of the pelvic area can take more than 60 minutes and the patient has to remain still in a long narrow tube for the entire duration. Any small motion can cause undesirable artifacts in the image. For this reason, there has been a lot of research in accelerating MRI. The MRI scanner measures the image in the spatial Fourier domain, also known as the “k-space”. In the ideal case, if we were to measure the entire k-space, then we could reconstruct the image exactly by taking the inverse Fourier transform of the measurements. However, in that case, the data acquisition would be extremely slow. In accelerated MRI, we acquire a relatively small number of samples in k-space and reconstruct the image from these few measurements. To increase the number of measurements without compromising the acquisition time, a popular technique is parallel MRI, in which k-space data is simultaneously acquired using multiple receiver coils. Each receiver coil has a different sensitivity map, which acts as a spatial encoding.

In this chapter, we will develop algorithms to reconstruct images in accelerated MRI using the Approximate Message Passing framework, for both single-coil and multicoil MRI. We will assume that the coil sensitivity maps are perfectly known. In practice, they can be estimated with techniques like ESPIRIT [94].

4.2 Measurement Model

With known coil maps, we can formulate MRI image reconstruction as a linear inverse problem. The measurements are noisy, subsampled, k-space data. In single-coil MRI, the measurements acquired by the coil can be mathematically modeled as

$$\mathbf{y} = \underbrace{\mathbf{S}\mathbf{F}}_{\mathbf{A}} \mathbf{x} + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w) \quad (4.1a)$$

$$\Rightarrow \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}, \quad (4.1b)$$

where \mathbf{F} is the $N \times N$ 2D or 3D DFT matrix, $\mathbf{S} \in \mathbb{R}^{M \times N}$ is a sampling mask constructed by selecting M random rows from the $N \times N$ identity matrix, and $\mathbf{w} \in \mathbb{C}^M$ is AWGN measurement noise. The ratio $R \triangleq N/M$ is called the *acceleration rate*.

In parallel MRI, we have $K > 1$ receiver coils acquiring subsampled k-space data points simultaneously. Each coil applies a spatial encoding represented by a diagonal coil sensitivity matrix \mathbf{C}_i and all the coils have the same sampling mask \mathbf{S} . The measurement \mathbf{y}_i from the i^{th} coil can be modeled as

$$\mathbf{y}_i = \mathbf{S}\mathbf{F}\mathbf{C}_i\mathbf{x} + \mathbf{w}_i \quad (4.2a)$$

$$\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w). \quad (4.2b)$$

We then stack all the coil measurements $\{\mathbf{y}_i\}_{i=1}^K$ into a single column vector \mathbf{y} which yields

$$\underbrace{\begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_K \end{bmatrix}}_{\mathbf{y}} = \underbrace{\begin{bmatrix} \mathbf{S}\mathbf{F}\mathbf{C}_1 \\ \vdots \\ \mathbf{S}\mathbf{F}\mathbf{C}_K \end{bmatrix}}_{\mathbf{A}} \mathbf{x} + \underbrace{\begin{bmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_K \end{bmatrix}}_{\mathbf{w}} \quad (4.3a)$$

$$\Rightarrow \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w} \in \mathbb{C}^{MK}. \quad (4.3b)$$

The sampling mask in k-space can take several forms; random uniformly chosen points, non-uniform variable density [57], or random vertical or horizontal lines in k-space with a fully sampled central region. The uniform and variable density sampling masks make \mathbf{A} in (4.3b) “incoherent”, which is helpful for compressive sensing

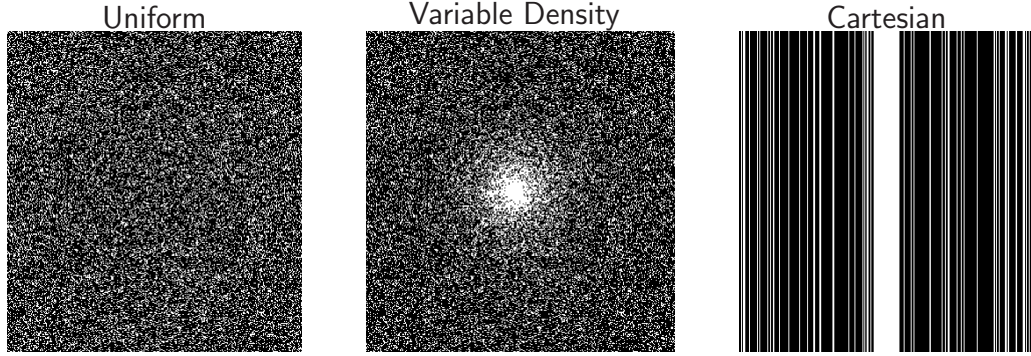


Figure 4.1: Example sampling masks for acceleration rate $R = 4$

based methods. Taking vertical lines in k-space sampling is called “Cartesian sampling” and is the most dominant method in clinical practice, mainly due its robustness to various system defects.

4.3 Signal Denoiser

From (4.1) and (4.3b), we see that image reconstruction in MRI is a linear inverse problem. Figure 4.2 shows that for high acceleration rates, the matrix \mathbf{A} is ill-conditioned, so using prior knowledge of MR images will lead to better reconstruction performance. If $p_{\mathbf{x}}$ is a prior density over \mathbf{x} , the *maximum a posteriori* (MAP) estimate of \mathbf{x} is

$$\hat{\mathbf{x}}_{\text{MAP}} \triangleq \arg \min_{\mathbf{x}} -\ln p_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) \quad (4.4a)$$

$$= \arg \min_{\mathbf{x}} -\ln p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) - \ln p_{\mathbf{x}}(\mathbf{x}). \quad (4.4b)$$

The likelihood function from (4.3b) is

$$p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \mathbf{I}/\gamma_w), \quad (4.5)$$

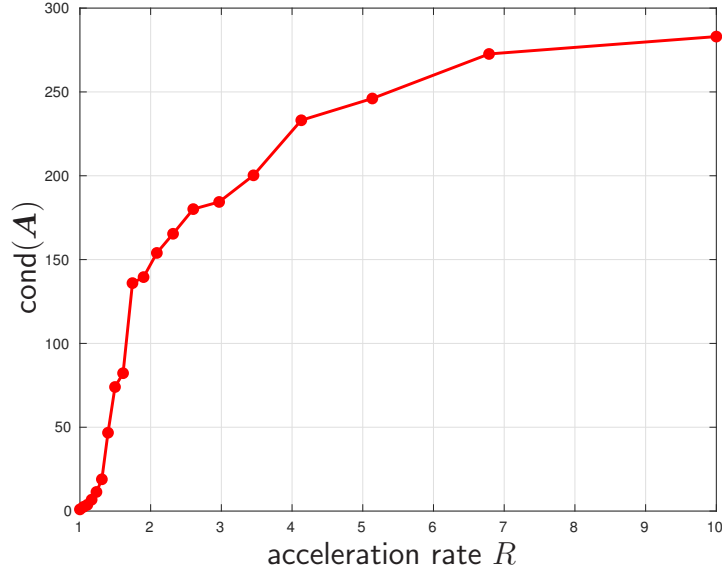


Figure 4.2: Condition number of \mathbf{A} vs acceleration rate R in multicoil MRI for Cartesian sampling and 4 receiver coils.

in which case Equation (4.4b) simplifies to

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{\gamma_w}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \lambda \rho(\mathbf{x}), \quad (4.6)$$

where $\rho(\mathbf{x}) \triangleq -\frac{1}{\lambda} \ln p_{\mathbf{x}}(\mathbf{x})$ and $\lambda > 0$. In CS based reconstruction methods [53], sparsity of transform coefficients for some transform Φ^H , e.g. a wavelet transform, is used as the prior, in which case $\rho(\cdot)$ typically takes the form

$$\rho(\mathbf{x}) = \|\Phi^H \mathbf{x}\|_1. \quad (4.7)$$

The ℓ_1 norm yields a convex regularization that encourages solutions with sparse transform coefficients.

As we increase the acceleration rate, we need more powerful priors to accurately reconstruct the MR image. Designing a prior distribution $p_{\mathbf{x}}$ for images is difficult in practice. However, most iterative optimization algorithms don't require an explicit

knowledge of $p_{\mathbf{x}}$. Rather, they use a signal denoiser, which implicitly uses some prior $p_{\mathbf{x}}$ and which is much easier to design in practice. A Gaussian denoiser is defined as a function that recovers \mathbf{x} from its AWGN corrupted version \mathbf{r} :

$$\mathbf{r} = \mathbf{x} + \mathbf{q}, \quad \mathbf{q} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma). \quad (4.8)$$

There have been many papers on image denoising, e.g., [15, 26, 52, 92, 104], some of which solve a MAP optimization problem of the form (4.6) with $\mathbf{A} = \mathbf{I}$, whereas some are algorithmic [26, 52, 92] or deep learning based [52, 104]. In the next section, we will describe in detail how a denoiser is used to solve linear inverse problems⁷.

4.4 Plug-and-Play Algorithms

To solve convex problems of the form (4.6), a commonly used algorithm is alternating direction method of multipliers (ADMM) [14] which converts (4.6) into a constrained optimization problem using the variable splitting technique,

$$\min_{\mathbf{x}, \mathbf{v}} \frac{\gamma_w}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \rho(\mathbf{v}) \quad (4.9a)$$

$$\text{subject to } \mathbf{x} = \mathbf{v}. \quad (4.9b)$$

and solves it using the augmented Lagrangian

$$\min_{\mathbf{x}, \mathbf{v}} \max_{\mathbf{u}} \left\{ \frac{\gamma_w}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \lambda \rho(\mathbf{v}) + \text{Re}\{\mathbf{u}^H (\mathbf{x} - \mathbf{v})\} + \frac{c}{2} \|\mathbf{x} - \mathbf{v}\|^2 \right\}, \quad (4.10)$$

⁷We saw an example of this earlier in Sec. 2.5

where \mathbf{u} is a Lagrange multiplier and $c > 0$ is a tuning parameter that does not affect the fixed point. ADMM solves the following sequence of sub-problems at iteration t ,

$$\mathbf{x}^{t+1} = \arg \min_{\mathbf{x}} \frac{\gamma_w}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \frac{c}{2} \|\mathbf{x} - \mathbf{v}^t + \tilde{\mathbf{u}}^t\|^2 \quad (4.11a)$$

$$\mathbf{v}^{t+1} = \arg \min_{\mathbf{v}} \lambda\rho(\mathbf{v}) + \frac{c}{2} \|\mathbf{v} - \mathbf{x}^{t+1} - \tilde{\mathbf{u}}^t\|^2 \quad (4.11b)$$

$$\tilde{\mathbf{u}}^{t+1} = \tilde{\mathbf{u}}^t + (\mathbf{x}^{t+1} - \mathbf{v}^{t+1}). \quad (4.11c)$$

where $\tilde{\mathbf{u}}^t \triangleq (1/c)\mathbf{u}^t$. Under mild conditions on convex ρ , it can be shown that ADMM converges to the solution of (4.6). Using the definition of the proximal operator

$$\text{prox}_{\phi}(\mathbf{r}) \triangleq \arg \min_{\mathbf{x}} \phi(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{r}\|^2, \quad (4.12)$$

we can rewrite the ADMM iterations in (4.11) as

$$\mathbf{x}^{t+1} = \text{prox}_{\frac{\gamma_w}{c}f}(\mathbf{v}^t - \tilde{\mathbf{u}}^t) \quad (4.13a)$$

$$\mathbf{v}^{t+1} = \text{prox}_{\frac{\lambda}{c}\rho}(\mathbf{x}^{t+1} + \tilde{\mathbf{u}}^t) \quad (4.13b)$$

$$\tilde{\mathbf{u}}^{t+1} = \tilde{\mathbf{u}}^t + (\mathbf{x}^{t+1} - \mathbf{v}^{t+1}), \quad (4.13c)$$

where $f(\mathbf{x}) \triangleq \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$ is the data fidelity term.

We can interpret the proximal operator $\text{prox}_{\frac{\lambda}{c}\rho}$ as a Gaussian denoiser. For example, it is the MAP estimator of \mathbf{x} from the AWGN corrupted measurement $\mathbf{r} = \mathbf{x} + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma)$ with prior density $p_{\mathbf{x}}(\mathbf{x}) \propto e^{-\lambda\rho(\mathbf{x})}$. Using this interpretation, Bouman et. al. [95] suggested to replace (4.13b) with any suitable image denoiser, $\mathbf{g}(\mathbf{r}, \gamma)$. In this way, ADMM can be used with denoisers that are neither optimization based nor have any statistical interpretation, e.g., block-matching and 3D filtering (BM3D) [26], non-local means (NLM) [15], denoising convolutional neural network (DnCNN) [104], etc.

The variable splitting technique in ADMM decouples the prior image knowledge from data consistency. It solves two subproblems at every iteration; (4.13a) enforcing data consistency and (4.13b) enforcing the prior knowledge. Many other iterative solvers for linear inverse problems also decouple the two stages using some variable splitting method. Though the authors in [95] originally used ADMM, their approach can be easily generalized to other iterative algorithms, which led to the development of so called *plug-and-play* (PnP) algorithms. PnP-FISTA [44], RED [76, 79], PnP-PDS [62] are examples of some other popular PnP algorithms.

Another advantage of any PnP algorithm is that the denoiser is not tied to any particular choice of \mathbf{A} . This is very beneficial in computational imaging, where the forward operator \mathbf{A} depends on the application. We use this PnP approach to develop our reconstruction algorithm for MRI using the AMP framework.

4.5 Our Approach

In most PnP algorithms, the denoiser input variance (often adjustable) is fixed over all iterations. But the denoiser variance affects their fixed point and convergence speed [25]. AMP based algorithms don't have this issue, because the denoiser input noise variance is automatically estimated by the algorithm at every iteration, leading to fast convergence and accurate solutions with large random \mathbf{A} . However, AMP based algorithms are not designed to use the structured \mathbf{A} typically used in MRI. Thus, their noise-variance estimates may not be accurate and the algorithm may even diverge.

Some existing AMP based algorithms for MRI are Variable Density AMP (VDAMP) [57], which uses soft-thresholding of the wavelet coefficients for denoising, and BM3D-AMP-MRI [33] which uses a modified version of the BM3D [26] denoiser to handle complex-valued images. VDAMP neither supports multicoil MRI nor Cartesian sampling mask and BM3D-AMP-MRI is for single-coil MRI. To the best of our knowledge, we are the first to propose an AMP based algorithm using a generic denoiser for image reconstruction in multicoil MRI using generic k-space sampling masks.

4.5.1 AMP for MRI

In this section, we discuss problems with the conventional AMP algorithm [32] for MRI and propose steps to rectify them. AMP recovers $\mathbf{x} \in \mathbb{C}^N$ from measurements $\mathbf{y} \in \mathbb{C}^M$ in (4.6) using the following iterations

$$\mathbf{v}^t = \mathbf{y} - \mathbf{A}\mathbf{x}^t + \frac{N}{M}\mathbf{v}^{t-1} \langle \mathbf{g}'^t(\mathbf{r}^{t-1}) \rangle \quad (4.14a)$$

$$\mathbf{r}^t = \mathbf{x}^t + \beta \mathbf{A}^H \mathbf{v}^t \quad (4.14b)$$

$$\mathbf{x}^{t+1} = \mathbf{g}^t(\mathbf{r}^t) \quad (4.14c)$$

where $\mathbf{v}^{-1} = \mathbf{0}$ and for i.i.d. Gaussian zero-mean \mathbf{A} the scalar β is set to $\frac{N}{\|\mathbf{A}\|_F^2}$. We later describe the effect of β on AMP. The denoiser \mathbf{g}^t may depend on iteration t , and its divergence $\langle \mathbf{g}'^t(\mathbf{r}) \rangle$ is the average of diagonal elements of the Jacobian of \mathbf{g}^t , that can be approximated using the Monte-Carlo method [70]

$$\langle \mathbf{g}'^t(\mathbf{r}) \rangle \approx \frac{1}{N\epsilon} \mathbf{q}^H [\mathbf{g}^t(\mathbf{r} + \epsilon \mathbf{q}) - \mathbf{g}^t(\mathbf{r})], \quad (4.15)$$

using small $\epsilon > 0$ and random \mathbf{q} such that $\mathbb{E}[\mathbf{q}\mathbf{q}^H] = \mathbf{I}$.

The rightmost term in (4.14a) is known as the *Onsager correction term* from statistical physics, which ensures that for zero-mean i.i.d. Gaussian \mathbf{A} , as the dimensions M, N grow to infinity (fixed M/N), the denoiser input error satisfies the AWGN property,

$$\mathbf{r}^t = \mathbf{x}_0 + \mathbf{e}^t, \quad \mathbf{e}^t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma^t). \quad (4.16)$$

Additionally, the errors $\{\mathbf{e}^t\}$ are independent across iterations t . See [9] for details.

The version of AMP in (4.14) is usually seen in the literature and it expects $\|\mathbf{A}\|_F^2 \approx N$. When it is not the case, the scalar β in (4.14) compensates for the mismatch. To understand how this happens, let \mathbf{A} be drawn from i.i.d. $\mathcal{N}(0, \sigma_a^2)$ and \mathbf{x}_0 be the true vector, and consider the following error terms in AMP:

$$\mathbf{e}^t \triangleq \mathbf{r}^t - \mathbf{x}_0 \quad (\text{denoiser input error}) \quad (4.17a)$$

$$\boldsymbol{\epsilon}^t \triangleq \mathbf{x}^t - \mathbf{x}_0 \quad (\text{denoiser output error}), \quad (4.17b)$$

which evolve according to

$$\boldsymbol{\epsilon}^{t+1} = \mathbf{g}^t(\mathbf{x}_0 + \mathbf{e}^t) - \mathbf{x}_0 \quad (4.18a)$$

$$\mathbf{e}^{t+1} = (\mathbf{I} - \beta \mathbf{A}^H \mathbf{A}) \boldsymbol{\epsilon}^{t+1} + \beta \frac{N}{M} \mathbf{A}^H \mathbf{v}^t \langle \mathbf{g}'^t(\mathbf{x}_0 + \mathbf{e}^t) \rangle + \beta \mathbf{A}^H \mathbf{w}. \quad (4.18b)$$

In the first iteration, $t = 0$, since we initialize using $\mathbf{v}^{-1} = 0$, we see that the denoiser input error is

$$\mathbf{e}^0 = (\mathbf{I} - \beta \mathbf{A}^H \mathbf{A}) \boldsymbol{\epsilon}^0 + \beta \mathbf{A}^H \mathbf{w}. \quad (4.19)$$

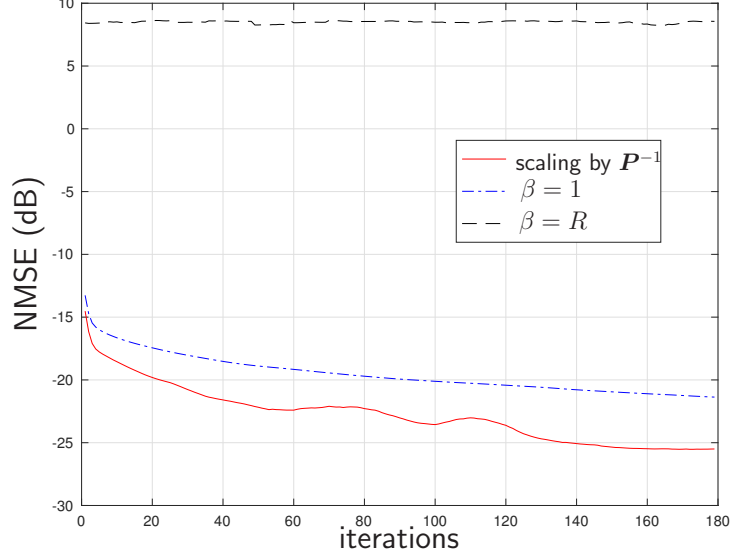


Figure 4.3: NMSE (dB) of AMP in MRI for different type of scaling in k-space

If we assume ϵ^0 is independent of $\{\mathbf{A}, \mathbf{w}\}$ and we take the expectation conditioned on ϵ^0 , then we get

$$\mathbb{E}[e^0 | \epsilon^0] = \mathbb{E}[(\mathbf{I} - \beta \mathbf{A}^H \mathbf{A}) \epsilon^0 | \epsilon^0] + \beta \underbrace{\mathbb{E}[\mathbf{A}^H \mathbf{w}]}_{\mathbf{0}} \quad (4.20a)$$

$$\stackrel{(a)}{=} \mathbb{E} \left[\left(\mathbf{I} - \frac{1}{\sigma_a^2 M} \mathbf{A}^H \mathbf{A} \right) \epsilon^0 \right] \quad (4.20b)$$

$$= \left(\mathbf{I} - \frac{1}{\sigma_a^2 M} \sigma_a^2 M \mathbf{I} \right) \epsilon^0 \quad (4.20c)$$

$$= \mathbf{0}, \quad (4.20d)$$

where in (a) we used $\beta = \frac{N}{\|\mathbf{A}\|_F^2} = \frac{1}{\sigma_a^2 M}$. Equation (4.20) shows that this choice of the scalar β is necessary for the denoiser input to be an unbiased estimate of \mathbf{x}_0 .

In MRI, \mathbf{A} is not i.i.d. Gaussian, so if we run AMP with $\beta = \frac{N}{\|\mathbf{A}\|_F^2} = R \triangleq \frac{N}{M}$, it doesn't work well and often diverges. An existing method BM3D-AMP-MRI [33] sets $\beta = 1$, which empirically helps to stabilize AMP but its final solution is not quite accurate. For single-coil knee MRI, Cartesian sampling of $R = 4$ and SNR=40 dB,

Fig. 4.3 shows the NMSE of AMP for $\beta = 1$, $\beta = R$, and another type of scaling in k-space that we propose next.

Scaling in k-space

It is not clear how to achieve (4.16) in MRI, but if one scales the k-space measurements in a certain way, one can still have the unbiased property satisfied. We are going to assume that the coil sensitivities satisfy:

$$\sum_{i=1}^K \mathbf{C}_i^H \mathbf{C}_i = \mathbf{I}. \quad (4.21)$$

Coil sensitivities estimated by ESPIRIT [94] guarantees the above scaling.

Consider generating the sampling mask \mathbf{S} by randomly selecting rows from the $N \times N$ identity matrix such that the probability of selecting the i^{th} row is p_i , where $p_i > 0 \forall i$. Let $\mathbf{P} \in \mathbb{R}^{N \times N}$ be a diagonal matrix with the probabilities $\{p_i\}_{i=1}^N$ along its diagonal, i.e., $\mathbf{P} \triangleq \text{Diag}(\mathbf{p})$. In the conventional AMP iterations, instead of scaling the k-space residual \mathbf{v}^t in line (4.14b) by the scalar β , we propose to scale it by $\mathbf{S}\mathbf{P}^{-1}\mathbf{S}^T$. In other words, we replace line (4.14b) by

$$\mathbf{r}^t = \mathbf{x}^t + \mathbf{A}^H(\mathbf{S}\mathbf{P}^{-1}\mathbf{S}^T\mathbf{v}^t). \quad (4.22)$$

In this case, at iteration $t = 0$, the denoiser input error becomes

$$\mathbf{e}^0 = (\mathbf{I} - \mathbf{A}^H\mathbf{S}\mathbf{P}^{-1}\mathbf{S}^T\mathbf{A})\mathbf{\epsilon}^0 + \mathbf{A}^H\mathbf{S}\mathbf{P}^{-1}\mathbf{S}^T\mathbf{w}. \quad (4.23)$$

And the expected error is

$$\mathbb{E}[\mathbf{e}^0 | \boldsymbol{\epsilon}^0] = \mathbb{E} \left[\mathbf{I} - \mathbf{A}^H \mathbf{S} \mathbf{P}^{-1} \mathbf{S}^T \mathbf{A} \right] \boldsymbol{\epsilon}^0 + \mathbb{E} \left[\mathbf{A}^H \mathbf{S} \mathbf{P}^{-1} \mathbf{S}^T \mathbf{w} \right] \quad (4.24a)$$

$$\begin{aligned} &= \mathbb{E} \left[\mathbf{I} - \sum_{i=1}^K \mathbf{C}_i^H \mathbf{F}^H \mathbf{S}^T \mathbf{S} \mathbf{P}^{-1} \mathbf{S}^T \mathbf{S} \mathbf{F} \mathbf{C}_i \right] \boldsymbol{\epsilon}^0 \\ &\quad + \mathbb{E} \left[\mathbf{A}^H \mathbf{S} \mathbf{P}^{-1} \mathbf{S}^T \right] \underbrace{\mathbb{E}[\mathbf{w}]}_{\mathbf{0}} \end{aligned} \quad (4.24b)$$

$$= \left(\mathbf{I} - \sum_{i=1}^K \mathbf{C}_i^H \mathbf{F}^H \mathbb{E} \left[\mathbf{S}^T \mathbf{S} \mathbf{P}^{-1} \mathbf{S}^T \mathbf{S} \right] \mathbf{F} \mathbf{C}_i \right) \boldsymbol{\epsilon}^0 \quad (4.24c)$$

$$= \left(\mathbf{I} - \sum_{i=1}^K \mathbf{C}_i^H \mathbf{F}^H \mathbb{E} \left[\overline{\mathbf{S}} \mathbf{P}^{-1} \overline{\mathbf{S}} \right] \mathbf{F} \mathbf{C}_i \right) \boldsymbol{\epsilon}^0 \quad \text{where } \overline{\mathbf{S}} \triangleq \mathbf{S}^T \mathbf{S} \quad (4.24d)$$

$$= \left(\mathbf{I} - \sum_{i=1}^K \mathbf{C}_i^H \mathbf{F}^H \mathbb{E} \left[\overline{\mathbf{S}} \right] \mathbf{P}^{-1} \mathbf{F} \mathbf{C}_i \right) \boldsymbol{\epsilon}^0 \quad (4.24e)$$

$$\stackrel{(a)}{=} \left(\mathbf{I} - \sum_{i=1}^K \mathbf{C}_i^H \mathbf{F}^H \mathbf{P} \mathbf{P}^{-1} \mathbf{F} \mathbf{C}_i \right) \boldsymbol{\epsilon}^0 \quad (4.24f)$$

$$= \left(\mathbf{I} - \sum_{i=1}^K \mathbf{C}_i^H \mathbf{C}_i \right) \boldsymbol{\epsilon}^0 \quad (4.24g)$$

$$\stackrel{(b)}{=} \mathbf{0}, \quad (4.24h)$$

where (a) follows from $\mathbb{E}[\overline{\mathbf{S}}] = \mathbf{P}$ and (b) follows from the property (4.21). Our proposed algorithm, Debaised Denoising AMP (DD-AMP) is summarized in Algorithm 6.

Another way to derive Alg. 6 is to left multiply (4.1) by $\mathbf{P}^{-\frac{1}{2}} \mathbf{S}^T$, in which case we get the following equivalent measurement model,

$$\tilde{\mathbf{y}} = \tilde{\mathbf{A}} \mathbf{x} + \tilde{\mathbf{w}}, \quad (4.25)$$

Algorithm 6 DD-AMP (Debiased Denoising AMP) for MRI

```

1: define:
    $\mathbf{B} = \mathbf{S}\mathbf{P}^{-1}\mathbf{S}^\top$ 
2: initialize:
    $\mathbf{x}^0, \mathbf{v}^{-1} = \mathbf{0}$ 
3: for  $t = 0, \dots, T_{\max}$  do
4:    $\mathbf{v}^t = \mathbf{y} - \mathbf{A}\mathbf{x}^t + \frac{N}{M}\mathbf{v}^{t-1} \langle \mathbf{g}^t(\mathbf{r}^{t-1}) \rangle$ 
5:    $\mathbf{r}^t = \mathbf{x}^t + \mathbf{A}^\mathbf{H}\mathbf{B}\mathbf{v}^t$ 
6:    $\mathbf{x}^{t+1} = \mathbf{g}^t(\mathbf{r}^t)$ 
7: end for

```

where $\tilde{\mathbf{y}} \triangleq \mathbf{P}^{-\frac{1}{2}}\mathbf{S}^\top\mathbf{y}$, $\tilde{\mathbf{A}} \triangleq \mathbf{P}^{-\frac{1}{2}}\mathbf{S}^\top\mathbf{A}$ and $\tilde{\mathbf{w}} \triangleq \mathbf{P}^{-\frac{1}{2}}\mathbf{S}^\top\mathbf{w}$. Taking expectation over the random quantities $\{\mathbf{S}, \mathbf{w}\}$, the noise in (4.25) is additive white noise,

$$\mathbb{E}[\tilde{\mathbf{w}}\tilde{\mathbf{w}}^\mathbf{H}] = \mathbb{E}\left[\mathbf{P}^{-\frac{1}{2}}\mathbf{S}^\top\mathbf{w}\mathbf{w}^\mathbf{H}\mathbf{S}\mathbf{P}^{-\frac{1}{2}}\right] \quad (4.26a)$$

$$= \mathbb{E}\left[\mathbf{P}^{-\frac{1}{2}}\mathbf{S}^\top\mathbb{E}[\mathbf{w}\mathbf{w}^\mathbf{H}]\mathbf{S}\mathbf{P}^{-\frac{1}{2}}\right] \quad (4.26b)$$

$$= \frac{1}{\gamma_w}\mathbb{E}\left[\mathbf{P}^{-\frac{1}{2}}\mathbf{S}^\top\mathbf{S}\mathbf{P}^{-\frac{1}{2}}\right] \quad (4.26c)$$

$$= \frac{1}{\gamma_w}\mathbb{E}\left[\mathbf{P}^{-\frac{1}{2}}\overline{\mathbf{S}}\mathbf{P}^{-\frac{1}{2}}\right], \quad \text{where } \overline{\mathbf{S}} \triangleq \mathbf{S}^\top\mathbf{S} \quad (4.26d)$$

$$= \frac{1}{\gamma_w}\mathbb{E}\left[\overline{\mathbf{S}}\mathbf{P}^{-1}\right] \quad (4.26e)$$

$$= \mathbf{I}/\gamma_w. \quad (4.26f)$$

If we run the traditional AMP algorithm (4.14) using $\beta = 1$ on the scaled measurement model (4.25), we get Algorithm 6 after a variable change.

VDAMP (Alg. 7) [57] also uses a scaling by \mathbf{P}^{-1} , but they do not apply it to the AMP or VAMP algorithms. They recover the wavelet coefficients of the image. In Alg. 7, $\mathbf{g}_{\text{soft-thresh}}$ is the soft-thresholding function, and \odot , \oslash , denote element-wise product and element-wise division respectively. The operator \mathcal{B} performs block-wise averaging, where one block is a wavelet band. In Alg. 7, line 5 is a linear estimator similar to AMP, whereas line 12 is an Onsager correction step similar to VAMP. Their

Algorithm 7 VDAMP from [57]

1: **define:**
 Ψ : wavelet transform
 $\mathbf{H} = |\mathbf{F}\Psi^H|^2$ (element-wise squared)
 $\mathcal{B}, \mathcal{B}^*$: block averaging operator, and its adjoint

2: **initialize:**
 $\tilde{\mathbf{r}}^0$

3: **for** $t = 0, \dots, T_{\max}$ **do**

4: $\mathbf{z}^t = \mathbf{y} - \mathbf{S}\mathbf{F}\Psi^H\tilde{\mathbf{r}}^t$

5: $\mathbf{r}^t = \tilde{\mathbf{r}}^t + \Psi\mathbf{F}^H\mathbf{P}^{-1}\mathbf{S}^T\mathbf{z}^t$

6: $\boldsymbol{\tau}^t = \mathbf{H}^T\mathbf{S}^T\mathbf{S}\mathbf{P}^{-1} [(\mathbf{P}^{-1} - \mathbf{I})\mathbf{S}^T|\mathbf{z}^t|^2 + \mathbf{1}/\gamma_w]$

7: $\bar{\boldsymbol{\tau}}^t = \mathcal{B}(\boldsymbol{\tau}^t)$

8: $\mathbf{w}^t = \mathbf{g}_{\text{soft-thresh}}(\mathbf{r}^t, \bar{\boldsymbol{\tau}}^t)$

9: $\bar{\boldsymbol{\alpha}}^t = \mathcal{B}(\text{diag}\{\partial\mathbf{g}'_{\text{soft-thresh}}(\mathbf{r}^t, \bar{\boldsymbol{\tau}}^t)/\partial\mathbf{r}\})$

10: $\boldsymbol{\alpha}^t = \mathcal{B}^*(\bar{\boldsymbol{\alpha}}^t)$

11: $\mathbf{c}^t = \mathbf{1} \oslash (\mathbf{1} - \boldsymbol{\alpha}^t)$

12: $\tilde{\mathbf{r}}^t = \mathbf{c}^t \odot (\mathbf{w}^t - \boldsymbol{\alpha}^t \odot \mathbf{r}^t)$

13: $\mathbf{x}^t = \Psi^H\mathbf{w}^t + \mathbf{F}^H\mathbf{S}^T(\mathbf{y} - \mathbf{S}\mathbf{F}\Psi^H\mathbf{w}^t)$

14: **end for**

algorithm is a hybrid of AMP and VAMP, and does not apply to the multicoil case.

We show in numerical experiments that our approach achieves superior results even for the single-coil case.

4.5.2 VAMP for MRI

VAMP has provable convergence guarantees when \mathbf{A} is right rotationally invariant (RRI) and the denoiser \mathbf{g}_1 is separable [73] and non-separable (chapter 2). For arbitrary \mathbf{A} and non-expansive denoiser \mathbf{g}_1 , VAMP is known to converge with appropriate damping [36]. In MRI, the RRI property does not hold, and empirically damping helps but we have no proof.

A damping scheme was proposed in [36] under which VAMP is guaranteed to converge with any non-expansive denoiser. They fixed the precisions in VAMP $(\gamma_1^t, \gamma_2^t, \eta_1^t, \eta_2^t) = (\gamma_1, \gamma_2, \gamma_1 + \gamma_2, \gamma_1 + \gamma_2)$ for arbitrary $\gamma_1, \gamma_2 > 0$, and then the damping factor ζ was calculated such that each iteration of damped VAMP was a contraction mapping. The expression for ζ is

$$\zeta = \frac{2}{1 + \max\left\{\frac{\gamma_1}{\gamma_2}, \frac{\gamma_2}{\gamma_1}\right\}}. \quad (4.27)$$

Though the damping factor was computed for fixed precisions, in practice ζ^t is calculated using the precisions at iteration t . Algorithm 8 shows the damped version of VAMP from [36]. However, this damping scheme doesn't work in some cases, especially in MRI with Cartesian sampling mask. VAMP can be unstable and diverge for high acceleration rates R and challenging sampling masks.

We propose a different damping scheme to stabilize VAMP and make it more robust. In Alg. 8, the precisions (γ_1, γ_2) control how much the denoiser \mathbf{g}_1 and linear estimator \mathbf{g}_2 act on their corresponding inputs in lines 10 and 5 respectively. VAMP usually becomes unstable when the precisions change rapidly between iterations. Thus it is important to damp the precisions too. The precisions γ_1 and γ_2 are calculated using the divergences α_1 and α_2 in lines 8 and 15 respectively. The divergences

(α_1, α_2) are computed using the Monte-Carlo method as shown in (4.15). Since \mathbf{g}_2 is linear, its divergence α_2 is quite accurate, whereas the divergence of \mathbf{g}_1 is not. To fix this, one could perform multiple Monte-Carlo draws to get a better estimate of α_1 or use a damping on α_1 . In practice, multiple Monte-Carlo draws are expensive, so we choose to damp α_1 .

Finally, it is important to note that the pseudo-measurements $(\mathbf{r}_1, \mathbf{r}_2)$ are amplitudes, the precisions (γ_1, γ_2) are inverse variances and the divergences (α_1, α_2) are variance scalings. This should be taken into account when damping, i.e.,

$$\mathbf{r}_2^{t+1} \leftarrow \zeta^t \mathbf{r}_2^{t+1} + (1 - \zeta^t) \mathbf{r}_2^t \quad (4.28a)$$

$$\gamma_2^{t+1} \leftarrow \left\{ \frac{\zeta^t}{\sqrt{\gamma_2^{t+1}}} + \frac{1 - \zeta^t}{\sqrt{\gamma_2^t}} \right\}^{-2} \quad (4.28b)$$

$$\alpha_1^t \leftarrow \left\{ \theta (\alpha_1^t)^{\frac{1}{2}} + (1 - \theta) (\alpha_1^{t-1})^{\frac{1}{2}} \right\}^2, \quad (4.28c)$$

where ζ^t is the damping factor computed similar to (4.27) using the precisions at iteration t , that can also be written in terms of the divergences as,

$$\zeta^t = \frac{2}{1 + \max \left\{ \frac{\gamma_1^t}{\gamma_2^t}, \frac{\gamma_2^t}{\gamma_1^t} \right\}} \quad (4.29a)$$

$$\stackrel{(a)}{=} \frac{2}{1 + \max \left\{ \frac{1 - \alpha_2^t}{\alpha_2^t}, \frac{1 - \alpha_1^t}{\alpha_1^t} \right\}} \quad (4.29b)$$

$$= \frac{2}{\max \left\{ \frac{1}{\alpha_2^t}, \frac{1}{\alpha_1^t} \right\}} \quad (4.29c)$$

$$= 2 \min \{ \alpha_1^t, \alpha_2^t \}, \quad (4.29d)$$

where (a) follows from lines 15 and 8 in Alg. 8.

Algorithm 8 Damped VAMP (dGEC) from [36]

```
1: require:  
   damping factor:  $\zeta \in (0, 1]$   
2: initialize:  
    $\mathbf{r}_2^0, \gamma_2^0$   
3: for  $t = 0, \dots, T_{\max}$  do  
4: //Linear Stage  
5:    $\mathbf{x}_2^t = \mathbf{g}_2(\mathbf{r}_2^t, \gamma_2^t)$   
6:    $\alpha_2^t = \langle \mathbf{g}'_2(\mathbf{r}_2^t, \gamma_2^t) \rangle$   
7:    $\mathbf{r}_1^t = (\mathbf{x}_2^t - \alpha_2^t \mathbf{r}_2^t) / (1 - \alpha_2^t)$   
8:    $\gamma_1^t = \gamma_2^t (1 - \alpha_2^t) / \alpha_2^t$   
9: //Denoising  
10:   $\mathbf{x}_1^t = \mathbf{g}_1(\mathbf{r}_1^t, \gamma_1^t)$   
11:   $\alpha_1^t = \langle \mathbf{g}'_1(\mathbf{r}_1^t, \gamma_1^t) \rangle$   
12:   $\zeta^t = 2 \min\{\alpha_1^t, \alpha_2^t\}$   
13:   $\bar{\mathbf{r}}_2^{t+1} = (\mathbf{x}_1^t - \alpha_1^t \mathbf{r}_1^t) / (1 - \alpha_1^t)$   
14:   $\mathbf{r}_2^{t+1} = \zeta^t \bar{\mathbf{r}}_2^{t+1} + (1 - \zeta^t) \mathbf{r}_2^t$   
15:   $\gamma_2^{t+1} = \gamma_1^t (1 - \alpha_1^t) / \alpha_1^t$   
16: end for
```

Using fast solvers in the Linear Stage

Recall that, in VAMP, the linear estimator \mathbf{g}_2 solves the following problem

$$\arg \min_{\mathbf{x}} \frac{\gamma_w}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \frac{\gamma_1^t}{2} \|\mathbf{x} - \mathbf{r}_1^t\|^2. \quad (4.30)$$

Its closed-form solution in (2.18) involves a matrix inverse that is expensive in high dimensions. So one can use a fast iterative solver such as LSQR [65] or conjugate gradient [7] to solve (4.30) approximately, as mentioned in [84]. Let $\mathbf{g}_{\text{linear-solver}}(\mathbf{r}_1^t, \gamma_1^t)$ denote the solution of a fast iterative solver for (4.30). Our proposed algorithm, Damped Denoising VAMP (DD-VAMP), is summarized in Alg. 9.

Algorithm 9 Damped Denoising-VAMP (DD-VAMP)

```
1: require:  
    damping factor:  $\theta = 0.1$   
     $\zeta_{\min} = 0.2, \zeta_{\max} = 0.3$   
2: initialize:  
     $\mathbf{r}_2^0, \gamma_2^0$   
3: for  $t = 0, \dots, T_{\max}$  do  
4: //Linear Stage  
5:    $\mathbf{x}_2^t = \mathbf{g}_{\text{linear-solver}}(\mathbf{r}_2^t, \gamma_2^t)$   
6:    $\alpha_2^t = \langle \mathbf{g}'_{\text{linear-solver}}(\mathbf{r}_2^t, \gamma_2^t) \rangle$   
7:    $\mathbf{r}_1^t = (\mathbf{x}_2^t - \alpha_2^t \mathbf{r}_2^t) / (1 - \alpha_2^t)$   
8:    $\gamma_1^t = \gamma_2^t (1 - \alpha_2^t) / \alpha_2^t$   
9: //Denoising  
10:   $\mathbf{x}_1^t = \mathbf{g}_1(\mathbf{r}_1^t, \gamma_1^t)$   
11:   $\bar{\alpha}_1^t = \langle \mathbf{g}'_1(\mathbf{r}_1^t, \gamma_1^t) \rangle$   
12:   $\alpha_1^t = \{\theta(\bar{\alpha}_1^t)^{\frac{1}{2}} + (1 - \theta)(\alpha_1^{t-1})^{\frac{1}{2}}\}^2$   
13:   $\bar{\zeta}^t = 2 \min\{\alpha_1^t, \alpha_2^t\}$   
14:   $\zeta^t = \max\{\zeta_{\min}, \min\{\zeta_{\max}, \bar{\zeta}^t\}\}$   
15:   $\bar{\mathbf{r}}_2^{t+1} = (\mathbf{x}_1^t - \alpha_1^t \mathbf{r}_1^t) / (1 - \alpha_1^t)$   
16:   $\mathbf{r}_2^{t+1} = \zeta^t \bar{\mathbf{r}}_2^{t+1} + (1 - \zeta^t) \mathbf{r}_2^t$   
17:   $\bar{\gamma}_2^{t+1} = \gamma_1^t (1 - \alpha_1^t) / \alpha_1^t$   
18:   $\gamma_2^{t+1} = \{\zeta^t (\bar{\gamma}_2^{t+1})^{-\frac{1}{2}} + (1 - \zeta^t) (\gamma_2^t)^{-\frac{1}{2}}\}^{-2}$   
19: end for
```



Figure 4.4: Ground truth, and Cartesian sampling and variable density sampling masks of acceleration $R = 4$

4.6 Numerical Experiments

In this section, we present the results of numerical experiments that study the behavior of our proposed algorithms, DD-AMP (Alg. 6) and DD-VAMP (Alg. 9), in comparison to PnP-ADMM [95], VDAMP [57] and BM3D-AMP-MRI [33] algorithms on image reconstruction in MRI. We considered both single-coil and multicoil MRI, and used real-valued images from the fastMRI [103] knee dataset. We considered real-valued images due to the lack of a good denoiser for complex valued images. As a future work, we aim to test our algorithms on the more challenging complex valued images.

For our experiments, we used a non fat-suppressed mid-slice knee MR image of size 128×128 , and Cartesian and variable density sampling masks of acceleration $R = 4$, shown in Fig. 4.4. The noise precision γ_w was set to achieve SNR=40 dB. We tuned PnP-ADMM to maximize the PSNR.

4.6.1 Single-Coil MRI

We first tested our algorithms on single-coil MRI, where the measurements \mathbf{y} are generated according to (4.1). Figure 4.5 shows the reconstruction PSNR and SSIM [99] of different plug-and-play algorithms for Cartesian sampling mask. Here, we used the BM3D [26] denoiser in DD-AMP, DD-VAMP and PnP-ADMM algorithms. VDAMP was designed for variable density sampling masks (Fig. 4.1) and doesn't work well with Cartesian masks. In Fig. 4.6, we used the DnCNN [104] denoiser included in the D-AMP toolbox⁸. Using a more powerful denoiser increased the accuracy of the algorithms, as expected. In Fig. 4.7, we used the DnCNN denoiser and variable density sampling mask. From Figures 4.5, 4.6 and 4.7, we see that the performance of DD-VAMP and PnP-ADMM is very similar and that DD-AMP performed a little worse.

⁸github.com/ricedsp/D-AMP_Toolbox

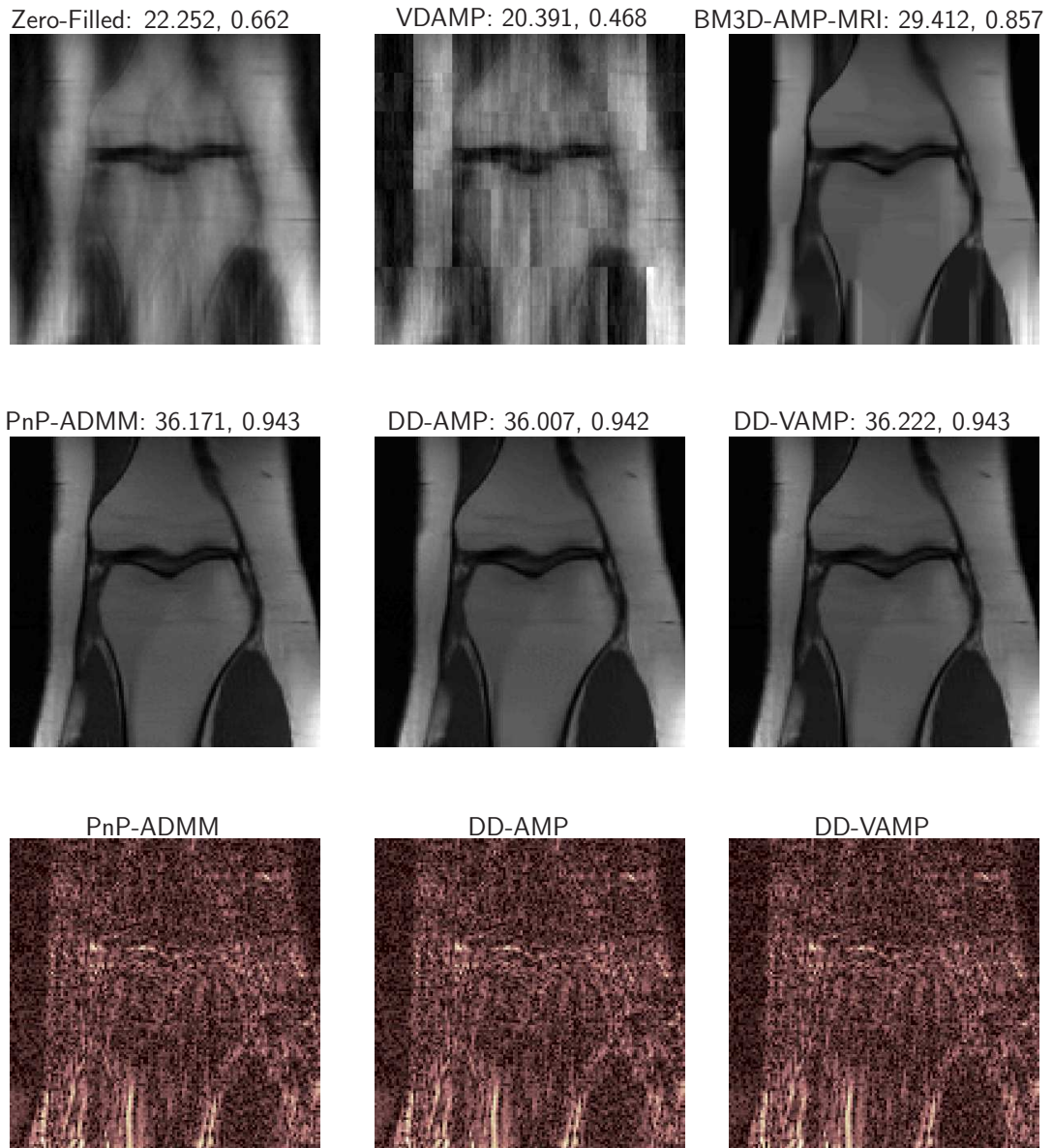


Figure 4.5: Reconstruction PSNR (dB) and SSIM of single-coil knee MRI for Cartesian sampling mask of acceleration $R = 4$. Using BM3D denoiser in DD-AMP, DD-VAMP and PnP-ADMM algorithms. Figures in last row are the error images.

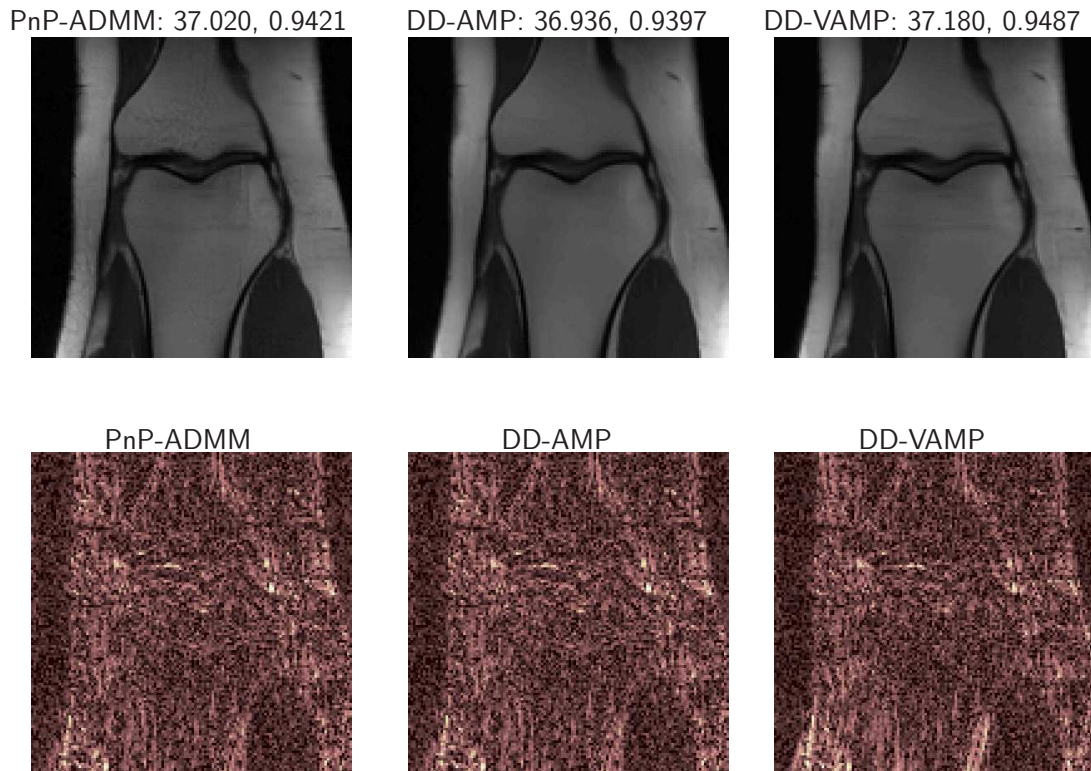


Figure 4.6: Reconstruction PSNR (dB) and SSIM of single-coil knee MRI for Cartesian sampling mask of acceleration $R = 4$. Using DnCNN denoiser in DD-AMP, DD-VAMP and PnP-ADMM algorithms. Figures in last row are error images.

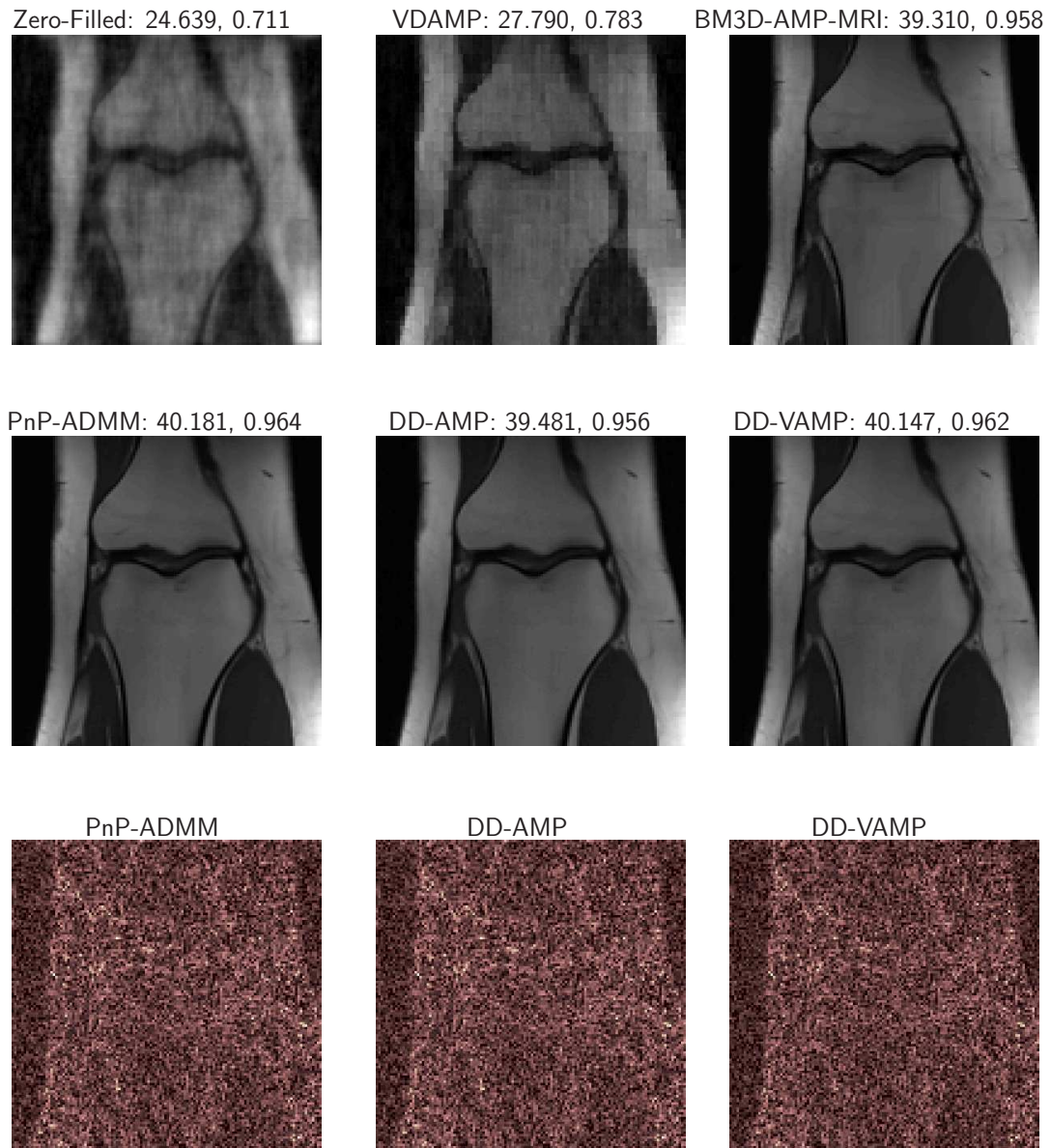


Figure 4.7: Reconstruction PSNR (dB) and SSIM of single-coil knee MRI using variable density sampling mask of acceleration $R = 4$. Using DnCNN denoiser in DD-AMP, DD-VAMP and PnP-ADMM algorithms. Figures in last row are the error images.

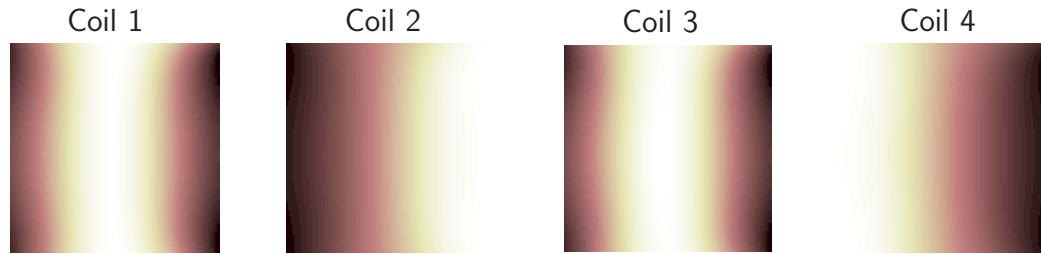


Figure 4.8: Coil sensitivities in multicoil MRI

4.6.2 Multicoil MRI

For our next experiment, we used $K = 4$ receiver coils and generated the coil sensitivity maps shown in Fig. 4.8 using the Biot-Savart law. VDAMP and BM3D-MRI-AMP are not applicable in this case, so we compared our algorithms against PnP-ADMM. We used Cartesian sampling and the BM3D denoiser in Fig. 4.9, and the DnCNN denoiser included in the D-AMP toolbox⁹ in Fig. 4.10. In Fig. 4.11, we used the variable density sampling and DnCNN denoiser. Similar to the single-coil case, the performance of DD-VAMP and PnP-ADMM is similar and that DD-AMP performed a little worse.

⁹github.com/ricedsp/D-AMP_Toolbox

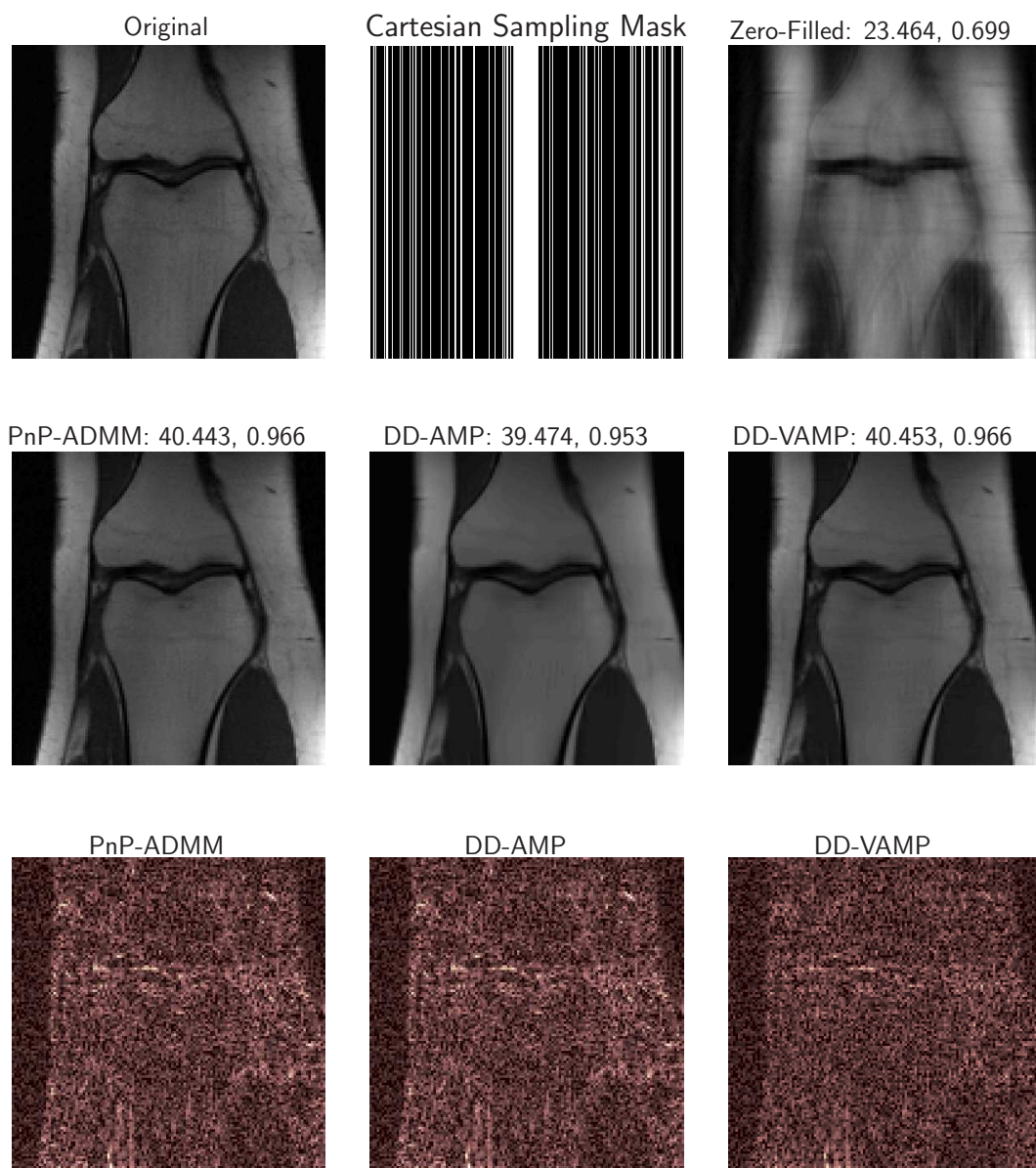


Figure 4.9: Reconstruction PSNR (dB) and SSIM of multicoil knee MRI for 4 receiver coils and Cartesian sampling of acceleration $R = 4$. Using BM3D denoiser in DD-AMP, DD-VAMP and PnP-ADMM. Figures in last row are the error images.

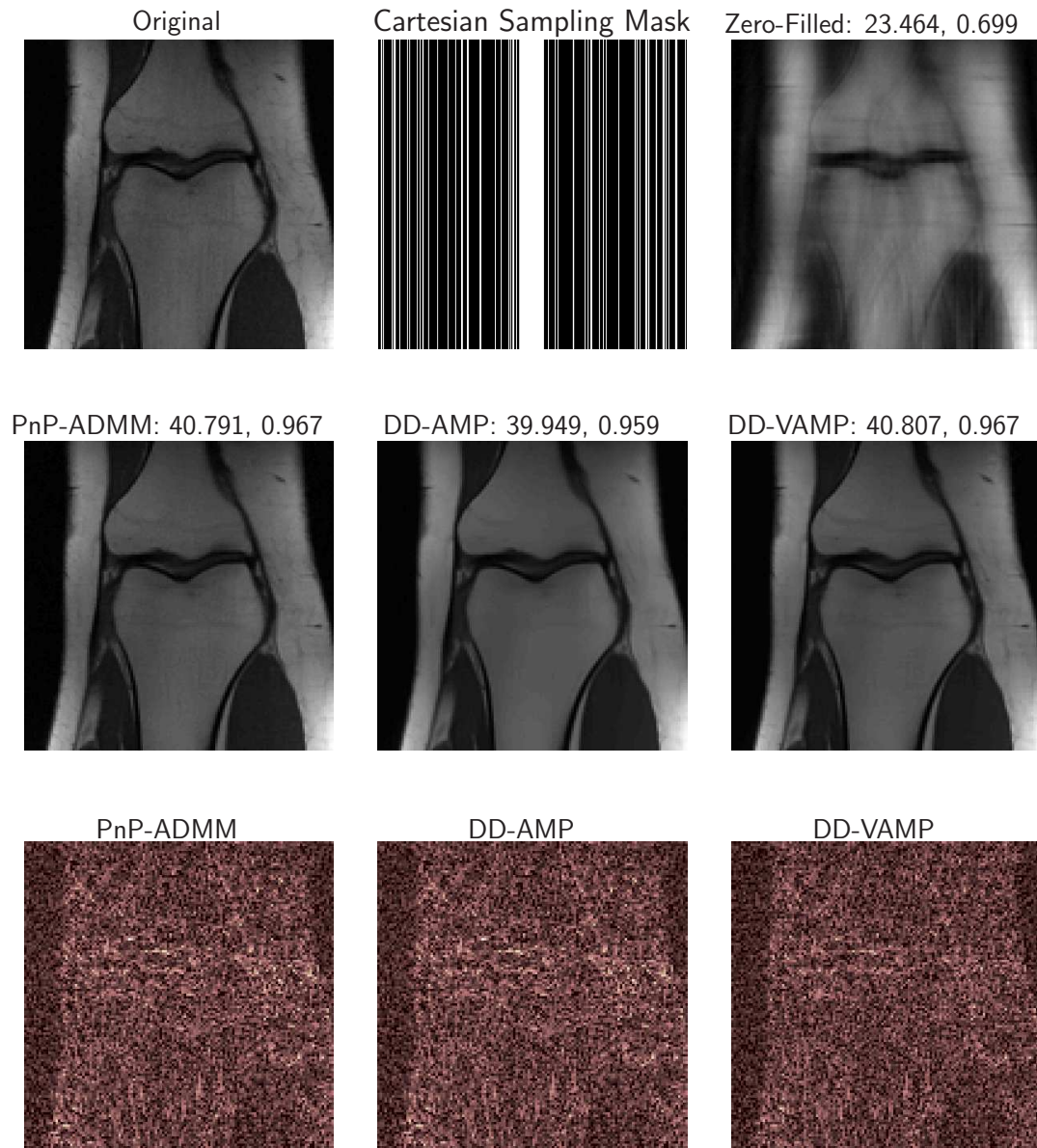


Figure 4.10: Reconstruction PSNR (dB) and SSIM of multicoil knee MRI for 4 receiver coils and Cartesian sampling of acceleration $R = 4$. Using DnCNN denoiser in DD-AMP, DD-VAMP and PnP-ADMM algorithms. Figures in last row are the error images.

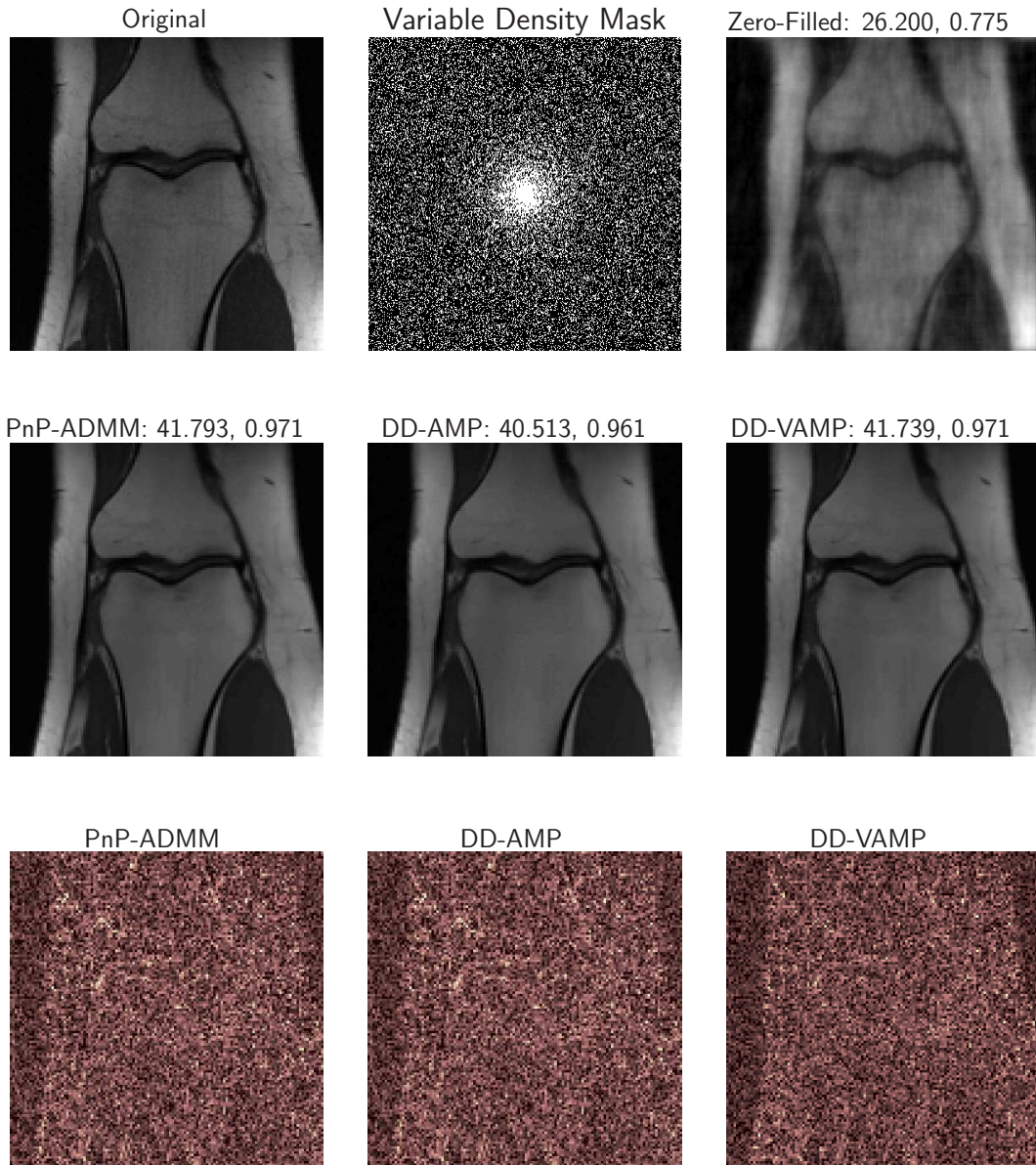


Figure 4.11: Reconstruction PSNR (dB) and SSIM of multicoil knee MRI for 4 receiver coils and variable density sampling of acceleration $R = 4$. Using DnCNN denoiser in DD-AMP, DD-VAMP and PnP-ADMM algorithms. Figures in last row are the error images.

4.6.3 Randomized Single-Coil MRI

In this experiment, we used a different measurement model in single-coil MRI,

$$\mathbf{y} = \mathbf{S}\mathbf{F}\mathbf{J}\mathbf{D}\mathbf{x} + \mathbf{w}, \quad (4.31)$$

where \mathbf{J} is a random permutation matrix, i.e., a random shuffling of the columns of \mathbf{I} , and $\mathbf{D} = \text{Diag}\{d_1, \dots, d_N\}$ for d_i uniformly sampled from $\{-1, 1\}$. The matrix $\mathbf{J}\mathbf{D}$ randomizes the measurements \mathbf{y} . This measurement model is not commonly seen in practice but we consider it to see whether the algorithms benefit from a simple randomization of the measurement operator.

We compared DD-AMP and DD-VAMP with PnP-ADMM using model (4.31). For single-coil MRI, Cartesian sampling mask of $R = 4$ and DnCNN denoiser, Fig. 4.12 shows the PSNR (dB) and SSIM performance of the algorithms. We see that the randomization in MRI helps AMP-type algorithms such as DD-AMP and DD-VAMP outperform PnP-ADMM. If one could implement the randomization $\mathbf{J}\mathbf{D}$ in a real-life MRI scanner, then it would be very beneficial to AMP-type algorithms.

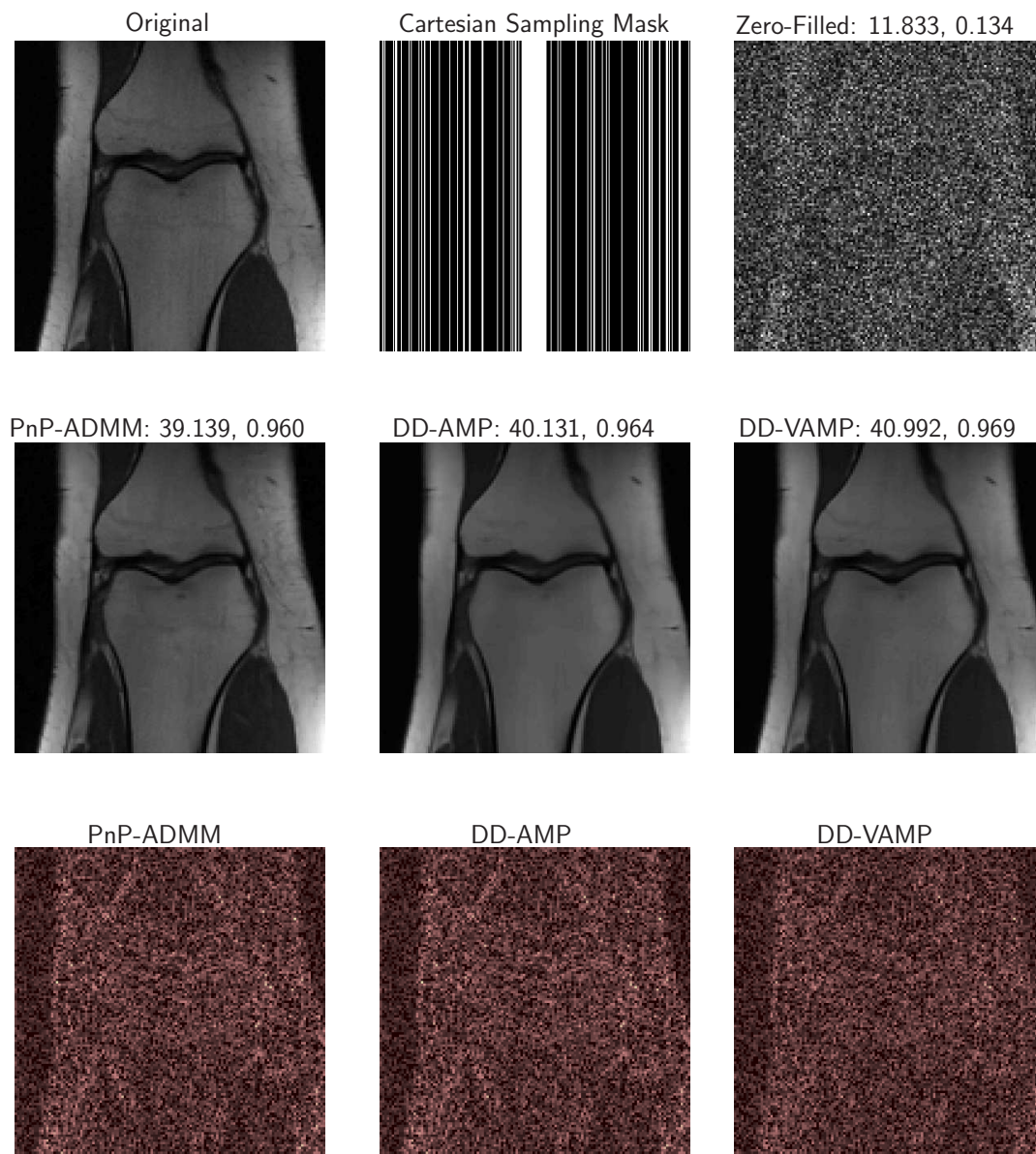


Figure 4.12: Reconstruction PSNR (dB) and SSIM of single-coil knee MRI with random coding, and Cartesian sampling of acceleration $R = 4$. Using DnCNN denoiser in DD-AMP, DD-VAMP and PnP-ADMM algorithms. Figures in last row are the error images.

Chapter 5: Conclusions

In this dissertation, we considered several linear and bilinear inverse problems and proposed algorithms to solve them using the Approximate Message Passing framework. We first proposed some theoretical results on the VAMP [73] algorithm, where we extended its state evolution (SE) analysis from separable to non-separable Lipschitz denoisers. Using the SE formalism, when the forward operator is right rotationally invariant (RRI), we can accurately predict the MSE performance of VAMP in high dimensions.

Next, we proposed the BAdVAMP algorithm for solving bilinear inverse problems. BAdVAMP jointly recovers the vector \mathbf{b} and the matrix \mathbf{C} from noisy measurements $\mathbf{Y} = \mathbf{A}(\mathbf{b})\mathbf{C} + \mathbf{W}$, where $\mathbf{A}(\cdot)$ is a known affine linear function of \mathbf{b} (i.e., $\mathbf{A}(\mathbf{b}) = \mathbf{A}_0 + \sum_{i=1}^Q b_i \mathbf{A}_i$ with known matrices \mathbf{A}_i). To solve this problem, we combined the VAMP algorithm [75], the EM algorithm [61], and variance auto-tuning [45] in a manner appropriate for bilinear recovery. We demonstrated numerically that the proposed approach has robustness advantages over other state-of-the-art bilinear recovery algorithms, including lifted VAMP [35] and EM-PBiGAMP [66]. As future work, we hope to rigorously analyze BAdVAMP through the state-evolution formalism.

Lastly, we proposed image reconstruction algorithms in MRI by modifying the existing AMP [32] and VAMP [73] algorithms. In AMP, we scaled the measurements in k-space by the inverse probability of the random sampling mask to debias the denoiser input error, which led to a significant increase in the reconstruction accuracy. In VAMP, we proposed an improved damping scheme that empirically stabilizes it for MRI. We numerically demonstrated that our proposed algorithms are more accurate than existing AMP-based algorithms in single-coil MRI. We showed that our algorithms are also applicable to multicoil MRI with known coil sensitivities. As a future work, we hope to do the following: rigorous analysis of the proposed damping scheme in VAMP; using ideas from BAdVAMP, learn the coil sensitivities and recover the image simultaneously in VAMP; design a good denoiser for complex-valued MR images and apply our algorithms.

Appendix A: Signal Denoisers

A.1 Singular Value Thresholding

Consider the estimation of a low-rank matrix \mathbf{X}_0 from linear measurements $\mathbf{y} = \mathcal{A}(\mathbf{X}_0)$, where \mathcal{A} is some linear operator [16]. Writing the SVD of \mathbf{R} as $\mathbf{R} = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$, the SVT denoiser is defined as

$$\mathbf{g}_1(\mathbf{R}, \gamma) \triangleq \sum_i (\sigma_i - \gamma)_+ \mathbf{u}_i \mathbf{v}_i^\top, \quad (\text{A.1})$$

where $(x)_+ \triangleq \max\{0, x\}$. To show that $\mathbf{g}_1(\cdot)$ in (A.1) is uniformly pseudo-Lipschitz, we first note that $\mathbf{g}_1(\cdot)$ is the proximal operator of the nuclear norm $\|\cdot\|_*$, i.e.,

$$\mathbf{g}_1(\mathbf{R}, \gamma) = \arg \min_{\mathbf{X} \in \mathbb{R}^{N_1 \times N_2}} \gamma \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X} - \mathbf{R}\|_F^2.$$

Define $\langle \cdot, \cdot \rangle$ to be the standard inner product, i.e., $\langle \mathbf{R}_1, \mathbf{R}_2 \rangle = \text{tr}\{\mathbf{R}_2^\top \mathbf{R}_1\}$. From [24], we have that $\mathbf{g}_1(\cdot)$ is non-expansive because the nuclear norm is convex and proper, i.e., we can write

$$\begin{aligned} \|\mathbf{g}_1(\mathbf{R}_1, \gamma) - \mathbf{g}_1(\mathbf{R}_2, \gamma)\|_F^2 &\leq \langle \mathbf{R}_1 - \mathbf{R}_2, \mathbf{g}_1(\mathbf{R}_1, \gamma) - \mathbf{g}_1(\mathbf{R}_2, \gamma) \rangle \\ \Rightarrow \|\mathbf{g}_1(\mathbf{R}_1, \gamma) - \mathbf{g}_1(\mathbf{R}_2, \gamma)\|_F &\leq \|\mathbf{R}_1 - \mathbf{R}_2\|_F. \end{aligned} \quad (\text{A.2})$$

Let the SVD of $\mathbf{R}_2 \in \mathbb{R}^{N_1 \times N_2}$ be $\sum_{i=1}^{\min\{N_1, N_2\}} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$. We can generalize the Lipschitz condition in (A.2) into

$$\begin{aligned}
\|\mathbf{g}_1(\mathbf{R}_1, \gamma_1) - \mathbf{g}_1(\mathbf{R}_2, \gamma_2)\|_F &= \|\mathbf{g}_1(\mathbf{R}_1, \gamma_1) - \mathbf{g}_1(\mathbf{R}_2, \gamma_1) + \mathbf{g}_1(\mathbf{R}_2, \gamma_1) - \mathbf{g}_1(\mathbf{R}_2, \gamma_2)\|_F \\
&\leq \|\mathbf{g}_1(\mathbf{R}_1, \gamma_1) - \mathbf{g}_1(\mathbf{R}_2, \gamma_1)\|_F + \|\mathbf{g}_1(\mathbf{R}_2, \gamma_1) - \mathbf{g}_1(\mathbf{R}_2, \gamma_2)\|_F \\
&\leq \|\mathbf{R}_1 - \mathbf{R}_2\|_F + \|\mathbf{g}_1(\mathbf{R}_2, \gamma_1) - \mathbf{g}_1(\mathbf{R}_2, \gamma_2)\|_F \\
&\stackrel{(a)}{=} \|\mathbf{R}_1 - \mathbf{R}_2\|_F \\
&\quad + \left\| \sum_{i=1}^{\min\{N_1, N_2\}} ((\sigma_i - \gamma_1)_+ - (\sigma_i - \gamma_2)_+) \mathbf{u}_i \mathbf{v}_i^\top \right\|_F \\
&\leq \|\mathbf{R}_1 - \mathbf{R}_2\|_F + \sum_{i=1}^{\min\{N_1, N_2\}} |(\sigma_i - \gamma_1)_+ - (\sigma_i - \gamma_2)_+| \\
&\leq \|\mathbf{R}_1 - \mathbf{R}_2\|_F + \min\{N_1, N_2\} |\gamma_1 - \gamma_2| \\
&\stackrel{(b)}{\leq} \|\mathbf{R}_1 - \mathbf{R}_2\|_F + \sqrt{N} |\gamma_1 - \gamma_2|,
\end{aligned}$$

where in (a) we have used the the definition of $\mathbf{g}_1(\cdot)$ from (A.1) and the SVD of \mathbf{R}_2 , and in (b) we used $\min\{N_1, N_2\} \leq \sqrt{N_1 N_2} = \sqrt{N}$. Next, we show that $\mathbf{g}_1(\cdot)$ also satisfies the convergence conditions in Definition 2. Let \mathbf{Z}_1 and \mathbf{Z}_2 be two sequences constructed according to Definition 1 and let \mathbf{x}_0 be the true signal. Assume that

$$\lim_{N \rightarrow \infty} \frac{1}{N} \|\mathbf{X}_0\|_F^2 \quad \text{and} \quad \lim_{N \rightarrow \infty} \frac{1}{N} \langle \mathbf{X}_0, \mathbf{Z}_1 \rangle \quad \text{exist almost surely.} \quad (\text{A.3})$$

If we write $\tilde{\mathbf{g}}(\mathbf{R}, \gamma) = \text{vec}(\mathbf{g}_1(\mathbf{R}, \gamma)) \triangleq [g_{1,1}(\mathbf{R}, \gamma), \dots, g_{1,N}(\mathbf{R}, \gamma)]^\top$, then the following series converges because it is bounded:

$$\begin{aligned}
&\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N |\tilde{g}_i(\mathbf{X}_0 + \mathbf{Z}_1, \gamma_1) \tilde{g}_i(\mathbf{X}_0 + \mathbf{Z}_2, \gamma_2)| \\
&\leq \lim_{N \rightarrow \infty} \frac{1}{N} \|\tilde{\mathbf{g}}(\mathbf{X}_0 + \mathbf{Z}_1, \gamma_1)\|_2 \|\tilde{\mathbf{g}}(\mathbf{X}_0 + \mathbf{Z}_2, \gamma_2)\|_2 \\
&\leq \lim_{N \rightarrow \infty} \sqrt{\frac{1}{N} \|\mathbf{X}_0 + \mathbf{Z}_1\|_F^2} \sqrt{\frac{1}{N} \|\mathbf{X}_0 + \mathbf{Z}_2\|_F^2} \\
&\stackrel{(a)}{<} \infty,
\end{aligned}$$

where (a) follows from the assumption (A.3). Since absolute convergence implies convergence, the following series also converges:

$$\begin{aligned} & \lim_{N \rightarrow \infty} \frac{1}{N} \langle \mathbf{g}_1(\mathbf{X}_0 + \mathbf{Z}_1, \gamma_1), \mathbf{g}_1(\mathbf{X}_0 + \mathbf{Z}_2, \gamma_2) \rangle \\ &= \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \tilde{g}_i(\mathbf{X}_0 + \mathbf{Z}_1, \gamma_1) \tilde{g}_i(\mathbf{X}_0 + \mathbf{Z}_2, \gamma_2). \end{aligned} \quad (\text{A.4})$$

If we choose the covariance matrix in Definition 1 to be $\mathbf{S} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$, then we get

$$\lim_{N \rightarrow \infty} \frac{1}{N} \langle \mathbf{g}_1(\mathbf{X}_0 + \mathbf{Z}_1, \gamma_1), \mathbf{X}_0 \rangle = \lim_{N \rightarrow \infty} \frac{1}{N} \langle \mathbf{g}_1(\mathbf{X}_0 + \mathbf{Z}_1, \gamma_1), \mathbf{g}_1(\mathbf{X}_0 + \mathbf{Z}_2, 0) \rangle. \quad (\text{A.5})$$

Thus, (A.5) also converges since it is a special case of (A.4). Similarly, it can be shown that $\frac{1}{N} \langle \mathbf{Z}_2, \mathbf{g}_1(\mathbf{X}_0 + \mathbf{Z}_1, \gamma_1) \rangle$ is uniformly Lipschitz.

Bibliography

- [1] A. Agarwal, S. Negahban, and M. J. Wainwright. Matrix decomposition via convex relaxation: Optimal rates in high dimensions. *Ann. Statist.*, 40(2):1171–1197, 2012.
- [2] Michal Aharon, Michael Elad, and Alfred Bruckstein. k-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Trans. Signal Process.*, 54(11):4311–4322, 2006.
- [3] A. Ahmed, B. Recht, and J. Romberg. Blind deconvolution using convex programming. *IEEE Trans. Inform. Theory*, 60(3):1711–1732, 2014.
- [4] S. D. Babacan, M. Luessi, R. Molina, and A. K. Katsaggelos. Sparse Bayesian methods for low-rank matrix estimation. *IEEE Trans. Signal Process.*, 60(8):3964–3977, August 2012.
- [5] Laura Balzano, Robert Nowak, and Benjamin Recht. Online identification and tracking of subspaces from highly incomplete information. In *Proc. Allerton Conf. Commun. Control Comput.*, pages 704–711, September 2010.
- [6] Jean Barbier, Mohamad Dia, Nicolas Macris, and Florent Krzakala. The mutual information in random linear estimation. In *Proc. Allerton Conf. Commun. Control Comput.*, pages 625–632, 2016.
- [7] Richard Barrett, Michael Berry, Tony F Chan, James Demmel, June Donato, Jack Dongarra, Victor Eijkhout, Roldan Pozo, Charles Romine, and Henk Van der Vorst. *Templates for the solution of linear systems: building blocks for iterative methods*. SIAM, 1994.
- [8] M. Bayati, M. Lelarge, and A. Montanari. Universality in polytope phase transitions and message passing algorithms. *Ann. App. Prob.*, 25(2):753–822, 2015.
- [9] M. Bayati and A. Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Trans. Inform. Theory*, 57(2):764–785, February 2011.

- [10] A. Beck and M. Teboulle. Fast gradient-based algorithms for constrained total variation image denoising and deblurring problem. *IEEE Trans. Image Process.*, 18(11):2419–2434, 2009.
- [11] S. Bellini. Bussgang techniques for blind deconvolution and equalization. In S. Haykin, editor, *Blind Deconvolution*, chapter 2, pages 8–59. Prentice Hall, 1994.
- [12] Raphael Berthier, Andrea Montanari, and Phan-Minh Nguyen. State evolution for approximate message passing with non-separable functions. *Inform. Inference*, 2019.
- [13] C. Bilen, G. Puy, and R. Gribonval. Convex optimization approaches for blind sensor calibration using sparsity. *IEEE Trans. Signal Process.*, 62(18):4847–4856, 2014.
- [14] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.*, 3(1):1–122, 2011.
- [15] A. Buades, B. Coll, and J.-M. Morel. A non-local algorithm for image denoising. In *Proc. IEEE Conf. Comp. Vision Pattern Recog.*, volume 2, pages 60–65, 2005.
- [16] J.-F. Cai, E. J. Candès, and Z. Shen. A singular value thresholding algorithm for matrix completion. *SIAM J. Optim.*, 20(4):1956–1982, January 2010.
- [17] Burak Çakmak and Ole Winther Bernard H. Fleury. S-AMP for non-linear observation models. In *Proc. IEEE Int. Symp. Inform. Thy.*, pages 2807–2811, 2015.
- [18] F. Caltagirone, F. Krzakala, and L. Zdeborová. On convergence of approximate message passing. In *Proc. IEEE Int. Symp. Inform. Thy.*, pages 1812–1816, July 2014.
- [19] E. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inform. Theory*, 52(2):489–509, February 2006.
- [20] E. J. Candès, X. Li, Yi Ma, and J. Wright. Robust principal component analysis? *J. ACM*, 58(3):11, May 2011.
- [21] E. J. Candès and Y. Plan. Matrix completion with noise. *Proc. IEEE*, 98(6):925–936, June 2010.

- [22] E. J. Candès and Y. Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Trans. Inform. Theory*, 57(4):2342–2359, 2011.
- [23] E. J. Candès, T. Strohmer, and V. Voroninski. PhaseLift: Exact and stable signal recovery from magnitude measurements via convex programming. *Commun. Pure & Appl. Math.*, 66(8):1241–1274, 2013.
- [24] Emmanuel J Candes, Carlos A Sing-Long, and Joshua D Trzasko. Unbiased risk estimates for singular value thresholding and spectral estimators. *IEEE Trans. Signal Process.*, 61(19):4643–4657, 2013.
- [25] Stanley H Chan, Xiran Wang, and Omar A Elgendy. Plug-and-play admm for image restoration: Fixed-point convergence and applications. *IEEE Trans. Comp. Imag.*, 3(1):84–98, 2016.
- [26] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian. Image denoising by sparse 3-D transform-domain collaborative filtering. *IEEE Trans. Image Process.*, 16(8):2080–2095, 2007.
- [27] I. Daubechies, M. Defrise, and C. D. Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint,. *Commun. Pure & Appl. Math.*, 57(11):1413–1457, November 2004.
- [28] M. A. Davenport and J. Romberg. An overview of low-rank matrix recovery from incomplete observations. *IEEE J. Sel. Topics Signal Process.*, 10(4):608–622, 2016.
- [29] A. Dempster, N. M. Laird, and D. B. Rubin. Maximum-likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc.*, 39:1–17, 1977.
- [30] Yash Deshpande and Andrea Montanari. Information-theoretically optimal sparse PCA. *Proc. IEEE Int. Symp. Inform. Thy.*, pages 2197–2201, 2014.
- [31] D. L. Donoho. Compressed sensing. *IEEE Trans. Inform. Theory*, 52(9):1289–1306, September 2006.
- [32] D. L. Donoho, A. Maleki, and A. Montanari. Message passing algorithms for compressed sensing. *Proc. Nat. Acad. Sci.*, 106(45):18914–18919, November 2009.
- [33] Ender M Eksioğlu and A Korhan Tanc. Denoising amp for mri reconstruction: Bm3d-amp-mri. *SIAM J. Imag. Sci.*, 11(3):2090–2109, 2018.
- [34] Y. C. Eldar and G. Kutyniok. *Compressed Sensing: Theory and Applications*. Cambridge Univ. Press, New York, 2012.

- [35] A. K. Fletcher, P. Pandit, S. Rangan, S. Sarkar, and P. Schniter. Plug-in estimation in high-dimensional linear inverse problems: A rigorous analysis. In *Proc. Neural Inform. Process. Syst. Conf.*, pages 7440–7449, 2018.
- [36] Alyson K. Fletcher, Mojtaba Sahraee-Ardakan, Sundeep Rangan, and Philip Schniter. Expectation consistent approximate inference: Generalizations and convergence. In *Proc. IEEE Int. Symp. Inform. Thy.*, pages 190–194, 2016.
- [37] Alyson K. Fletcher, Mojtaba Sahraee-Ardakan, Sundeep Rangan, and Philip Schniter. Rigorous dynamics and consistent estimation in arbitrarily conditioned linear systems. In *Proc. Neural Inform. Process. Syst. Conf.*, pages 2542–2551, 2017.
- [38] T. Goldstein, C. Studer, and R. Baraniuk. Forward-backward splitting with a FASTA implementation. *arXiv:1411.3406*, 2014.
- [39] Jun He, Laura Balzano, and Arthur Szlam. Incremental gradient on the Grassmannian for online foreground and background separation in subsampled video. In *Proc. IEEE Conf. Comp. Vision Pattern Recog.*, pages 1568–1575, 2012.
- [40] C. Hegde and R. G. Baraniuk. Sampling and recovery of pulse streams. *IEEE Trans. Signal Process.*, 59(14):1505–1517, 2011.
- [41] T. Heskes and O. Zoeter. Expectation propagation for approximate inference in dynamic Bayesian networks. In *Proc. Uncertainty Artif. Intell.*, pages 313–320, 2002.
- [42] Yoshiyuki Kabashima, Florent Krzakala, Marc Mézard, Ayaka Sakata, and Lenka Zdeborová. Phase transitions and sample complexity in Bayes-optimal matrix factorization. *IEEE Trans. Inform. Theory*, 62(7):4228–4265, 2016.
- [43] G. K. Kaleh and R. Vallet. Joint parameter estimation and symbol detection for linear or nonlinear unknown channels. *IEEE Trans. Commun.*, 42:2406–2413, Jul. 1994.
- [44] U. Kamilov, H. Mansour, and B. Wohlberg. A plug-and-play priors approach for solving nonlinear imaging inverse problems. *IEEE Signal Process. Lett.*, 24(12):1872–1876, May 2017.
- [45] U. S. Kamilov, S. Rangan, A. K. Fletcher, and M. Unser. Approximate message passing with consistent parameter estimation and applications to sparse learning. *IEEE Trans. Inform. Theory*, 60(5):2969–2985, May 2014.
- [46] A. Kyrillidis and V. Cevher. Matrix recipes for hard thresholding methods. *J. Math. Imaging Vis.*, 48:235–265, 2014.

- [47] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Proc. Neural Inform. Process. Syst. Conf.*, pages 556–562, 2001.
- [48] Thibault Lesieur, Florent Krzakala, and Lenka Zdeborova. Phase transitions in sparse PCA. In *Proc. IEEE Int. Symp. Inform. Thy.*, pages 1635–1639, 2015.
- [49] Thibault Lesieur, Florent Krzakala, and Lenka Zdeborov. Constrained low-rank matrix estimation: Phase transitions, approximate message passing and applications. *J. Stat. Mech.*, 2017(7):073403, 2017.
- [50] Z. Lin, M. Chen, L. Wu, and Y. Ma. The augmented Lagrange multiplier method for exact recovery of corrupted low-rank matrices. *arXiv:1009.5055*, 2010.
- [51] S. Ling and T. Strohmer. Self-calibration and biconvex compressive sensing. *Inverse Problems*, 31(11):115002, 2015.
- [52] Pengju Liu, Hongzhi Zhang, Kai Zhang, Liang Lin, and Wangmeng Zuo. Multi-level wavelet-cnn for image restoration. In *Proc. IEEE Conf. Comp. Vision Pattern Recog.*, pages 773–782, 2018.
- [53] M. Lustig, D. L. Donoho, J. M. Santos, and J. M. Pauly. Compressed sensing MRI. *IEEE Signal Process. Mag.*, 25(2):72–82, March 2008.
- [54] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *J. Mach. Learn. Res.*, 11:19–60, January 2010.
- [55] Ryosuke Matsushita and Toshiyuki Tanaka. Low-rank matrix reconstruction and clustering via approximate message passing. In *Proc. Neural Inform. Process. Syst. Conf.*, pages 917–925, 2013.
- [56] Christopher A. Metzler, Arian Maleki, and Richard G. Baraniuk. From denoising to compressed sensing. *IEEE Trans. Inform. Theory*, 62(9):5117–5144, 2016.
- [57] Charles Millard, Aaron T Hess, Boris Mailhé, and Jared Tanner. Approximate message passing with a colored aliasing model for variable density fourier sampled images. *arXiv:2003.02701*, 2020.
- [58] T. Minka. *A Family of Approximate Algorithms for Bayesian Inference*. PhD thesis, Dept. Comp. Sci. Eng., MIT, Cambridge, MA, Jan. 2001.
- [59] Léo Miolane. Fundamental limits of low-rank matrix estimation: The non-symmetric case. *arXiv:1702.00473*, 2017.

- [60] B. K. Natarajan. Sparse approximate solutions to linear systems. *SIAM J. Scientific Comput.*, 24(2):227–234, 1995.
- [61] Radford Neal and Geoffrey Hinton. A view of the EM algorithm that justifies incremental, sparse, and other variants. In Michael I. Jordan, editor, *Learning in Graphical Models*, pages 355–368. MIT Press, 1998.
- [62] Shunsuke Ono. Primal-dual plug-and-play image restoration. *IEEE Signal Process. Lett.*, 24(8):1108–1112, 2017.
- [63] M. Opper and O. Winther. Expectation consistent free energies for approximate inference. In *Proc. Neural Inform. Process. Syst. Conf.*, pages 1001–1008, 2005.
- [64] Manfred Opper and Ole Winther. Adaptive and self-averaging thouless-anderson-palmer mean-field theory for probabilistic modeling. *Physical Rev. E*, 64(5):056131, 2001.
- [65] Christopher C Paige and Michael A Saunders. Lsqr: An algorithm for sparse linear equations and sparse least squares. *ACM Transactions on Mathematical Software (TOMS)*, 8(1):43–71, 1982.
- [66] J. T. Parker and P. Schniter. Parametric bilinear generalized approximate message passing. *IEEE J. Sel. Topics Signal Process.*, 10(4):795–808, 2016.
- [67] J. T. Parker, P. Schniter, and V. Cevher. Bilinear generalized approximate message passing—Part I: Derivation. *IEEE Trans. Signal Process.*, 62(22):5839–5853, November 2014.
- [68] J. T. Parker, P. Schniter, and V. Cevher. Bilinear generalized approximate message passing—Part II: Applications. *IEEE Trans. Signal Process.*, 62(22):5854–5867, November 2014.
- [69] Lawrence R Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
- [70] Sathish Ramani, Thierry Blu, and Michael Unser. Monte-Carlo SURE: A black-box optimization of regularization parameters for general denoising algorithms. *IEEE Trans. Image Process.*, 17(9):1540–1554, 2008.
- [71] S. Rangan. Generalized approximate message passing for estimation with random linear mixing. In *Proc. IEEE Int. Symp. Inform. Thy.*, pages 2168–2172, August 2011. (full version at *arXiv:1010.5141*).
- [72] S. Rangan and A. K. Fletcher. Iterative estimation of constrained rank-one matrices in noise. In *Proc. IEEE Int. Symp. Inform. Thy.*, pages 1246–1250, 2012.

- [73] S. Rangan, P. Schniter, and A. K. Fletcher. Vector approximate message passing. *IEEE Trans. Inform. Theory*, to appear (see also arXiv:1610.03082).
- [74] S. Rangan, P. Schniter, E. Riegler, A. Fletcher, and V. Cevher. Fixed points of generalized approximate message passing with arbitrary matrices. *IEEE Trans. Inform. Theory*, 62(12):7464–7474, December 2016.
- [75] Sundeep Rangan, Philip Schniter, and Alyson K Fletcher. Vector approximate message passing. In *Proc. IEEE Int. Symp. Inform. Thy.*, pages 1588–1592, 2017.
- [76] E. T. Reehorst and P. Schniter. Regularization by denoising: Clarifications and new interpretations. *IEEE Trans. Comp. Imag.*, 5(1):52–67, March 2019.
- [77] G. Reeves. Additivity of information in multilayer networks via additive Gaussian noise transforms. In *Proc. Allerton Conf. Commun. Control Comput.*, pages 1064–1070, 2017.
- [78] Galen Reeves and Henry D. Pfister. The replica-symmetric prediction for compressed sensing with Gaussian matrices is exact. In *Proc. IEEE Int. Symp. Inform. Thy.*, 2016.
- [79] Yaniv Romano, Michael Elad, and Peyman Milanfar. The little engine that could: Regularization by denoising (RED). *SIAM J. Imag. Sci.*, 10(4):1804–1844, 2017.
- [80] Elad Romanov and Matan Gavish. Near-optimal matrix recovery from random linear measurements. *Proc. Nat. Acad. Sci.*, 2018.
- [81] Ron Rubinstein, Alfred M Bruckstein, and Michael Elad. Dictionaries for sparse representation modeling. *Proc. IEEE*, 98(6):1045–1057, 2010.
- [82] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D*, 60:259–268, 1992.
- [83] Cynthia Rush and Ramji Venkataramanan. Finite-sample analysis of approximate message passing. In *Proc. IEEE Int. Symp. Inform. Thy.*, pages 755–759, 2016.
- [84] Philip Schniter, Sundeep Rangan, and Alyson Fletcher. Plug-and-play image recovery using vector amp. *Proc. Intl. Biomed. Astronom. Signal Process. (BASP) Frontiers Workshop*, 2017.
- [85] C. Schülke, P. Schniter, and L. Zdeborová. Phase diagram of matrix compressed sensing. *Physical Rev. E*, 94(6):062136(1–16), December 2016.

- [86] M. Seeger. Expectation propagation for exponential families. Technical Report 161464, EPFL, 2005.
- [87] Daniel A Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. In *Proc. Conf. Learning Thy.*, pages 37.1–37.18, 2012.
- [88] P. Sun, Z. Wang, and P. Schniter. Joint channel-estimation and equalization of single-carrier systems via bilinear AMP. *IEEE Trans. Signal Process.*, 66(10):2772–2785, 2018.
- [89] K. Takeuchi. Rigorous dynamics of expectation-propagation-based signal recovery from unitarily invariant measurements. In *Proc. IEEE Int. Symp. Inform. Thy.*, pages 501–505, 2017.
- [90] D. J. Thouless, P. W. Anderson, and R. G. Palmer. Solution of solvable model of a spin glass. *Phil. Mag.*, 35:593–601, 1977.
- [91] A. N. Tikhonov and V. Arsenin. *Solutions of Ill-Posed Problems*. V.H. Winston and Sons, Washington, 1977.
- [92] Carlo Tomasi and Roberto Manduchi. Bilateral filtering for gray and color images. In *Proc. IEEE Int. Conf. Comput. Vis.*, pages 839–846. IEEE, 1998.
- [93] Antonia M. Tulino, Giuseppe Caire, Sergio Verdú, and Shlomo Shamai (Shitz). Support recovery with sparsely sampled free random matrices. *IEEE Trans. Inform. Theory*, 59(7):4243–4271, July 2013.
- [94] Martin Uecker, Peng Lai, Mark J Murphy, Patrick Virtue, Michael Elad, John M Pauly, Shreyas S Vasanawala, and Michael Lustig. Espiritan eigenvalue approach to autocalibrating parallel mri: where sense meets grappa. *Magnetic Resonance Med.*, 71(3):990–1001, 2014.
- [95] S. V. Venkatakrisnan, C. A. Bouman, and B. Wohlberg. Plug-and-play priors for model based reconstruction. In *Proc. IEEE Global Conf. Signal Info. Process.*, pages 945–948, 2013.
- [96] J. Vila, P. Schniter, S. Rangan, F. Krzakala, and L. Zdeborová. Adaptive damping and mean removal for the generalized approximate message passing algorithm. In *Proc. IEEE Int. Conf. Acoust. Speech & Signal Process.*, pages 2021–2025, 2015.
- [97] J. P. Vila and P. Schniter. Expectation-maximization Gaussian-mixture approximate message passing. *IEEE Trans. Signal Process.*, 61(19):4658–4672, Oct. 2013.

- [98] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.*, 1, May 2008.
- [99] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [100] A. E. Waters, A. C. Sankaranarayanan, and R. G. Baraniuk. SpaRCS: Recovering low-rank and sparse matrices from compressive measurements. In *Proc. Neural Inform. Process. Syst. Conf.*, pages 1089–1097, 2011.
- [101] Zaiwen Wen, Wotao Yin, and Yin Zhang. Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm. *Mathematical Programming Computation*, 4:333–361, 2012.
- [102] J. Wright, A. Ganesh, K. Min, and Y. Ma. Compressive principal component pursuit. *Inform. Inference*, 2(1):32–68, 2013.
- [103] Jure Zbontar, Florian Knoll, Anuroop Sriram, Matthew J. Muckley, Mary Bruno, Aaron Defazio, Marc Parente, Krzysztof J. Geras, Joe Katsnelson, Hersh Chandarana, Zizhao Zhang, Michal Drozdal, Adriana Romero, Michael Rabbat, Pascal Vincent, James Pinkerton, Duo Wang, Nafissa Yakubova, Erich Owens, C. Lawrence Zitnick, Michael P. Recht, Daniel K. Sodickson, and Yvonne W. Lui. fastMRI: An open dataset and benchmarks for accelerated MRI. *arXiv:1811.08839*, 2018.
- [104] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising. *IEEE Trans. Image Process.*, 26(7):3142–3155, 2017.
- [105] H. Zhu, G. Leus, and G. B. Giannakis. Sparsity-cognizant total least-squares for perturbed compressive sampling. *IEEE Trans. Signal Process.*, 59(5):2002–2016, May 2011.