



THE UNIVERSITY *of* EDINBURGH

This thesis has been submitted in fulfilment of the requirements for a postgraduate degree (e.g. PhD, MPhil, DClinPsychol) at the University of Edinburgh. Please note the following terms and conditions of use:

- This work is protected by copyright and other intellectual property rights, which are retained by the thesis author, unless otherwise stated.
- A copy can be downloaded for personal non-commercial research or study, without prior permission or charge.
- This thesis cannot be reproduced or quoted extensively from without first obtaining permission in writing from the author.
- The content must not be changed in any way or sold commercially in any format or medium without the formal permission of the author.
- When referring to this work, full bibliographic details including the author, title, awarding institution and date of the thesis must be given.

Compressed Sensing with Approximate Message Passing: Measurement Matrix and Algorithm Design

Chunli Guo



A thesis submitted for the degree of Doctor of Philosophy.
The University of Edinburgh.
November 2013

To Mom, Dad and Xiaohu

Abstract

Compressed sensing (CS) is an emerging technique that exploits the properties of a sparse or compressible signal to efficiently and faithfully capture it with a sampling rate far below the Nyquist rate. The primary goal of compressed sensing is to achieve the best signal recovery with the least number of samples. To this end, two research directions have been receiving increasing attention: customizing the measurement matrix to the signal of interest and optimizing the reconstruction algorithm. In this thesis, contributions in both directions are made in the Bayesian setting for compressed sensing. The work presented in this thesis focuses on the approximate message passing (AMP) schemes, a new class of recovery algorithm that takes advantage of the statistical properties of the CS problem.

First of all, a complete sample distortion (SD) framework is presented to fundamentally quantify the reconstruction performance for a certain pair of measurement matrix and recovery scheme. In the SD setting, the non-optimality region of the homogeneous Gaussian matrix is identified and the novel zeroing matrix is proposed with an improved performance. With the SD framework, the optimal sample allocation strategy for the block diagonal measurement matrix are derived for the wavelet representation of natural images. Extensive simulations validate the optimality of the proposed measurement matrix design.

Motivated by the zeroing matrix, we extend the seeded matrix design in the CS literature to the novel modulated matrix structure. The major advantage of the modulated matrix over the seeded matrix lies in the simplicity of its state evolution dynamics. Together with the AMP based algorithm, the modulated matrix possesses a 1-D performance prediction system, with which we can optimize the matrix configuration. We then focus on a special modulated matrix form, designated as the two block matrix, which can also be seen as a generalization of the zeroing matrix. The effectiveness of the two block matrix is demonstrated through both sparse and compressible signals. The underlining reason for the improved performance is presented through the analysis of the state evolution dynamics.

The final contribution of the thesis explores improving the reconstruction algorithm. By taking the signal prior into account, the Bayesian optimal AMP (BAMP) algorithm is demonstrated to dramatically improve the reconstruction quality. The key insight for its success is that it utilizes the minimum mean square error (MMSE) estimator for the CS denoising. However, the prerequisite of the prior information makes it often impractical. A novel SURE-AMP algorithm is proposed to address the dilemma. The critical feature of SURE-AMP is that the Stein's unbiased risk estimate (SURE) based parametric least square estimator is used to replace the MMSE estimator. Given the optimization of the SURE estimator only involves the noisy data, it eliminates the need for the signal prior, thus can accommodate more general sparse models.

Declaration of originality

I declare that this thesis was composed by myself and that the work contained therein is my own, except where explicitly stated otherwise in the text. This work has not been submitted for any other degree or professional qualification.

Chunli Guo

Acknowledgements

The four years PhD study has been a wonderful and overwhelming experience. I am indebted to many people who made this four years an unforgettable experience.

First of all, I offer my sincerest gratitude to my supervisor, Professor Mike Davies, for his support, patience, encouragement and advice throughout my entire PhD. I am extremely lucky to have Mike as my supervisor. His huge enthusiasm on research and fathomless knowledge in many areas has deeply inspired me. With his guidance, I have learnt to do research properly, to follow high scientific standards for my work and above all, to be an engineer. I fully appreciate his patience during the first year of my PhD, when I was so clueless about structured programming, scientific writing, let alone tackling research problems. I have lost count of the number of times when I came to him for all sorts of questions, and left of his office with a mind full of new insights. I wholeheartedly enjoyed the individual and group meetings with him for the past four years. His valuable feedback and expert counsel helped shape much of the work in this thesis. I am also very grateful for his generous financial support for sending me to all workshops and conferences, to meet all leading experts in the field, which definitely enriched the whole PhD experience.

I would like to thank Dr. Mehrdad Yaghoobi Vaighan for his valuable suggestion and insightful comments on both academic and administrative problems. I am grateful for all our discussion about solving various optimization problems. My thanks also go to Dr. Sinan Sinanović, who took the time and effort to find reference for my numerical problems.

Special thanks go to my awesome friends and colleagues Chaoran Du, Bogomil Shtarkalev, Dobroslav Tsonev, Abdelhamid Alhassi, Stefan Videv, Harald Burchardt, Nicola Serafimovski, Sakis Stavridis, Nataša Utješanović and Yuchang Wong. I enjoyed all the eye-opening discussion on every aspects of the life. You guys let me know the world better and made my life in Edinburgh full of happy and unforgettable moments. The past four years would not be so colorful without any of you.

I am forever grateful to Xiaohu for his endless love and sacrifice throughout this entire journey. Sharing my life with him for the past seven years is the sweetest and most unforgettable period of my life. Thank you for coming to UK and pursuing a PhD with me. Thank you for the

countless travelling from York to Edinburgh for the past four years. Most of all, thank you for being there for me for all my ups and downs in various stages of my life. My life would not be complete without you.

Last but not least, I would like to thank my parents for their unflagging love and support throughout my life, especially for not raising me up in a traditional way and restricting my life choices. To them I dedicate this thesis.

Contents

Declaration of originality	iv
Acknowledgements	v
Contents	vii
List of figures	x
List of tables	xiv
Acronyms and abbreviations	xv
Nomenclature	xvii
1 Introduction	1
1.1 Introduction	1
1.2 Original Contributions	3
1.3 Thesis Organization	4
1.4 Publications	6
2 Background	8
2.1 Introduction	8
2.2 Low Dimensional Signal Models	9
2.2.1 Sparse and Compressible Signals: Deterministic Model	9
2.2.2 Sparse and Compressible Signals: Stochastic Model	11
2.3 Sensing Matrices	13
2.4 Compressed Sensing Reconstruction	15
2.4.1 ℓ_1 -Minimization	15
2.4.2 Greedy methods	16
2.4.3 Approximate Message Passing Based Methods	17
2.5 Phase Transitions	20
2.6 Graphical Model for CS and AMP Derivation	21
2.6.1 Factor Graph and Sum-Product Algorithm Review	22
2.6.2 Relaxed Message Passing for CS	25
2.6.3 From Relaxed Message Passing to AMP	27
2.7 AMP Based Algorithm Summary	30
2.7.1 Bayesian optimal AMP	30
2.7.2 ℓ_1 -AMP	32
2.7.3 Extensions for AMP	33
2.8 State Evolution Dynamics	34
2.8.1 State Evolution Heuristics	35
2.8.2 State Evolution Formula	37
2.9 Summary	37
3 Sample Distortion Framework for Compressed Sensing	39
3.1 Introduction	39
3.2 Sample Distortion Framework	42
3.2.1 Definition	42
3.2.2 Lower bounds	43

3.2.3	Convex property	46
3.3	Measurement Matrix for Multi-resolution Image Model	48
3.3.1	Compressible Distributions	48
3.3.2	Band-wise Independent Image Model	49
3.3.3	Band-wise Sampling	53
3.4	Sample Allocation with Tree Structure	57
3.4.1	Hidden Markov Tree Model	57
3.4.2	Turbo Decoding	59
3.4.3	Sample Allocation with Tree Structure	61
3.5	Simulations	62
3.5.1	Sample Allocation with Oracle Image Statistics	63
3.5.2	Sample Allocation with General Image Statistics: The GSA	65
3.6	Summary	69
4	Modulated Matrix Design	71
4.1	Introduction	71
4.2	Seeded Matrix Review	73
4.2.1	Seeded Matrix Structure	74
4.2.2	When Does Seeding Work?	75
4.2.3	State Evolution for Seeded Matrix	75
4.3	Modulated Matrix Framework	76
4.3.1	Modulated Matrix Structure	76
4.3.2	1-D State Evolution	77
4.3.3	Two Block Matrix	79
4.4	First Order Phase Transition	80
4.4.1	Analysis via Statistical Physics	81
4.4.2	Analysis via State Evolution	84
4.4.3	FOPT Condition	86
4.4.4	Two Block Matrix Effect on FOPT	86
4.5	Simulations	88
4.5.1	Two Block Matrix for Sparse Signal	89
4.5.2	Two Block Matrix for a Compressible Signal	91
4.5.3	Two Block Matrix for Dense Signals	91
4.6	Summary	93
5	Bayesian optimal reconstruction without priors: parametric SURE-AMP algorithm	95
5.1	Introduction	96
5.2	Parametric SURE-AMP Framework	99
5.2.1	Parametric SURE-AMP algorithm	99
5.2.2	SURE based denoiser selection	100
5.2.3	State evolution	102
5.3	Construction of the Parametric Denoiser	104
5.3.1	Kernel families	105
5.3.2	Non-linear parameter tuning for kernel functions	107
5.3.3	Linear parameter optimization	109
5.3.4	Denoising performance	110
5.4	Numerical Results	112

5.4.1	Noisy signal recovery	112
5.4.2	Runtime comparison	116
5.5	Conclusion	117
6	Conclusions and Future work	120
6.1	Conclusion	120
6.2	Open Problem and Further Work	121
A	Approximation of the factor-to-variable Message	125
B	Derivation of Equation (2.53)	128
C	Relationship of Conditional Mean and Variance for Gaussian Corrupted Data	129
D	Proof of the Entropy Based Bound	131
E	Derivation of the Hierarchical Bayesian Model for GGD	133

List of figures

2.1	Compressible representation of the cameraman image via a multiscale wavelet transform and its best k-term approximation. (a) Original image. (b) db2 wavelet decomposition. Large coefficients are represented by light pixels, while small coefficients are represented by dark pixels. Note that most of the wavelet coefficients are close to zero. (c) Approximation of the image obtained by keeping the largest 10% of the wavelet coefficients. (d) Sorted wavelet coefficients in a descend order for each scale.	10
2.2	Solid lines illustrate the sorted magnitude of the db2 wavelet coefficients of the cameraman image in five different scale. Dashed lines show the expected order statistics of the GGD and GMD models with the parameters estimated directly from cameraman. (a) Wavelet/GGD. (b)Wavelet/GMD.	12
2.3	The QQ plot for the IST (a) and AMP (b) residual at the 10 th iteration. The linearity of the QQ plot indicates the Gaussian behaviour.	18
2.4	Theoretical phase transition for ℓ_1 -minimization [1, 2].	20
2.5	The factor graph associated with the probability distribution in (2.26). Empty circles represent variables $x_i, i \in [n]$ and $y_j, j \in [m]$. Squares correspond to measurement function a_j . $m_{i \rightarrow j}(x_i)$ and $m_{j \rightarrow i}(x_i)$ are messages representing interaction among nodes.	23
2.6	A factor graph featuring one variable node and one factor node, which is used to explain the updating rule for the sum-product algorithm.	24
2.7	General linear mixing dealt with GAMP algorithm	33
3.1	SD functions for GMD data $p(x) = 0.38 \mathcal{N}(0, 1.198) + 0.62 \mathcal{N}(0, 0.004)$ and lower bounds. The critical sampling ratio (defined later in page 47) to convexify this SD function is $\gamma_c = 0.61$	43
3.2	SD functions for GGD data $\alpha = 0.4, \sigma = 1$ and lower bounds. The critical sampling ratio (defined later in Section. 3.2.3.1)to convexify this SD function is $\gamma_c = 0.15$	44
3.3	Hybrid zeroing Gaussian matrix as the convex combination of a trivial decoder $\hat{x} = 0$ and a BAMP decoder Δ . Elements equal to 0 are represented with white blocks.	47
3.4	EBB for GGD model with $\alpha = 0.4$ (left most curve), $0.5, \dots, 1.0$ and $\alpha = 2$	48
3.5	Ten 512×1024 HDR Images. From left to right, top to bottom: Chapel, Dog, Pine, Sea, Man, Wedding, Hill, Penguin, Room, Sign.	49
3.6	GGD parameters for six wavelet bands of HDR images in Fig. 3.5.	50
3.7	Two-state GMD parameters for 6 wavelet bands of HDR images in Fig. 3.5.	51
3.8	Nine 256×256 Natural Images: (a)Concordant (b) football (c) Gantry Crane (d) M83 (e) Spine (f) Kids (g) Rice (h) Peppers (i) Cameraman	52
3.9	GGD parameters for 6 wavelet bands of natural images in Fig. 3.8.	52
3.10	Distortion reduction function of six bands Daubechies 2 wavelet decomposition of cameraman image using GMD model (including the low-pass band). The statistics is reported in Table 3.1 in page 63.	55

3.11	(a) An illustration of the image quad-tree structure. (b) A zoomed in factor graph of the HMT structure featuring a typical variable node (the hidden state) $s_{j,k}$ connected with its four children $\{s_{j+1,c_{ki}}\}_{i=1}^4$ and parent node s_{j-1,p_k} by the factor nodes (the transition matrix).	58
3.12	Top: factor graph for the compressive imaging with HMT structure. Bottom: two sub-graphs for the turbo decoding. The likelihood from one sub-graph is used as the prior for the other sub-graph.	60
3.13	Factor graph for band-wise sampling with HMT decoding. The upper graph illustrates a quad-tree structure of the wavelet hidden states. The lower graph is the band-wise independent random mixing.	62
3.14	Sample allocation per band for Daubechies 2 wavelet with the GMD model. SA: sample allocation based on the bandwise independent model. HSA: sample allocation based on the empirical SD functions for BAMP decoder with soft information. ESA: empirically optimized sample allocation for turbo decoding.	64
3.15	PSNR comparison of different encoder-decoder pairs for cameraman Daubechies 2 wavelet with the GGD model. The lines are theoretical predictions with the SD function. While dots represent simulations with the cameraman image.	66
3.16	PSNR comparison of different encoder-decoder pairs for cameraman Daubechies 2 wavelet with the GMD model. The lines are theoretical predictions with the SD function. While dots represent simulations with the cameraman image.	66
3.17	Reconstruction using 10000 (15%) samples of the 256×256 cameraman image with different encoder-decoder pairs. The GMD is used to model the Daubechies 2 wavelet coefficients statistics. The encoding matrices for the cameraman simulations are explained in details in Sec. 3.5.1.	67
3.18	Ten test images from the Berkeley dataset [3]. From left to right, top to bottom are: car, plane, eagle, sculpture, surfer, tourists, building, castle, man and fish.	68
4.1	Construction of the seeded matrix for compressed sensing [4].	74
4.2	Construction of the modulated matrix. Gaussian random elements with different variances are indicated by different shade.	77
4.3	Construction of the two block matrix.	79
4.4	The potential function $\Lambda(E)$ for different signals at various sampling ratios with the homogeneous Gaussian matrix: (a) Bernoulli-Gaussian data with FOPT, $p_{BG}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\delta(x)$; (b) Gaussian-mixture data with FOPT, $p_{GM_1} = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 5e - 4)$ (c) Gaussian-mixture data without FOPT, $p_{GM_2} = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 0.003)$. The diamond-shaped dots represent the global maximums, while the star-shaped dots are the secondary local maximums.	82
4.5	The schematic plot of three types of SE behaviour to explain FOPT. The dash line is the baseline $\gamma\tau^t$. The solid lines are $S(\tau^t)$ for BAMP with the homogeneous Gaussian matrix. The number of non-zero intersection points with the baseline varies for different types of signals.	84

4.6	Fixed points of the SE evolution for both homogeneous Gaussian matrix and the two block matrix. (a) The compressible prior $p_{\text{GM}_2}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 3 \times 10^{-3})$ with $\gamma = 0.58$. For the homogeneous Gaussian matrix, the SE function has only one non-zero fixed point at $\tau_1^* = 0.1181$. Applying the two block matrix $J_2 = 1e - 3$, $\gamma_1 = 0.9206$ will not alter the shape of the SE function, it only shrinks the function so that the fix point is moved to $\hat{\tau}^* = 0.0222$. (b) The sparse prior $p_{\text{BG}}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\delta(x)$ with $\gamma = 0.55$. For the homogeneous Gaussian matrix, the SE function has two non-zero fixed points at $\tau_{3,1}^* = 0.1619$ and $\tau_{3,2}^* = 0.01020$. With the two block matrix $\gamma_1 = 0.847$, $J_2 = 10^{-3}$, the SE evolution successfully removes the spurious fixed points and leads to perfect reconstruction. (c) The compressible signal is $p_{\text{GM}_1}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 5 \times 10^{-4})$ with $\gamma = 0.58$. For the homogeneous Gaussian matrix, the SE equation has three non-zero fixed points at $\tau_{2,1}^* = 0.1006$, $\tau_{2,2}^* = 0.017$ and $\tau_{2,3}^* = 3.4 \times 10^{-3}$. With the two block matrix $\gamma_1 = 0.847$, $J_2 = 10^{-3}$, the fix point is moved to $\hat{\tau}^* = 0.0008$	87
4.7	The normalized SD function for the sparse signal $p_{\text{BG}}(x)$ with different measurement matrix configuration. For two-block matrix, $\gamma_1 = \gamma/\gamma_c$ for $\gamma < \gamma_c$, where $\gamma_{c_1} = 0.59$ is the perfect reconstruction ratio for the homogeneous Gaussian matrix. The three-block matrix is the achieved by convexify the SD function of the two block matrix with $\gamma_1 = \frac{\gamma}{\gamma_{c_1}}$, $\gamma_2 = \frac{\gamma}{\gamma_{c_2}} - \frac{\gamma}{\gamma_{c_1}}$, $\gamma_3 = 1 - \gamma_2 - \gamma_3$, where $\gamma_{c_2} = 0.45$ is the perfect reconstruction ratio achieved by the two block matrix.	89
4.8	The normalized SD function for the compressible signal $p_{\text{GM}_1}(x)$ with different measurement matrix configuration. For two-block matrix, $\gamma_1 = \gamma/\gamma_c$ for $\gamma < \gamma_c$, where $\gamma_c = 0.6$ is the critical sampling ratio for the homogeneous Gaussian matrix.	90
4.9	The normalized SD function for the compressible signal $p_{\text{GM}_2}(x)$ with different measurement matrix configuration. For the two-block matrix, $\gamma_1 = \gamma/\gamma_c$ for $\gamma < \gamma_c$, where $\gamma_c = 0.63$ is the critical sampling ratio for the homogeneous Gaussian matrix.	92
4.10	The SD function for the dense signal with different measurement matrix configuration. With $\gamma_1 = 0.6$, $J_2 = 10^{-5}$, the two block matrix is able to move the perfect reconstruction sampling ratio from 0.49 to 0.29.	93
5.1	QQ plots tracking the effective noise of the AMP algorithm under various iterations while reconstruction a 40% sampled Bernoulli-Gaussian data with pdf in (5.35). The residual of AMP remains Gaussian because of the Onsager reaction term. Decreasing slope as the iteration increasing indicates the decreasing standard deviation.	96
5.2	The actual MSE for the noiseless Bernoulli-Gaussian data reconstruction at each parametric SURE-AMP iteration versus the state evolution prediction. The signal is generated i.i.d. according to eq. (5.35). The first piecewise linear kernel family is utilized within the parametric SURE-AMP algorithm, which will be discuss in section 5.3.1.1. The reconstruction MSE is an average over 100 Monte Carlo realizations.	103
5.3	Kernel families used for linear parameterization of the SURE based denoiser: (a) the first piecewise linear kernel family (b) The second piecewise linear kernel family. (c) The exponential kernel family.	104

5.4	MMSE estimator and parametric SURE for the noisy Bernoulli-Gaussian data. The noise variance c is 0.1. The reconstruction error for the MMSE estimator, the SURE estimator with the first piecewise linear kernel and the SURE estimator with the exponential kernel are 0.020615, 0.020788 and 0.022047, respectively.	110
5.5	MMSE estimator and parametric SURE for the noisy k -dense data. The noise variance c is 0.1. The reconstruction error for the MMSE estimator, the SURE with the second piecewise linear kernel, the SURE estimator with the first piecewise linear kernel and the IDR denoiser are 0.0243 and 0.0248, 0.0251 and 0.0315 respectively.	111
5.6	SNR_x versus sampling ratio for CS recovery of noisy Bernoulli-Gaussian data.	114
5.7	SNR_x versus sampling ratio for CS recovery of noisy k -dense data.	115
5.8	SNRx versus sampling ratio for CS recovery of noisy student-t data.	116
5.9	Runtime versus signal dimension for CS recovery of noisy Bernoulli-Gaussian data.	118
5.10	Runtime versus signal dimension for CS recovery of noisy k -dense data.	119
5.11	Runtime versus signal dimension for CS recovery of noisy student-t data.	119

List of tables

3.1	Statistics for Daubechies 2 wavelet coefficients of cameraman	63
3.2	Average Statistics for Daubechies 2 wavelet coefficients of 200 test images from the Berkeley dataset [3]	68
3.3	Image reconstruction results for ten 256×256 test images from the berkeley image database [3] with $\gamma = 0.1$. Entries are the peak signal-to-noise ratio (PSNR) in decibels, $\text{PSNR} := 10 \log_{10}(N/\ \hat{x} - x\ _2^2)$. All results use the aveargae image statistics reported in Table 3.2 and the BAMP decoder.	69
3.4	Reconstruction PSNR for test images with $\gamma = 0.2$	70
3.5	Reconstruction PSNR for test images with $\gamma = 0.3$	70
5.1	Denoising comparison for noisy Student's-t signal with various denoisers . . .	112

Acronyms and abbreviations

1-D	One Dimensional
AMP	Approximate Message Passing
AWGN	Additive White Gaussian Noise
BAMP	Bayesian optimal Approximate Message Passing
BP	Basis Pursuit
BPDN	Basis Pursuit Denoising
BG	Bernoulli-Gaussian
CoSamp	Compressive Sampling matching pursuit
CS	Compressed Sensing
CT	Computed Tomography
D-AMP	Denoising-based Approximate Message Passing
DCT	Discrete Cosine Transform
DOG	Derivatives of Gaussians
DR	Distortion Reduction
EBB	Entropy Based Bound
EM	Expectation Maximization
ESA	Empirically optimized Sample Allocation
FFT	Fast Fourier Transform
FOPT	First Order Phase Transition
GAMP	Generalized Approximate Message Passing
GGD	Generalized Gaussian Distribution
GM	Gaussian Mixture
GMD	Gaussian Mixture Distribution
GSA	General Sample Allocation
HDR	High Dynamic Range
HMT	Hidden Markov Tree
HSA	HMT based Sample Allocation
IDR	Iterative Dense Recovery
IHT	Iterative Hard Thresholding
i.i.d	independent identically distributed

InforSA	Informative Sensing Matrix
IST	Iterative Soft Thresholding
KT	Kuhn-Tucker
LASSO	Least Absolution Shrinkage and Selection Operator
LBP	Loopy Belief Propagation
LDPC	Low Density Parity Check
LP	Linear Programming
MBB	Model Based Bound
MBSA	Multi-scale Sensing Matrix
MCMC	Markov-chain Monte-Carlo
ML	Maximum-Likelihood
MRI	Magnetic Resonance Imaging
MSE	Mean Squared Error
MMSE	Mimimum Mean Squared Error
OMP	Orthogonal Matching Pursuit
PAS	Persistence Across Scale
pdf	probability density function
PSNR	Peak Signal-to-Noise Ratio
QQ	Quantile-Quantile
RIP	Restricted Isometry Property
SA	Sample Allocation
s-BP	seeded Belief Propagation
SD	Sample Distortion
SE	State Evolution
StOMP	Stagewise Orthogonal Matching Pursuit
SURE	Stein's Unbiased Risk Estimate
TAP	Thouless-Anderson-Palmer

Nomenclature

$\cdot, (\cdot)$	operand, which can be scalar, vector or matrix
$\hat{\cdot}$	Estimated symbol
$(\cdot)^T$	transpose
$\langle \cdot \rangle$	mean of an vector
$\langle \cdot, \cdot \rangle$	inner product of two vectors
\approx	approximately equal
$ \cdot $	absolute value
$\mathbb{E}\{\cdot\}$	expectation operator
$\text{Var}\{\cdot\}$	variance operator
$\mathcal{N}(\cdot)$	Gaussian distribution
$\delta(\cdot)$	Dirac delta function
Δ	General CS reconstruction algorithm or decoder
$\ \cdot\ _0$	ℓ_0 norm
m	number of CS measurements
n	signal length
Φ	measurement matrix
Φ_i	i th column of the matrix
$\Phi_{i,j}$	element of the matrix
γ	sampling ratio
ρ	sparsity ratio for sparse signals
λ	density ratio for compressible signals
\mathbf{x}	signal vector
\mathbf{y}	measurement vector
ξ	additive white Gaussian noise vector

Chapter 1

Introduction

1.1 Introduction

Compressed sensing (CS) has become a popular topic for signal processing since around 2004. As a non-conventional technique, it advocates acquiring a compressed representation of a signal at a sub-Nyquist rate during the sampling stage and obtaining a faithful recovery with the computational power afterwards. A re-examination of the “arms race” between the camera manufactures and the image compression software engineers is often used to explain the motivation of the CS paradigm [5]: while the hardware engineers are passionate about making multi-megapixel cameras, the software engineers are racking their brains for developing clever algorithms for image compression, because storing and transmitting the original enormous computer files are often impractical. Moreover, with proper compression algorithms, it is usually impossible to tell the differences between the original and compressed images with the naked eye. Since only a few data is required to adequately describe a signal and most of the finely sampled data would end up being discarded, one would like to have the data compression built directly into the data acquisition procedure. That is what CS mainly about.

The group testing example is also frequently referred to explain the CS concept in layman’s terms. The group testing problem appears in many forms, one of which dates back to World War II when a huge blood sample population needs to be tested for syphilis where very few patients had the disease. Instead of performing the individual and often expensive testing, the blood samples were partitioned into groups and mixed according to some pre-designed rule. One measurement per group was then obtained for testing. Given the low possibility of infection, this strategy largely reduced the number of testing required thus the cost. In the CS context, the blood testing results for all samples are the original signal. The measurements are the samples drawn from the mixed blood. The design of the blood mixing rule plays the same role as the measurement matrix in compressed sensing, in the sense that it enables an efficient set of test results containing enough information to determine a small subset of items of interest [6].

The first step in CS is the sensing mechanism to obtain the information of a signal $\mathbf{x} \in \mathbb{R}^n$,

which can be written mathematically as

$$y_i = \langle \Phi_i, \mathbf{x} \rangle, \quad i = 1, \dots, m \quad (1.1)$$

where $\langle \cdot, \cdot \rangle$ is the inner product operator. That is, m measurements of the signal \mathbf{x} are obtained by taking the inner product with the sensing vectors Φ_1, \dots, Φ_m . For the standard CS setup, the number of measurements is much less than the signal dimension $m \ll n$. Assembling all sensing vectors together, we will have an under-determined linear system

$$\mathbf{y} = \Phi \mathbf{x} \quad (1.2)$$

where $\mathbf{y} = [y_1, \dots, y_m]$ is the observation vector and $\Phi \in \mathbb{R}^{m \times n}$ is the measurement or sensing matrix with the vectors $\Phi_1^T \dots \Phi_m^T$ as the rows. The CS reconstruction task is to obtain a unique solution $\hat{\mathbf{x}}$ that matches the observed vector \mathbf{y} and some additional prior information.

The seemingly magic power of compressed sensing to obtain original signals with a sub-Nyquist rate relies on two principles. Firstly, it exploits the fact that most signals that we are interested depend on a much smaller number of degrees of freedom than its bandwidth or signal length suggests. Given a proper basis Ψ , they can have a concise representation with a few numbers without losing much information. Such signals are said to be sparse or compressible and are the targets for CS techniques.

Secondly, the sensing vectors Φ_i must have a dense representation in the sparse basis Ψ [7]. A example of signal with dense representation is a Dirac function spreading out in the frequency domain. This is where the randomness enters the picture. With high probability, a measurement matrix with independent identically distributed (i.i.d.) Gaussian entries is largely incoherent with any fixed Ψ . Thus the homogeneous Gaussian measurement matrix is widely considered in CS works. However, it is not necessarily the optimal choice for CS reconstruction. An interesting research avenue is to design the structured measurement matrix to obtain better reconstruction with few measurements. For practical signals, one would expect more properties than just sparsity. Tailoring the measurement matrix with the additional prior information would also benefit reconstruction. In this thesis, designing and optimizing the structured measurement matrix in accordance with the original signal form a major contribution.

Another important ingredient for CS is the reconstruction algorithm. For signal recovery, CS leverages the highly nonlinear methods. The conventional tractable algorithms include the lin-

ear programming (LP) and greedy methods. Extensive research for the reconstruction performance and convergence analysis for both types of algorithms can be found in the CS literature. In 2009, Donoho and co-authors introduced the novel approximate message passing (AMP) algorithm, which utilized the graphic model approach for the CS reconstruction. As an iterative algorithm, AMP is particularly of interest in two aspects. First of all, a distinguishable feature of AMP is the state evolution (SE) formalism, which can accurately predict the asymptotic algorithm behaviour in the large system. Secondly, AMP is able to incorporate the signal prior information and thus deliver improved recovery in comparison with the conventional algorithms. As a relatively new CS algorithm, there are a lot unanswered questions and possible applications to be explored. In this thesis, AMP is used as the primary reconstruction tool. The SE dynamics are deployed for optimizing the structured measurement matrix. Modification and enhancement of the generic AMP algorithm is also investigated and leads to novel AMP based algorithms in this thesis. More recent study reveals the recursive Gaussian denoising nature of the AMP reconstruction. It means that there are opportunities to marry various off-the-shelves denoising methods with the AMP framework to solve practical signal processing problems.

1.2 Original Contributions

The contribution in this thesis is twofold: the structured measurement matrix design and enhancing the generic AMP algorithm. In particular, three main points are listed below to give a short overview.

- **Optimized sample allocation for the compressed imaging**

The non-optimality of the homogeneous Gaussian matrix for the compressed imaging has been identified in the literature for a long time. Driven by the need for an analytical bandwise sampling scheme, we establish a sample distortion (SD) function for the wavelet multi-resolution image model and introduce a tractable sample allocation method assuming the independence of the wavelet bands. Essentially we address the following problem: given a fixed number of CS measurements, how many samples should be allocated for each wavelet band to achieve the optimal reconstruction. To our knowledge, the work presented in this thesis is the first analytical result for optimizing the bandwise sampling of CS imaging. Furthermore, the novel sample distortion framework provides us with the accurate prediction for the reconstruction error associated with the optimized measurement matrix.

- **Modulated matrix design with the one dimensional state evolution dynamics**

Apart from the bandwise independent measurement matrix, the spatial coupling structure has also been applied for the CS reconstruction to reduce the recovery error. The main contribution in this respect is the introduction of the modulated matrix design, which can be seen as a concatenation of Gaussian random sub-matrices with different variances. In this thesis we show that such measurement matrix is able to achieve the perfect reconstruction with a sampling ratio approaching the theoretical limit for sparse signals. For compressible signals, it also offers improved reconstruction quality. More importantly, we introduce a simple one dimensional (1-D) state evolution formalism to characterize the AMP behaviour with the modulated matrix, with which we can predict the reconstruction error and optimize the matrix configuration.

- **Parametric SURE-AMP algorithm with the Bayesian optimal reconstruction**

The final main contribution of the thesis lies in the introduction of the novel parametric SURE-AMP algorithm, which is a variant for the generic AMP algorithm. Leveraging the intrinsic signal denoising nature of the AMP iterations, an adaptive parametric denoising module is introduced to the AMP framework. At each iteration, the denoiser is optimized by minimizing the Stein's unbiased risk estimate (SURE) of the recovered signal. Since SURE is the unbiased estimate of the mean squared error (MSE), the parametric SURE-AMP progresses by directly minimizing the least squared error of the reconstructed signal. The proposed parametric SURE-AMP algorithm improves the generic AMP performance and is able to achieve the Bayesian optimal recovery as if the true signal prior is known for reconstruction.

1.3 Thesis Organization

The rest of the thesis is organised as follows:

Chapter 2 presents the background information related to the topic of this thesis. It starts with a summary of the key concepts and theoretical results of compressed sensing. This is followed by a detailed overview of AMP, from the derivation of the algorithm to a summary of the AMP variants. This chapter finishes with a short overview of the state evolution dynamics for the AMP algorithms. Overall the AMP algorithms have been used throughout the work for both the measurement matrix and the reconstruction algorithm design. Thus this chapter lays the foundation for the rest of the thesis.

Chapter 3 establishes the sample distortion framework for quantitatively evaluating the reconstruction error for certain pairs of measurement matrix and recovery scheme. The intrinsic convex property of the SD function leads to the hybrid zeroing matrix design. After a brief discussion of the bandwise independent wavelet model, the SD framework is applied to natural images to optimize the sample allocation for the block diagonal measurement matrices. With the convexified SD function for the wavelet image model, a reversed water-filling scheme is applied to achieve the optimal sample allocation. Finally, the wavelet tree structure is incorporated with the bandwise sampling and a more general measurement matrix for natural images based on the average image statistics is derived.

Chapter 4 introduces a novel dense measurement matrix to improve the SD performance for both sparse and compressible signals. The proposed matrix, designated as the modulated matrix, is inspired by the hybrid zeroing matrix presented in Chapter 3 and the seeded matrix in the literature. After a description of the matrix structure, a simple 1-D SE equation is derived to characterize its asymptotic behaviour when used together with the AMP algorithm. Under the modulated matrix framework, the two block matrix is then presented as a special realization. Finally the chapter concludes with a simulation for both sparse and compressible signals to illustrate the effectiveness of the modulated matrix design.

Chapter 5 presents the novel parametric SURE-AMP algorithm. The Gaussian behaviour of the AMP residual and the recursive denoising nature of AMP are first revisited in this chapter. The pros and cons of the existing AMP based algorithms are also analysed to motivate the parametric SURE-AMP algorithm. The proposed algorithm incorporates an adaptive denoising function family with the AMP iteration and select the denoiser with the minimum mean squared error (MMSE) at each step. Three different kernel families are then introduced as the base functions to form the denoisers. Simulation with both sparse and compressible signals demonstrate that with proper design of the denoiser family, the parametric SURE-AMP algorithm is able to deliver state-of-the-art performance in terms of both reconstruction quality and computational complexity.

Chapter 6 concludes this thesis with a discussion of the limitations of the presented work and the directions for potential future research.

1.4 Publications

Work presented in this thesis has previously submitted or published in the peer reviewed journals and conference proceedings. A full list of the publications is as follows,

Peer Reviewed Journal Articles:

1. Chunli Guo and Mike E. Davies, “Near optimal compressed sensing without priors: Parametric SURE Approximate Message Passing” is accepted by *IEEE Transactions on Signal Processing*.

Part of this paper has found its way into the background in Chapter 2 and the whole of Chapter 5 has been taken from this publication.

2. Norbert Goertz, Chunli Guo, Alexander Jung, Mike E. Davies and Gerhard Doblinger, “Iterative recovery of dense signals from incomplete measurements”, in *IEEE Signal Processing Letters*, vol 21, pp.1059-1063, 2014.

The k -dense signal model from this publication is introduced in Chapter 2 and used in both Chapter 4 and Chapter 5. This paper was a joint paper and the MATLAB code for iterative dense recovery has been supplied by the first author of the paper.

3. Chunli Guo and Mike E. Davies, “Sample distortion for compressed imaging”, in *IEEE Transactions on Signal Processing*, vol. 61, No. 24, pp 6431-6442, 2013.

The work presented in Chapter 3 has mainly been taken from this paper.

Conference Proceedings:

1. Chunli Guo and Mike E. Davies, “Bayesian optimal compressed sensing without priors: parametric SURE approximate message passing”, in *Proc. European Signal Processing Conference (EUSIPCO)*, September 2014.

This publication contributes part of Chapter 5.

2. Chunli Guo and Mike E. Davies, “Modulated measurement matrix design for compressed sensing”, in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, May 2014.

This paper presented the work that can be found in Chapter 4 of this thesis.

3. Chunli Guo and Mike E. Davies, “Sample allocation for statistical multiresolution com-

pressed sensing”, in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, May 2013.

This paper discuss the issues that can be found in the Chapter 3.

4. Mike E. Davies and Chunli Guo, “Sample-distortion functions for compressed sensing”, in *Annual Allerton Conference on Communication, Control and Computing* (Allerton), 2011 (Invited Paper).

This early publication contributes to the theoretical work in Chapter 3.

Chapter 2

Background

2.1 Introduction

The compressed sensing problem we consider throughout this thesis is formulated as the following under-determined linear system

$$\mathbf{y} = \Phi \mathbf{x} + \boldsymbol{\xi} \quad (2.1)$$

where $\mathbf{x} \in \mathbb{R}^n$ is the sparse or compressible signal that we would like to reconstruct, $\mathbf{y} \in \mathbb{R}^m$ is the CS observation vector, $\Phi \in \mathbb{R}^{m \times n}$ is the measurement matrix with the sampling ratio $\gamma = \frac{m}{n} < 1$, and $\boldsymbol{\xi} \in \mathbb{R}^m$ is the additive white Gaussian noise (AWGN) vector with i.i.d. entries $\xi_i \sim \mathcal{N}(0, \sigma_\xi^2)$. The noiseless case is incorporated in this model by setting $\boldsymbol{\xi} = 0$. Given Φ , \mathbf{y} and the signal prior $p(\mathbf{x})$ in some scenario, the goal of CS is to reconstruct \mathbf{x} as best as possible through appropriate algorithm and measurement matrix design.

In this chapter, we review some basic knowledge of compressed sensing, with the emphasis on the AMP algorithm. We start with one of the most important principles that compressed sensing relies on, the low dimensional signal structure. Then the CS measurement matrix properties are shortly discussed. Three different types of existing CS algorithms are briefly summarized. We then focus on the presentation of the AMP related information. Since AMP is derived from the canonical message passing algorithm over the graphic model for CS, the key concept of the factor graph and the sum-product algorithm are reviewed. It is followed by a detailed re-derivation of the Bayesian optimal AMP (BAMP) algorithm as a working example. Three types of AMP-based algorithms: BAMP, ℓ_1 -AMP and generalized AMP (GAMP) are summarized. Finally one of the most distinguishable features of AMP-base algorithms, the state evolution formalism, is discussed with the basic intuition and detailed formula. It will be shown later in Chapter 3 that the SE dynamics provides a theoretical basis for the sample distortion framework. Further it can be used as a tool for optimizing the measurement matrix configuration as shown in Chapter 4.

2.2 Low Dimensional Signal Models

As stated in Chapter 1, the success of compressed sensing relies heavily on the fact that the number of degrees of freedom for high-dimensional signals is often much smaller than their ambient dimensionality. In this section, we explain the most common low-dimensional structures encountered in compressed sensing from both a deterministic and stochastic perspective.

2.2.1 Sparse and Compressible Signals: Deterministic Model

Signals can often be well-approximated by a linear combination of just a few elements from a known basis. When the approximation is exact, we say the signal is sparse. From the deterministic perspective, the sparsity is often quantified by the ℓ_0 -norm in the CS literature. Suppose $\mathbf{x} \in \mathbb{R}^n$ is the signal to be acquired, we say \mathbf{x} is k -sparse when it has at most k non-zero components

$$\|\mathbf{x}\|_0 \leq k, \quad k < n \quad (2.2)$$

Typically in CS we are dealing with signals that are not sparse in the time domain but a transformed domain. Suppose \mathbf{x} can be expressed as a linear combination of $\boldsymbol{\theta} \in \mathbb{R}^n$ in some orthonormal basis $\Psi \in \mathbb{R}^{n \times n}$, which is $\mathbf{x} = \Psi\boldsymbol{\theta}$, we still refer to \mathbf{x} as k -sparse if $\|\boldsymbol{\theta}\|_0 \leq k$.

In practice, few real-world signals are exactly sparse. Most of them are only well approximated by a sparse signal. Such signals are denoted as compressible signals. A typical compressible signal example is the natural image represented with a multi-resolution wavelet transform. As shown in Fig. 2.1, most of the wavelet coefficients of the cameraman image are so small that we can hardly tell the difference between the original and the approximated image, which is obtained by setting the small coefficients to zero. This procedure yields the best k -term approximation of the image, i.e. the best approximation of the signal using only k basis elements.

One possible definition of compressible signal is the one whose coefficients, when sorted in a descending order, satisfies the following inequality

$$|x_i| \leq ci^{-q} \quad (2.3)$$

where $c, q > 0$ are constants [20]. Fig. 2.1(d) displays the sorted wavelet coefficients for each wavelet scale of the cameraman image in the log-log scale. It is clear that its wavelet coefficients are compressible within each wavelet scale.

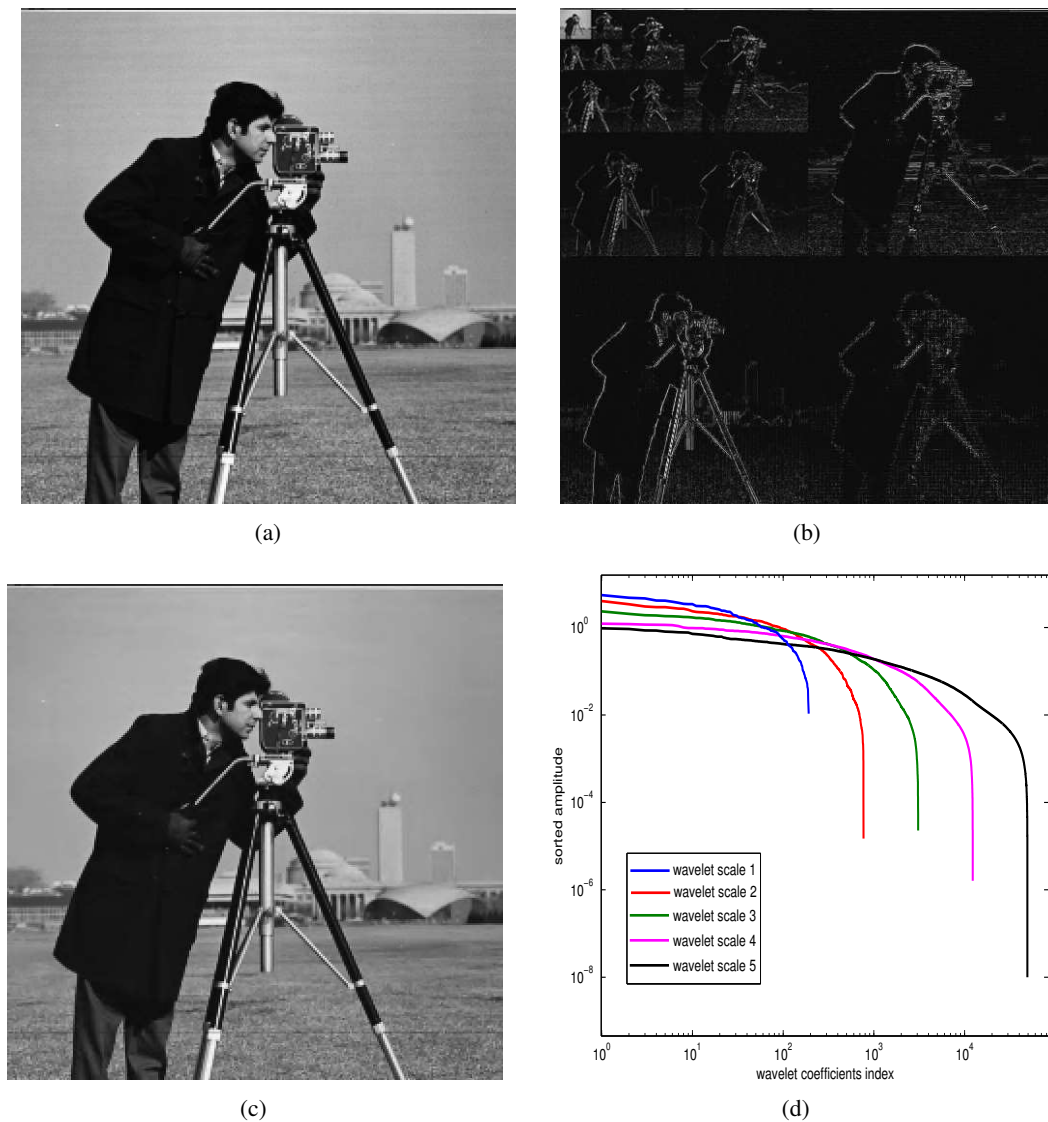


Figure 2.1: Compressible representation of the cameraman image via a multiscale wavelet transform and its best k -term approximation. (a) Original image. (b) db2 wavelet decomposition. Large coefficients are represented by light pixels, while small coefficients are represented by dark pixels. Note that most of the wavelet coefficients are close to zero. (c) Approximation of the image obtained by keeping the largest 10% of the wavelet coefficients. (d) Sorted wavelet coefficients in a descend order for each scale.

2.2.2 Sparse and Compressible Signals: Stochastic Model

When considering the CS reconstruction problem in the stochastic setting, probabilistic Bayesian models to characterize the signal sparsity/compressibility are naturally required. To be specific, we seek distributions whose i.i.d. realizations are strictly sparse or can be well approximated as sparse.

Based on the deterministic description of sparse signals, it is straightforward to model sparsity with the following distribution

$$p_X(x_i) = \rho F(x_i) + (1 - \rho)\delta(x_i) \quad (2.4)$$

where $\delta(\cdot)$ is the Dirac delta function and $F(\cdot)$ characterizes the statistical property of non-zero coefficients. Given the form in (2.4), the signal sparsity level is invariant to $\Gamma(\cdot)$ but controlled by the sparsity ratio $\rho = \frac{k}{n}$. In this thesis, we broadly use the Bernoulli-Gaussian (BG) distribution as an exemplary sparse signal model with $\Gamma(\cdot)$ being the Gaussian distribution.

For compressible signals, the appropriate distribution should be 'peaky' around zero to capture the concentration of small magnitude components and have heavy tail to represent the large magnitude components. In [8], a specific definition of compressible distribution is given and the way of identifying compressible distributions is discussed. Here we present two specific non-Gaussian distributions that we will use in this thesis to model compressible signals.

First, a popular probabilistic model for heavy-tailed non-Gaussian distributions is the generalized Gaussian distribution (GGD) [9, 10]. The pdf for the GGD can be written as

$$p_{\text{GGD}}(x) = \frac{\alpha}{2\sqrt{\beta}\sigma\Gamma(\frac{1}{\alpha})} \exp\left(-\left|\frac{x}{\sqrt{\beta}\sigma}\right|^\alpha\right) \quad (2.5)$$

where $\beta = \Gamma(1/\alpha)/\Gamma(3/\alpha)$, σ is the standard deviation and α is the shape parameter. As α goes to zero the distribution has increasingly heavy tails. For the special cases of $\alpha = 1$ we have the Laplace distribution and when $\alpha = 2$ we have the Gaussian distribution. The GGD provides a good approximation to the distribution of the wavelet coefficients for natural images (at a fixed wavelet scale) with $\alpha \sim [0.3, 1]$.

Another commonly used distribution to model compressible signals is the two-state Gaussian

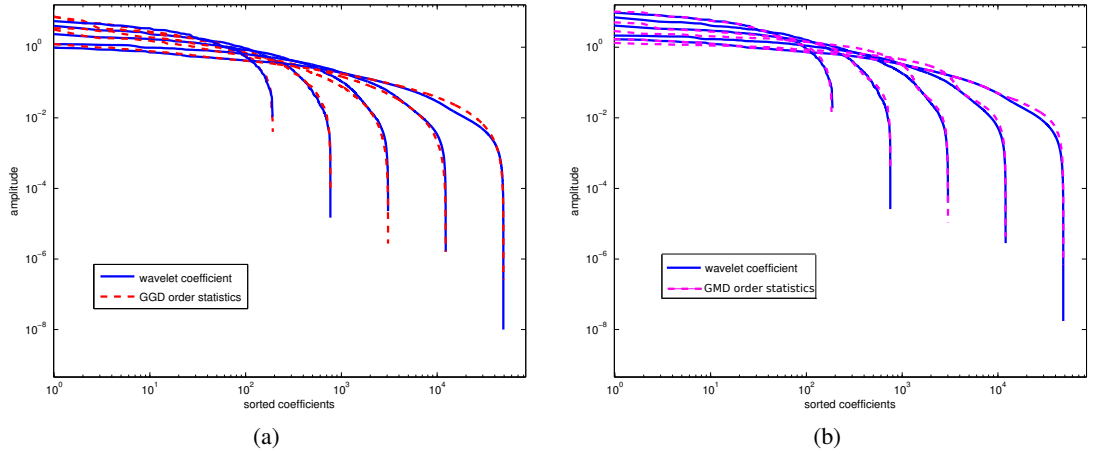


Figure 2.2: Solid lines illustrate the sorted magnitude of the db2 wavelet coefficients of the cameraman image in five different scale. Dashed lines show the expected order statistics of the GGD and GMD models with the parameters estimated directly from cameraman. (a) Wavelet/GGD. (b) Wavelet/GMD.

mixture distribution (GMD). The pdf for the GMD is written as

$$\begin{aligned}
 p_{\text{GMD}}(x) &= p(x|s=1) + p(x|s=0) \\
 &= p(s=1)\mathcal{N}(x; 0, \sigma_L^2) + p(s=0)\mathcal{N}(x; 0, \sigma_S^2) \\
 &= \lambda\mathcal{N}(x; 0, \sigma_L^2) + (1-\lambda)\mathcal{N}(x; 0, \sigma_S^2)
 \end{aligned} \tag{2.6}$$

where $s = \{0, 1\}$ are the hidden states, σ_L^2 and σ_S^2 are the large and small Gaussian variance, respectively. The density ratio for compressible signals is λ , which represents the portion of the significant elements. The two-state GMD model is quite effective at capturing the heavy tailed nature of compressible signals by adjusting λ . A random vector with i.i.d. two-state GMD components can be seen as generated either from the small or large variance Gaussian distribution, depending on the hidden states s . Since coefficients with small magnitude are expected to dominate the signal domain for compressive signals, we normally observe $\lambda < 0.5$.

In Fig. 2.2, we plot the magnitude-ordered wavelet coefficients for the cameraman image in the log-log scale. For comparison, we also present the expected order statistics¹ of both GGD and GMD models. The specific parameters for both distributions are estimated directly from the wavelet coefficients of cameraman via moment matching. It is clear that both GGD and GMD are able to well capture the marginal statistical properties of the image.

¹The i.i.d realizations of signal models are generated and expected magnitudes of the signal coefficients are sorted in a descending order [11]

Finally, we introduce another canonical CS signal model, the k -dense distribution. As opposed to the k -sparse concept, the k -dense signal has most of its elements taking their value from the discrete set $\mathcal{D} \equiv \{-\beta, +\beta\}$ with β being a real positive constant. The remaining k elements are real valued and taken from the open continuous set $\mathcal{C} \equiv (-\beta, +\beta)$. The pdf of the k -dense signal is written as

$$p_{\text{KD}}(x) = \frac{1 - \lambda_k}{2} \delta(x + \beta) + \frac{1 - \lambda_k}{2} \delta(x - \beta) + \lambda_k \mathcal{U}(-\beta, \beta) \quad (2.7)$$

where $\lambda_k = \frac{k}{n}$ and \mathcal{U} represents the pdf of the continuous components. In this thesis, we consider \mathcal{U} being the uniform distribution. The k -dense signal may stem from a source with real components that are clipped. Another example is a binary modulated signal received at a relay in cooperative communication [12]. In the CS literature, the k -dense signal has been considered before as the k -simple signal in [13]. The face counting theory has been established to bound the minimum sampling ratio for its perfect reconstruction via the convex optimization. In [14], the soft thresholding function with the adaptive thresholding level is suggested for the AMP algorithm. In this thesis, we will take it as a special non-sparse signal model to test our measurement matrix design and the proposed CS reconstruction algorithm in Chapter 4 and Chapter 5.

The concept of sparsity/compressibility of a signal cannot be discussed without the representation domain. With the orthonormal basis *Psi* fixed, the sparsity/compressibility can be checked by the definition. However, finding the sparse basis may not be trivial.

2.3 Sensing Matrices

As stated in Chapter 1, one of the main contributions in this thesis is the measurement matrix design. Before embarking on that topic, it is important to understand the CS measurement matrices properties that preserve the signal information to enable practical algorithms to accurately and efficiently recover the original signal. A key measurement matrix condition, used to study the general system's robustness, is known as the restricted isometry property (RIP) [15].

Definition 1. A matrix Φ satisfies the RIP of order k if for all k -sparse vectors \mathbf{x} there exists a constant $0 < \delta_k < 1$ such that

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\Phi \mathbf{x}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2 \quad (2.8)$$

The smallest constant δ_k (as a function of k) for which (2.8) holds is defined as the RIP constant.

For any two distinct k -space vectors \mathbf{x}_1 and \mathbf{x}_2 , denote the difference as $\mathbf{e} = \mathbf{x}_1 - \mathbf{x}_2$. Consequently, the support size of \mathbf{e} is at most $2k$ and $\|\mathbf{e}\|_2 > 0$. If a matrix Φ satisfies the RIP of order $2k$, we have

$$\|\Phi\mathbf{x}_1 - \Phi\mathbf{x}_2\|_2^2 = \|\Phi\mathbf{e}\|_2^2 \geq (1 - \delta_{2k})\|\mathbf{e}\|_2^2 > 0 \quad (2.9)$$

It implies that the CS observation $\mathbf{y}_1 = \Phi\mathbf{x}_1$ and $\mathbf{y}_2 = \Phi\mathbf{x}_2$ are also distinct. In other words, if the measurement matrix Φ satisfies the RIP of order $2k$, the distance between any pair of k -sparse signals in the high dimension can be approximately preserved when projected into a lower dimension. Thus RIP is a very useful condition to guarantee the existence of practical algorithms for reconstructing sparse and compressible signals from noisy measurements. See [15] for more details.

Although checking the validation of RIP for a given matrix is difficult, it has been proved that many random measurement matrices satisfy RIP with high probability, which includes the measurement matrices whose entries following the i.i.d. Gaussian distribution, Bernoulli distribution, and the partial Fourier matrix. For these matrices, the order $2k$ RIP condition is satisfied with overwhelming probability if the number of measurements satisfies the inequality

$$m \geq Ck \log(n/k) \quad (2.10)$$

where C is a constant depending on the specific measurement matrix instance [16, 17]. In other words, for the above mentioned matrices, there exists an algorithm with which the exact recovery of the sparse signal is achievable with overwhelmingly high probability.

Another more computable measurement matrix condition for analysing the CS recovery guarantee is the coherence [18].

Definition 2. The coherence of a matrix Φ , $\mu(\Phi)$, is given by the largest absolute inner product between any two of its columns Φ_i, Φ_j

$$\mu(\Phi) = \max_{1 \leq i < j \leq n} \frac{|\langle \Phi_i, \Phi_j \rangle|}{\|\Phi_i\|_2 \|\Phi_j\|_2} \quad (2.11)$$

It can be shown that the coherence of a matrix is bounded by $\mu(\Phi) \in [\sqrt{\frac{n-m}{m(n-1)}}, 1]$ [19]. The connection of coherence and RIP is explained in the following lemma.

lemma 1. ([20]) *If a matrix Φ has unit-norm columns and coherence $\mu(\Phi)$. Then Φ satisfies the RIP of order k with $\delta_k = (k - 1)\mu(\Phi)$ for all $k < \mu(\Phi)^{-1}$.*

For random measurement matrices with elements generated i.i.d. from the Bernoulli, Gaussian and sub-Gaussian distributions, their coherence is roughly

$$\mu \sim \sqrt{\frac{\ln n}{m}} \quad (2.12)$$

Then lemma 1 implies that for these matrices, the exact recovery happens with high probability when the number of measurements satisfies

$$m \geq ck^2 \ln n \quad (2.13)$$

where c is a constant depending on the matrix ensemble.

2.4 Compressed Sensing Reconstruction

There exists a wide variety of algorithmic approaches to the problem of recovering a sparse or compressible signal from an under-determined linear system. We now briefly review three typical types of methods in the literature.

2.4.1 ℓ_1 -Minimization

To retrieve the unknown signal \mathbf{x} as well as preserve its sparsity from the noiseless CS measurements, it is natural to consider the following optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{s.t.} \quad \mathbf{y} = \Phi \mathbf{x} \quad (2.14)$$

Unfortunately, the ℓ_0 minimization problem in (2.14) is not convex and lacks a practical procedure for even finding a solution that approximates the true minimum [21]. Alternatively, a more computationally tractable approach which relaxes the ℓ_0 -norm objective to the ℓ_1 -norm has been proposed and is called the *Basis Pursuit* (BP) [22]:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \mathbf{y} = \Phi \mathbf{x} \quad (2.15)$$

Equation 2.15 can then be solved by many off-the-shelf linear programming (LP) solvers. The use of ℓ_1 -minimization to promote sparsity has a long history, dating back to the work of Beurling on Fourier transform extrapolation from partial observations [20]. In [22], extensive empirical evidence suggests that the solution of (2.15) indeed recovers the sparsest solution in many cases. Donoho and co-authors further established a condition on the measurement matrix Φ , for which the BP solution is equivalent to solving the ℓ_0 -minimization problem in (2.14) [23]. Further work studying the relationship between (2.14) and (2.15) was conducted by several research groups, see [24–30].

In the presence of noise, another important problem broadly considered in the compressed sensing community is to solve

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \|\mathbf{y} - \Phi\mathbf{x}\|_2 \leq \varepsilon \quad (2.16)$$

Provided that the constraint in (2.16) is convex, the minimization is computationally feasible. See [31, 32] for some good solvers for (2.16). In the CS literature, more effort has actually been put into considering the unconstrained version of the optimization problem in (2.16):

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{y} - \Phi\mathbf{x}\|_2^2 + \kappa \|\mathbf{x}\|_1 \quad (2.17)$$

It is also known as the Basis pursuit denoising (BPDN) or LASSO (Least Absolution Shrinkage and Selection Operator). With appropriate choice of κ , the solution of LASSO coincides with that of the constraint minimization in (2.16). Several approaches for choosing κ are discussed in [33, 34].

The convex optimization technique is a powerful framework for recovering sparse signals since there exists accurate numerical solvers. The potential drawback of applying the ℓ_1 -minimization for the CS reconstruction though is that it may not be very efficient for large-scale problems.

2.4.2 Greedy methods

Another important class of CS reconstruction algorithms is the greedy method. They attempt to directly approximate the solution for (2.14) by iteratively identifying the support and value of the signal until a convergence criteria is met. Prominent examples of greedy methods include orthogonal matching pursuit (OMP) [35], stagewise OMP (StOMP) [36], compressive sam-

pling matching pursuit (CoSamp) [37] and iterative hard thresholding (IHT) [38]. Many greedy algorithms have been shown to have the similar performance guarantee as the ℓ_1 -minimization method [37, 38]. They also sometimes outperform the ℓ_1 -minimization based methods in terms of speed, storage and ease of implementation requirement for algorithms. Since the greedy methods are not the primary focus of this thesis, we refer the interested readers to [39–46] for a variety of existing algorithms.

2.4.3 Approximate Message Passing Based Methods

A very recent development for CS reconstruction is the approximate message passing algorithms, which is closely related to the approximate belief propagation for the CDMA multi-user detection problem [47–49]. The AMP was first introduced by Donoho and co-authors in [14]. It generally takes the iterative form

$$\begin{aligned}\hat{\mathbf{x}}^{t+1} &= \eta_t(\hat{\mathbf{x}}^t + \Phi^T \mathbf{z}^t) \\ \mathbf{z}^t &= \mathbf{y} - \Phi \hat{\mathbf{x}}^t + \frac{1}{\gamma} \mathbf{z}^{t-1} \langle \eta'_{t-1}(\hat{\mathbf{x}}^{t-1} + \Phi^T \mathbf{z}^{t-1}) \rangle\end{aligned}\tag{2.18}$$

where $\{\eta_t(\cdot)\}_{t \geq 0}$ is a sequence of scalar non-linear functions applied elementwise to the vector $\hat{\mathbf{x}}^t + \Phi^T \mathbf{z}^t$ with t indicating the iterations and $\eta'_t(\cdot)$ is the derivative of $\eta_t(\cdot)$ with respect to its first argument. With different selection of the non-linear function $\eta_t(\cdot)$ and possible extension of the algorithm, one would end up with different AMP variants. With a slight abuse of terminology, the term AMP is used to refer to both the class of AMP algorithms and the generic form in (2.18).

In the original AMP paper [14], the non-linear function takes the soft thresholding form

$$\eta_S(x; b) = \begin{cases} x - b & \text{if } x > b \\ 0 & \text{if } -b \leq x \leq b \\ x + b & \text{if } x < -b \end{cases}\tag{2.19}$$

The generic AMP algorithm has a very similar structure with the iterative soft thresholding (IST) algorithm. The only difference is that the IST does not have the additional term $\frac{1}{\gamma} \mathbf{z}^{t-1} \langle \eta'_{t-1}(\hat{\mathbf{x}}^{t-1} + \Phi^T \mathbf{z}^{t-1}) \rangle$. The IST approach is an iterative thresholding method for solving the LASSO problem in (2.17). The application of the soft thresholding function was first introduced for compressed sensing in [42]. It has also been proved in [42] that for $\|\Phi\|_2 < 1$,

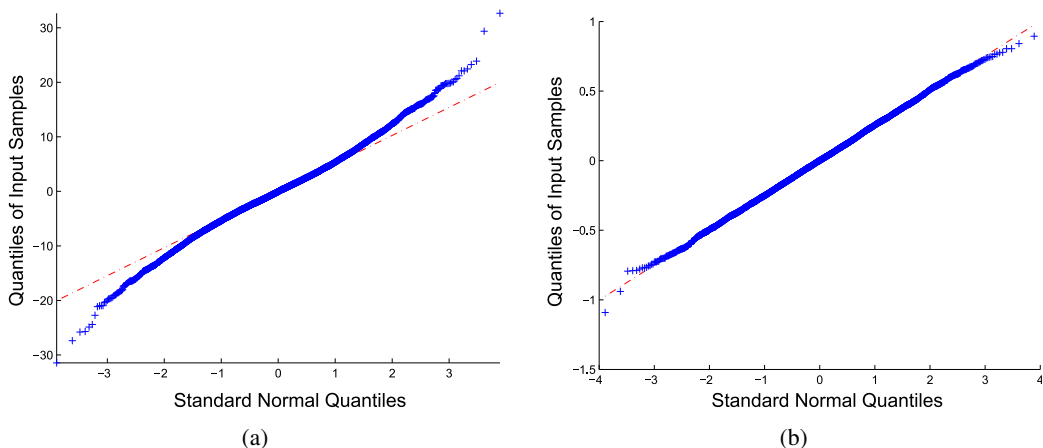


Figure 2.3: *The QQ plot for the IST (a) and AMP (b) residual at the 10th iteration. The linearity of the QQ plot indicates the Gaussian behaviour.*

the IST is guaranteed to converge towards to the solution of a LASSO minimizer.

With Φ being the measurement matrix with Gaussian distributed elements, the inclusion of the extra term $\frac{1}{\gamma} \mathbf{z}^{t-1} \langle \eta'_{t-1} (\hat{\mathbf{x}}^{t-1} + \Phi^T \mathbf{z}^{t-1}) \rangle$, which is designated as the ‘‘Onsager’’ reaction term, has fundamentally altered the IST behaviour. With the Onsager term, the AMP reconstruction can be interpreted as a recursive Gaussian denoising problem [14, 50, 51]. To be specific, the residual $\hat{\mathbf{x}}^t + \Phi^T \mathbf{z}^t - \mathbf{x}$ can be well modelled as an AWGN vector at each AMP iteration. The Gaussian behaviour of the AMP residual is demonstrated in the quantile-quantile (QQ) plot alongside the one for the IST residual in Fig. 2.3. In Fig. 2.3, the sample quantile is plotted against the theoretical quantile from a Gaussian distribution. The linear behaviour of the sample quantile in Fig. 2.3(b) implies the Gaussian nature of the AMP residual. In contrast, the IST residual does not demonstrate such behaviour. The difference is exactly introduced by the Onsager term.

Given such observation for AMP, the non-linearity $\eta_t(\cdot)$ essentially acts as the denoising function to remove the Gaussian noise to obtain a clearer data estimate at each AMP iteration. It also implies that better denoising function could be employed for this purpose. When the signal prior is available for reconstruction, the MMSE estimator is undoubtedly the best denoising choice. The corresponding BAMP algorithm is thus able to deliver better reconstruction than the generic AMP with the soft thresholding function [52]. When the signal prior is unknown, it has been proposed that a expectation-maximization (EM) learning procedure can be combined with the generic AMP algorithm as discussed in [53, 54]. In Chapter 5, we also propose an alternative which marries the Stein’s unbiased risk estimate with the AMP formula and delivers

a BAMP-like performance without the knowledge of the signal prior. A short review of AMP based algorithms will be presented later in section 2.7.

Another distinguishable feature that AMP possesses is the state evolution formalism, which analyses the asymptotic behaviour of AMP in the large system limit. Given the generic AMP algorithm defined in (2.18), the state evolution is a recursive function of the state variable τ_t^2 , which is the rescaled MSE for the signal estimate at iteration t

$$\tau_t^2 = \sigma_\xi^2 + \frac{1}{\gamma} \mathbb{E}\{[\eta_t(X_0 + \tau_{t-1}Z) - X_0]^2\} \quad (2.20)$$

The expectation is taken with respect to the independent random variable $Z \sim \mathcal{N}(0, 1)$ and X_0 , whose distribution coincides with the signal of interest \mathbf{x} . There will be more detailed discussion about SE later in Section 2.8. In the large system limit, the convergence point τ_*^2 of (2.20) accurately predicts the MSE of the AMP reconstructed signal $\hat{\mathbf{x}}$ [14, 55]

$$\lim_{n \rightarrow \infty} \frac{1}{n} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \mathbb{E}\{[\eta_t(X_0 + \tau_*^2 Z) - X]^2\} = \gamma(\tau_*^2 - \sigma_\xi^2) \quad (2.21)$$

Numerical evidence to support the SE dynamics can be found in [14] for i.i.d. Gaussian, Rademacher and partial Fourier matrices. The agreement between the SE prediction and the Monte Carlo simulation is remarkably good for signal dimension of the order of a few hundreds. In [55], the authors proved that in the large system limit, SE holds for random measurement matrices with i.i.d. Gaussian entries. They also commented that although the proof technique heavily relies on the Gaussian assumption, the SE is expected to hold for a broader range of random matrices. The benefits of having the SE dynamics is twofold. First, the expected MSE of an AMP-based algorithm can be obtained without running Monte Carlo simulations. We will see later in Chapter 3 how it can be used to quantify the reconstruction performance for a certain pair of measurement matrix and recovery scheme. Second, the SE dynamics provides a systematic way for optimizing the non-linear function $\eta_t(\cdot)$ in the AMP iteration. In Chapter 5, the parametric SURE-AMP is proposed by choosing $\eta_t(\cdot)$ that minimizes the right hand side of (2.20). A detailed summary of the SE dynamics for different AMP variants and the intuitive explanation for SE will be presented in section 2.8.

Theoretically and empirically speaking, AMP is a class of computationally efficient algorithms with the state-of-the-art performance [14, 50, 56]. It will be the main reconstruction tool that we resort to throughout the thesis.

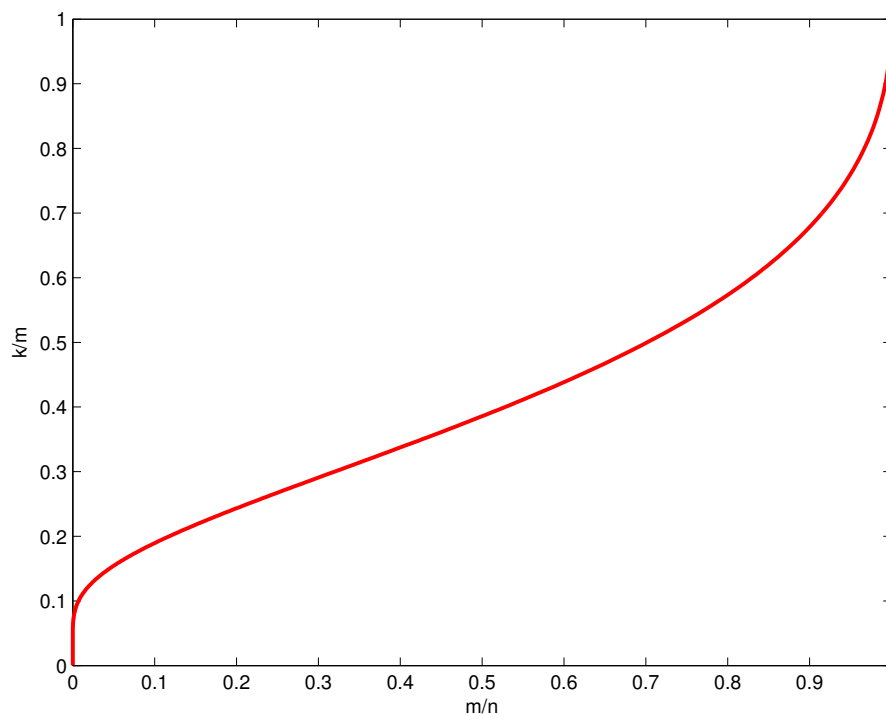


Figure 2.4: *Theoretical phase transition for ℓ_1 -minimization [1, 2].*

2.5 Phase Transitions

In section 2.3, the number of measurements required for exact recovery of sparse signals is analysed by considering the RIP and coherence condition of the measurement matrix. Although the order of bound derived from RIP for random measurement matrices in (2.10) is optimal, the unknown constant C makes the bound of little practical use for real engineering applications. The phase transition is then proposed to provide a more specific guidance for sampling as well as a fair scheme to compare the undersampling-sparsity tradeoff for various CS algorithms.

The phase transition phenomenon was first empirically observed and rigorously characterized in [1, 2] for the ℓ_1 -minimization method. We assume γ and ρ are fixed as $m, n \rightarrow \infty$. For a fixed signal support size, there is a well-defined 'break-down' sampling ratio above which ℓ_1 -minimization can successfully recover the original signal with overwhelmingly high probability. The phase diagram indicating the probability of success and failure of the ℓ_1 -minimization as a function of ρ and γ is shown in Fig. 2.4. The red line is the theoretical phase transition curve for ℓ_1 -minimization method assuming large system, whose derivation can be found in [1, 2]. In the 'upper' region of the plot, the probability of exact recovery tends to zero exponentially fast. While in the 'lower' region, the probability of successful recovery tends to one exponentially fast [14]. In practice with finite problem size, the transition zone between failure

and success becomes narrower as n increases and corresponds to the theoretical curve in the large n limit.

The existence of the phase transition has also been observed for other algorithms. For example, the phase transition of StOMP is shown to be comparable to the ℓ_1 -minimization methods in [36]. The lower bound of the phase transition for IHT and CoSamp with Gaussian measurement matrices and a certain distribution for non-zero coefficients are considered in [57]. In [14, 58], the phase transition for AMP is rigorously derived with explicit expressions from the state evolution perspective. Remarkably, the result coincides with the phase transition derived earlier for the ℓ_1 -minimization in [1]. For more information on phase transition, please refer to [2, 13, 18, 44, 59].

2.6 Graphical Model for CS and AMP Derivation

After the brief summary of the CS related knowledge, the rest of this chapter is devoted to a more detailed introduction for AMP. The AMP algorithm considers the CS reconstruction problem from a probabilistic perspective using the graphical model approach. Essentially it postulates a joint probability distribution $p(\mathbf{x}, \mathbf{y})$ on (\mathbf{x}, \mathbf{y}) and infers \mathbf{x} from \mathbf{y} by approximating the posterior distribution $p(\mathbf{x}|\mathbf{y})$. The graphical model approach manipulates the distributions involved and factorizes them into a specific graph model to aid the inference procedure.

We take the stochastic CS prior in section 2.2.2 and model \mathbf{x} as a vector with i.i.d. entries that does not depend on Φ

$$p(\mathbf{x}|\Phi) = \prod_{i=1}^n [(1 - \rho)\delta(x_i) + \rho\Gamma(x_i)] \quad (2.22)$$

The compressible signal case is included by setting $\rho = 1$. From here on, $\Gamma(\cdot)$ refers to the distribution of the non-zero coefficients. In practice, the estimation of the prior can be obtained with moment matching method as illustrated later in Chapter 3. We assume the noise vector has i.i.d. Gaussian random entries

$$p(\boldsymbol{\xi}) = \prod_{i=1}^m \mathcal{N}(\xi_i; 0, \sigma_\xi^2) \quad (2.23)$$

with $\sigma_\xi^2 = 0$ corresponding to the noiseless scenario. Then the conditional distribution of \mathbf{y}

given \mathbf{x} and Φ is calculated as

$$p(\mathbf{y}|\mathbf{x}, \Phi) = \prod_{j=1}^m \mathcal{N}(y_j; \sum_{i=1}^n \Phi_{ji}x_i, \sigma_\xi^2) \quad (2.24)$$

From Bayes' theorem, the posterior is calculated as

$$p(\mathbf{x}|\mathbf{y}, \Phi) = \frac{p(\mathbf{y}|\mathbf{x}, \Phi)p(\mathbf{x}|\Phi)}{p(\mathbf{y}|\Phi)} = \frac{p(\mathbf{y}|\mathbf{x}, \Phi)p(\mathbf{x}|\Phi)}{\int p(\mathbf{y}|\mathbf{x}, \Phi)p(\mathbf{x}|\Phi)d\mathbf{x}} \quad (2.25)$$

Applying (2.22) and (2.24) to (2.25), we have

$$p(\mathbf{x}|\mathbf{y}, \Phi) = \frac{1}{Z(\mathbf{y}, \Phi)} \prod_{i=1}^n [(1-\rho)\delta(x_i) + \rho\Gamma(x_i)] \prod_{j=1}^m \frac{1}{\sqrt{2\pi\sigma_\xi^2}} e^{-\frac{1}{2\sigma_\xi^2}(y_j - \sum_{k=1}^n \Phi_{jk}x_k)^2} \quad (2.26)$$

where $Z(\mathbf{y}, \Phi) = p(\mathbf{y}|\Phi)$ is the normalization constant. The MMSE estimate of \mathbf{x} can be then extracted using the conditional expectation

$$\hat{\mathbf{x}} = \int \mathbf{x}p(\mathbf{x}|\mathbf{y}, \Phi) d\mathbf{x} \quad (2.27)$$

An important problem with the estimator in (2.27) is that its exact computation in general is very hard. In this section, we will re-derive the BAMP algorithm which approximately infers this estimate.

2.6.1 Factor Graph and Sum-Product Algorithm Review

Before embarking on the actual derivation, it is convenient to go through the basic knowledge of the factor graph and the sum-product algorithm. A close observation of the complicated global function in (2.26) reveals that it is a product of several simpler ‘‘local’’ functions, each of which depends only on a subset of the variables. This factorized structure can be conveniently described by its factor graph, a bipartite graph that connects local functions with their related argument variables. In Fig. 2.5, the factor graph for the posterior in (2.26) is illustrated: there is a ‘‘variable node’’ $x_i, i \in [n]$ for each signal entry and a ‘‘factor node’’ $a_j, j \in [m]$ for each term $a_j(\mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma_\xi^2}} e^{-\frac{1}{2\sigma_\xi^2}(y_j - \sum_{k=1}^n \Phi_{jk}x_k)^2}$. A variable node x_i and a factor node a_j are connected by an edge if and only if x_i is the argument of a_j .

In [60], a generic message passing algorithm designated as the sum-product algorithm is pro-

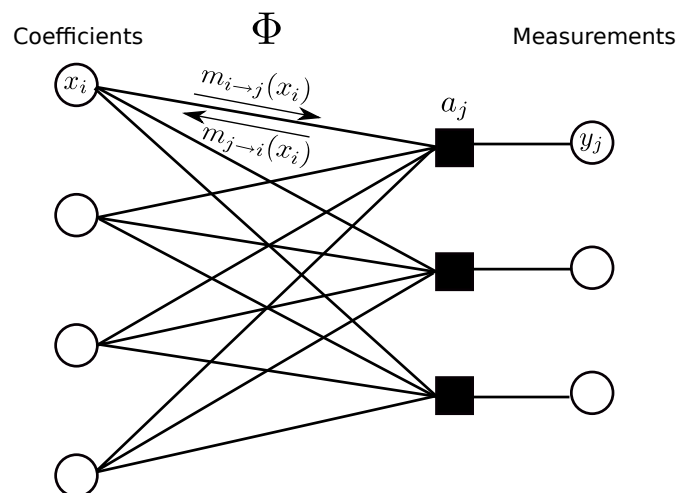


Figure 2.5: The factor graph associated with the probability distribution in (2.26). Empty circles represent variables x_i , $i \in [n]$ and y_j , $j \in [m]$. Squares correspond to measurement function a_j . $m_{i \rightarrow j}(x_i)$ and $m_{j \rightarrow i}(x_i)$ are messages representing interaction among nodes.

posed, which operates on the factor graph and calculates the “message” associated to each directed edge in the factor graph. In Fig. 2.5, we denote the message sending from a variable node to a factor node as $m_{i \rightarrow j}(x_i)$ and $m_{j \rightarrow i}(x_i)$ vice versa.

Algorithm 1 : Sum-Product Update Rule [60]

- 1: The message sent from a node ν on an edge e is the product of the local function at ν (or the unit function if ν is a variable node) with all messages received at ν on edges other than e , summarized for the variable associated with e .
-

According to [60], the update rule for the sum-product algorithm is summarized in Algorithm 1. To better explain it, we present a portion of the factor graph featuring one variable node x , one factor node f and the related messages in Fig. 2.6. Here $n(x) \setminus f$ are all neighbour nodes of x except the node f . Similarly, $n(f) \setminus x$ is the set of neighbours of f except x . Then the messages exchanging between x and f can be computed as following

variable to local function:

$$m_{x \rightarrow f}(x) = \prod_{h_i \subset n(x) \setminus f} m_{h_i \rightarrow x}(x) \quad (2.28)$$

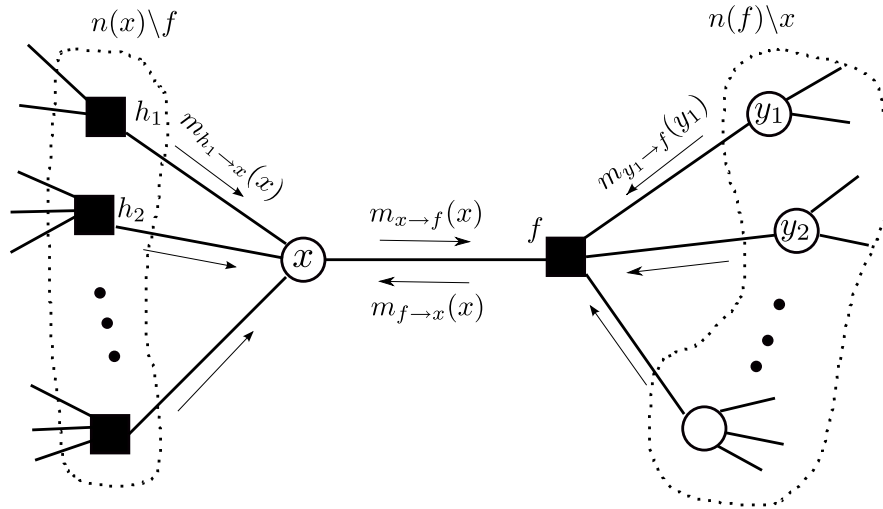


Figure 2.6: A factor graph featuring one variable node and one factor node, which is used to explain the updating rule for the sum-product algorithm.

local function to variable:

$$m_{f \rightarrow x}(x) = \sum_{\sim x} \left(f(\kappa) \prod_{y_i \in n(f) \setminus x} m_{y_i \rightarrow f}(y_i) \right) \quad (2.29)$$

where $\kappa = n(f)$ is the set of all arguments of the function f . As defined in [60], the operator $\sum_{\sim(x)}$ is the “not-sum” operation. Instead of specifying the variables being summed over, the “not-sum” operation indicates the variables that are not summed over. For example, for a function f with three variables x_1, x_2 and x_3 , the “not-sum for x_2 ” is computed as

$$\sum_{\sim x_2} f(x_1, x_2, x_3) \equiv \sum_{x_1 \in A_1} \sum_{x_3 \in A_3} f(x_1, x_2, x_3) \quad (2.30)$$

marginal function for variable node

$$\begin{aligned} m(x) &= m_{f \rightarrow x}(x) \prod_{h_i \in n(x) \setminus f} m_{h_i \rightarrow x}(x) \\ &= m_{f \rightarrow x}(x) m_{x \rightarrow f}(x) \end{aligned} \quad (2.31)$$

Together with (2.28), (2.29) and (2.31), we are well-equipped to derive a wide variety of algorithms developed in signal processing, digital communication and artificial intelligence over the appropriate graphical models. In fact, the sum-product algorithm can be seen as the generalization of the forward/backward algorithm, Kalman filter, Turbo decoding algorithm and certain FFT etc [60]. Pearl’s belief propagation algorithm for Bayesian networks can also be

derived as a special instance [61–67]. One thing worth noting is that although the sum-product algorithm is not generally guaranteed to converge for graphs with closed loops, for example Fig. 2.5, sometimes favourable results are obtained by performing the message passing in a recursive manner until certain convergence criteria is satisfied [68].

2.6.2 Relaxed Message Passing for CS

There are several research groups that provided independent derivations of AMP from the standard message passing equations [53, 58, 69, 70]. In this section, we stick with the notation in [53] and reproduce their derivation of the BAMP algorithm, which is widely used in this thesis.

Given the factor graph representation for $p(\mathbf{x}|\mathbf{y}, \Phi)$ and the sum-product algorithm review in section 2.6.1, we can explicitly write the canonical message passing equations, which consists of $2mn$ probability functions or messages, namely $m_{i \rightarrow j}(x_i)$ and $m_{j \rightarrow i}(x_i)$ with $i \in [n]$, $j \in [m]$.

$$m_{i \rightarrow j}(x_i) = \frac{1}{Z_{i \rightarrow j}} [(1 - \rho)\delta(x_i) + \rho\Gamma(x_i)] \prod_{p \neq j} m_{p \rightarrow i}(x_i) \quad (2.32)$$

$$m_{j \rightarrow i}(x_i) = \frac{1}{Z_{j \rightarrow i}} \int \prod_{q \neq i} dx_q e^{-\frac{1}{2\sigma^2\xi}(y_j - \Phi_{ji}x_i - \sum_{q \neq i} \Phi_{jq}x_q)^2} \prod_{q \neq i} m_{q \rightarrow j}(x_q) \quad (2.33)$$

where $Z^{i \rightarrow j}$, $Z^{j \rightarrow i}$ are normalization factors so that $\int m_{i \rightarrow j}(x_i) dx_i = \int m_{j \rightarrow i}(x_i) dx_i = 1$. The CS reconstruction of x_i is obtained through the expectation of the local belief on x_i , which is the product of all messages directed towards x_i as in (2.31).

$$m_i(x_i) = \frac{1}{Z_i} [(1 - \rho)\delta(x_i) + \rho\Gamma(x_i)] \prod_j m_{j \rightarrow i}(x_i) \quad (2.34)$$

Unfortunately the exact implementation of the message passing algorithm to propagate the pdfs, i.e. from (2.32) to (2.34), is intractable. Hence a relaxed message passing system where the messages are real numbers instead of the pdfs is derived to approximate the dynamics. For the CS reconstruction problem, the messages that we are interested in are the mean and variance of the marginal distributions for the desirable variables. The relaxation is valid in the large system limit and by assuming all measurement matrix elements are Gaussian distributed and scale as $\mathcal{O}(1/\sqrt{m})$. First let us define the mean and variance of the variable-to-factor message

$m_{i \rightarrow j}(x_i), i \in [n]$ as the following

$$\alpha_{i \rightarrow j} \equiv \int dx_i \quad x_i m_{i \rightarrow j}(x_i) \quad (2.35)$$

$$\nu_{i \rightarrow j} \equiv \int dx_i \quad x_i^2 m_{i \rightarrow j}(x_i) - \alpha_{i \rightarrow j}^2 \quad (2.36)$$

With the definition (2.35) and (2.36), the factor-to-variable message $m_{j \rightarrow i}(x_i)$ can be approximated with the following Gaussian format in the large system limit

$$m_{j \rightarrow i}(x_i) = \frac{1}{\tilde{Z}_{j \rightarrow i}} e^{-\frac{A_{j \rightarrow i}}{2} x_i^2 + B_{j \rightarrow i} x_i}, \quad \tilde{Z}_{j \rightarrow i} = \sqrt{\frac{2\pi}{A_{j \rightarrow i}}} e^{\frac{B_{j \rightarrow i}^2}{2A_{j \rightarrow i}}} \quad (2.37)$$

where $\tilde{Z}_{j \rightarrow i}$ is the normalization factor containing all x_i -independent terms and $A_{j \rightarrow i}, B_{j \rightarrow i}$ are defined as:

$$A_{j \rightarrow i} = \frac{\Phi_{ji}^2}{\sigma_\xi^2 + \sum_{q \neq i} \Phi_{jq}^2 \nu_{q \rightarrow j}} \quad (2.38)$$

$$B_{j \rightarrow i} = \frac{\Phi_{ji}(y_j - \sum_{q \neq i} \Phi_{jq} \alpha_{q \rightarrow j})}{\sigma_\xi^2 + \sum_{q \neq i} \Phi_{jq}^2 \nu_{q \rightarrow j}} \quad (2.39)$$

The detailed derivation from (2.37) to (2.39) is given in the Appendix A, in which the second order Taylor expansions of some exponential terms in $m_{j \rightarrow i}(x_i)$ are used and components that are above the order of $\mathcal{O}(1/m)$ are assumed vanishing as $m \rightarrow \infty$. This simplified form (2.37) basically shows that a pair of real numbers, namely $(A_{j \rightarrow i}, B_{j \rightarrow i})$, is enough to characterize the factor-to-variable message $m_{j \rightarrow i}(x_i)$.

With (2.37), the variable-to-factor message $m_{i \rightarrow j}(x_i)$ in (2.32) becomes

$$m_{i \rightarrow j}(x_i) = \frac{1}{\tilde{Z}_{i \rightarrow j}} [(1 - \rho)\delta(x_i) + \rho\Gamma(x_i)] e^{-\frac{x_i^2}{2} \sum_{p \neq j} A_{p \rightarrow i} + x_i \sum_{p \neq j} B_{p \rightarrow i}} \quad (2.40)$$

To obtain the closed message passing iteration, we would like to obtain the mean and variance of (2.40) to characterize the message as well. To this end, let us first define a general probability distribution

$$p(x, R, \Sigma^2) = \frac{1}{\hat{Z}(R, \Sigma^2)} [(1 - \rho)\delta(x) + \rho\Gamma(x)] \frac{1}{\sqrt{2\pi\Sigma}} e^{-\frac{(x-R)^2}{2\Sigma^2}} \quad (2.41)$$

where $\hat{Z}(R, \Sigma^2)$ is the normalization factor for the distribution. Then the mean and variance

for $p(x, R, \Sigma^2)$ are calculated as

$$f_a(R, \Sigma^2) = \int dx \quad xp(x, R, \Sigma^2) \quad (2.42)$$

$$f_c(R, \Sigma^2) = \int dx \quad x^2p(x, R, \Sigma^2) - f_a^2(R, \Sigma^2) \quad (2.43)$$

Given the above definition, the mean and variance of the message $m_{i \rightarrow j}(x_i)$ can thus be expressed as

$$\alpha_{i \rightarrow j}(x_i) = f_a\left(\frac{\sum_{p \neq j} B_{p \rightarrow i}}{\sum_{p \neq j} A_{p \rightarrow i}}, \frac{1}{\sum_{p \neq j} A_{p \rightarrow i}}\right) \quad (2.44)$$

$$\nu_{i \rightarrow j}(x_i) = f_c\left(\frac{\sum_{p \neq j} B_{p \rightarrow i}}{\sum_{p \neq j} A_{p \rightarrow i}}, \frac{1}{\sum_{p \neq j} A_{p \rightarrow i}}\right) \quad (2.45)$$

Finally the local belief on x_i is approximated as

$$\begin{aligned} m_i(x_i) &= \frac{1}{\tilde{Z}_i} [(1 - \rho)\delta(x_i) + \rho\Gamma(x_i)] \prod_j m_{j \rightarrow i}(x_i) \\ &= \frac{1}{\tilde{Z}_i} [(1 - \rho)\delta(x_i) + \rho\Gamma(x_i)] e^{-\frac{x_i^2}{2} \sum_j A_{j \rightarrow i} + x_i \sum_j B_{j \rightarrow i}} \end{aligned} \quad (2.46)$$

Noticing the similarity between (2.46) and (2.40), we can easily obtain the corresponding mean and variance of the local belief for x_i as

$$\alpha_i(x_i) = f_a\left(\frac{\sum_j B_{j \rightarrow i}}{\sum_j A_{j \rightarrow i}}, \frac{1}{\sum_j A_{j \rightarrow i}}\right) \quad (2.47)$$

$$\nu_i(x_i) = f_c\left(\frac{\sum_j B_{j \rightarrow i}}{\sum_j A_{j \rightarrow i}}, \frac{1}{\sum_j A_{j \rightarrow i}}\right) \quad (2.48)$$

Equations (2.38), (2.39), (2.44), (2.45), (2.47) and (2.48) all together form a closed message passing algorithm, which is considerably simpler than the original sum-product iteration to propagate pdfs. Yet fundamentally it is still a $2mn$ -message system. In the next section, we will further reduce the complexity to obtain a system with only $m + n$ messages.

2.6.3 From Relaxed Message Passing to AMP

Close observation of the relaxed message passing equations in (2.44) and (2.45) reveals that the messages $\alpha_{i \rightarrow j}$ and $\nu_{i \rightarrow j}$ are almost independent of j . It is reasonable to think that in the large

system limit, the dependencies on the instance are so weak that can be neglected. Although it is tempting to directly replace $\alpha_{i \rightarrow j}$ and $\nu_{i \rightarrow j}$ with α_i and ν_i respectively to further simplify the algorithm, special care must be taken when discarding the negligible terms. The success of the further relaxation of the message passing system depends on whether we keep the correct ‘‘Onsager’’ term as we mentioned in Section 2.4.3. Again we follow the procedure in [53] and complete the final step to obtain the $m + n$ approximated message passing algorithm. Remember that the following approximation is always within the assumption that the measurement matrix is dense and all its element scale as $\mathcal{O}(1/\sqrt{m})$.

We begin by defining some scalars

$$c_i = \frac{1}{\sum_j A_{j \rightarrow i}}, \quad \varepsilon_i = \frac{\sum_j B_{j \rightarrow i}}{\sum_j A_{j \rightarrow i}} \quad (2.49)$$

$$\omega_j = \sum_i \Phi_{ji} \alpha_{i \rightarrow j}, \quad V_j = \sum_i \Phi_{ji}^2 \nu_{i \rightarrow j} \quad (2.50)$$

In the large system limit, c_i and ε_i can be approximately expressed as

$$\begin{aligned} c_i &= \left[\sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + \sum_{q \neq i} \Phi_{jq}^2 \nu_{q \rightarrow j}} \right]^{-1} \\ &= \left[\sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + V_j - \Phi_{ji}^2 \nu_{i \rightarrow j}} \right]^{-1} \approx \left[\sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + V_j} \right]^{-1} \end{aligned} \quad (2.51)$$

and

$$\begin{aligned} \varepsilon_i &= \left[\sum_j \frac{\Phi_{ji}(y_j - \sum_{q \neq i} \Phi_{jq} \alpha_{q \rightarrow j})}{\sigma_\xi^2 + \sum_{q \neq i} \Phi_{jq}^2 \nu_{q \rightarrow j}} \right] \left[\sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + \sum_{q \neq i} \Phi_{jq}^2 \nu_{q \rightarrow j}} \right]^{-1} \\ &= \left[\sum_j \frac{\Phi_{ji}(y_j - \sum_{q \neq i} \Phi_{jq} \alpha_{q \rightarrow j})}{\sigma_\xi^2 + V_j - \Phi_{ji}^2 \nu_{i \rightarrow j}} \right] \left[\sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + V_j - \Phi_{ji}^2 \nu_{i \rightarrow j}} \right]^{-1} \\ &\approx \left[\sum_j \Phi_{ji} \frac{y_j - \omega_j}{\sigma_\xi^2 + V_j} + \alpha_i \sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + V_j} \right] \left[\sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + V_j} \right]^{-1} \\ &= \alpha_i + \frac{\sum_j \Phi_{ji} \frac{y_j - \omega_j}{\sigma_\xi^2 + V_j}}{\sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + V_j}} \end{aligned} \quad (2.52)$$

Next we are going to approximate $\alpha_{i \rightarrow j}$ in terms of α_i . By doing so, we omit all terms that are not linear with Φ_{ji}

$$\alpha_{i \rightarrow j} \approx \alpha_i - B_{j \rightarrow i} \nu_i \quad (2.53)$$

The detailed derivation is given in Appendix B.

With (2.53), the scalar ω_j can be further approximated as

$$\begin{aligned}
 \omega_j &= \sum_i \Phi_{ji} [\alpha_i - B_{j \rightarrow i} \nu_i] \\
 &= \sum_i \Phi_{ji} \alpha_i - \sum_i \frac{\Phi_{ji}^2 (y_j - \sum_{q \neq i} \Phi_{jq} \alpha_{q \rightarrow j})}{\sigma_\xi^2 + \sum_{q \neq i} \Phi_{jq}^2 \nu_{q \rightarrow j}} \nu_i \\
 &\approx \sum_i \Phi_{ji} \alpha_i - \frac{y_j - \omega_j}{\sigma_\xi^2 + V_j} \sum_i \Phi_{ji}^2 \nu_i
 \end{aligned} \tag{2.54}$$

The approximation of V_j is similar to ω_j . This time all the correction terms are assumed negligible in the large system limit. Therefore, we have

$$V_j = \sum_i \Phi_{ji}^2 \nu_i \tag{2.55}$$

Equations (2.47), (2.48) together with the approximation terms in (2.51), (2.52), (2.54) and (2.55) form a complete AMP algorithm. As mentioned before, messages are exchanged iteratively until convergence for loopy factor graphs. In the context of statistical physics, the resulting iterative system corresponds to the Thouless-Anderson-Palmer (TAP) equations used in the study of spin glasses [71]. It is thus designated as the TAP-AMP algorithm in [53] and summarized here in Algorithm 2.

Algorithm 2 : TAP-AMP [53]

- 1: **initialization:** $\hat{\mathbf{x}}^0 = \mathbf{0}$, $\mathbf{z}^0 = \mathbf{y}$, $c^0 = \sigma_x^2$
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: $V_j^{t+1} = \sum_i \Phi_{ji}^2 \nu_i^t$
 - 4: $\omega_j^{t+1} = \sum_i \Phi_{ji} \hat{x}_i^t - \frac{y_j - \omega_j^t}{\sigma_\xi^2 + V_j^t} \sum_i \Phi_{ji}^2 \nu_i^t$
 - 5: $c_i^{t+1} = \left[\sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + V_j^{t+1}} \right]^{-1}$
 - 6: $\varepsilon_i^{t+1} = \hat{x}_i^t + \frac{\sum_j \Phi_{ji} \frac{y_j - \omega_j^{t+1}}{\sigma_\xi^2 + V_j^{t+1}}}{\sum_j \frac{\Phi_{ji}^2}{\sigma_\xi^2 + V_j^{t+1}}}$
 - 7: $\hat{x}_i^{t+1} = f_a(\varepsilon_i^{t+1}, c_i^{t+1})$
 - 8: $\nu_i^{t+1} = f_c(\varepsilon_i^{t+1}, c_i^{t+1})$
 - 9: **end for**
-

So far we have completed the major approximation steps to go from the standard message passing equations to the TAP-AMP iteration which involves only matrix multiplication. Throughout

the derivation for the TAP-AMP, we leveraged on the assumption that the measurement matrix is a dense matrix. Indeed, the simplification of the message passing for the system emerges in the large system limit. As commented in [50], “this is one instance of the blessings of dimensionality”. One thing worth noting is that to perform TAP-AMP, the signal prior $p(\mathbf{x})$ is assumed known. Moreover, TAP-AMP is applicable for a general form of Gaussian measurement matrix whose entries does not necessarily come from the same distribution. This special feature will come into use later in Chapter 4 as the reconstruction algorithm for the modulated measurement matrix.

2.7 AMP Based Algorithm Summary

In this section, we will summarize three different types of AMP variants. Previously we reviewed the approximation steps to obtain the TAP-AMP algorithm. We start with some further simplification of TAP-AMP with the homogeneous assumption for the Gaussian measurement matrix and present the BAMP algorithm, which utilizes the signal pdf for reconstruction. When the signal prior is unknown, the ℓ_1 -AMP algorithm approximately optimizes the solution in the maximin framework and obtains the same phase transition performance as the ℓ_1 -minimization method for sparse signal reconstruction [72]. Finally, we will give a brief review of the GAMP algorithm, which is capable of dealing with a wide range of noise models.

2.7.1 Bayesian optimal AMP

Evolving from the TAP-AMP to the BAMP we mainly leverage on the assumption that Φ is a homogeneous matrix with i.i.d. Gaussian random entries $\Phi_{ij} \sim \mathcal{N}(0, 1/m)$. In the large system limit, the Φ_{ij}^2 terms in Algorithm 2 can be effectively replaced by $1/m$. We can hence neglect the dependence on the index j and consider all V_j to be the same. Consequently, c_i^t in line 5 of TAP-AMP can also be replaced with a scalar c^t independent of i . With a change of variable, TAP-AMP then transforms to BAMP as summarized in Algorithm 3.

Recall that $f_a(\cdot)$ and $f_c(\cdot)$ are defined as the mean and variance for the general probability function $p(x, R, \Sigma^2)$ in (2.41). Actually this general form has a posterior interpretation: let x be a random variable with $p(x) = (1 - \rho)\delta(x) + \rho\Gamma(x)$ and $r = x + \omega$ be the noisy data corrupted by the Gaussian noise $\omega \sim \mathcal{N}(\omega; 0, \Sigma^2)$. The likelihood then takes the form of $p(r|x) \sim \mathcal{N}(r; x, \Sigma^2)$. From the Bayes' rule, it is straightforward to show (2.41) is the posterior $p(x|r)$. As a consequence, the function $f_a(\cdot)$ and $f_c(\cdot)$ are the MMSE estimator and

Algorithm 3 : BAMP [52, 69]

```

1: initialization:  $\hat{\mathbf{x}}^0 = \mathbf{0}, \mathbf{z}^0 = \mathbf{y}, c^0 = \sigma_x^2$ 
2: for  $t = 1, 2, \dots$  do
3:    $\boldsymbol{\varepsilon}^t = \boldsymbol{\Phi}^T \mathbf{z}^t + \mathbf{x}^t$ 
4:    $\hat{\mathbf{x}}^{t+1} = f_a(\boldsymbol{\varepsilon}^t, c^t)$ 
5:    $v^{t+1} = f_c(\boldsymbol{\varepsilon}^t, c^t)$ 
6:    $\mathbf{z}^{t+1} = \mathbf{y} - \boldsymbol{\Phi} \hat{\mathbf{x}}^{t+1} + \frac{1}{\gamma} \mathbf{z}^t \langle f'_a(\boldsymbol{\varepsilon}^t, c^t) \rangle$ 
7:    $c^{t+1} = \sigma_\xi^2 + \frac{1}{\gamma} \langle v^{t+1} \rangle$ 
8: end for

```

the conditional variance of x given r respectively.

$$f_a(R, \Sigma^2) = \mathbb{E}_{x|r}(x|r = R) \quad (2.56)$$

$$f_c(R, \Sigma^2) = \text{Var}_{x|r}(x|r = R) \quad (2.57)$$

With the correct signal prior, we achieve the maximum denoising amount at each step using the MMSE estimator and eventually obtain the optimal reconstruction in the least squared sense. For some special cases, i.e. BG or GMD, a closed form expression can be obtained for $f_a(\cdot)$ and $f_c(\cdot)$. For more general signal distributions, the calculation of the conditional mean and variance can be conducted through numerical integration over x .

When the explicit expression for $f_a(\cdot)$ is not available, numerical calculation of its derivative $f'_a(\cdot)$ can sometimes be non-trivial and introduce unnecessary error if not treated properly. With the relationship proved in Appendix C, we can express $f'_a(\cdot)$ as a function of $f_c(\cdot)$ and denote the corresponding algorithm as BAMP-V2. The difference between BAMP and BAMP-V2 lies

Algorithm 4 : BAMP-V2

```

1: initialization:  $\hat{\mathbf{x}}^0 = \mathbf{0}, \mathbf{z}^0 = \mathbf{y}, c^0 = \sigma_x^2$ 
2: for  $t = 1, 2, \dots$  do
3:    $\boldsymbol{\varepsilon}^t = \boldsymbol{\Phi}^T \mathbf{z}^t + \mathbf{x}^t$ 
4:    $\hat{\mathbf{x}}^{t+1} = f_a(\boldsymbol{\varepsilon}^t, c^t)$ 
5:    $v^{t+1} = f_c(\boldsymbol{\varepsilon}^t, c^t)$ 
6:    $\mathbf{z}^{t+1} = \mathbf{y} - \boldsymbol{\Phi} \hat{\mathbf{x}}^{t+1} + \frac{v^{t+1}}{\gamma c^2} \mathbf{z}^t$ 
7:    $c^{t+1} = \sigma_\xi^2 + \frac{1}{\gamma} \langle v^{t+1} \rangle$ 
8: end for

```

in line 6 for updating \mathbf{z}^{t+1} . In general, the BAMP algorithm benefits from its simple iterative form and the ability of making use of the signal prior for reconstruction. It is the most efficient CS reconstruction algorithm with the state-of-art recovery quality as far as the CS literature is concerned.

2.7.2 ℓ_1 -AMP

The counterpart of the ℓ_1 -minimization approach in the AMP family is the ℓ_1 -AMP algorithm, which deals with the CS problem when sparsity/compressibility is the only prior information we know for the signal. At each ℓ_1 -AMP iteration, the signal denoising is performed by using the soft shrinkage function in (2.19). Theoretical analysis and extensive simulations confirm that the ℓ_1 -AMP has the identical phase transition curve as the ℓ_1 -minimization based algorithms for sparse signals reconstruction, but runs faster than conventional ℓ_1 -solvers [14].

We first present the BP-AMP algorithm, which amounts to solving the basis pursuit problem in (2.15). The BP-AMP possesses both the low complexity feature of the iterative thresholding

Algorithm 5 : BP-AMP [52]

- 1: **initialization:** $\hat{\mathbf{x}}^0 = \mathbf{0}$, $\mathbf{z}^0 = \mathbf{y}$, $c^0 \leftarrow \sigma_x^2$
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: $\boldsymbol{\varepsilon}^t = \mathbf{\Phi}^T \mathbf{z}^t + \mathbf{x}^t$
 - 4: $\hat{\mathbf{x}}^{t+1} = \eta_S(\boldsymbol{\varepsilon}^t; \zeta c^t)$
 - 5: $\mathbf{z}^{t+1} = \mathbf{y} - \mathbf{\Phi} \hat{\mathbf{x}}^{t+1} + \frac{1}{\gamma} z^t \langle \eta'_S(\boldsymbol{\varepsilon}^t; \zeta c^t) \rangle$
 - 6: $c^{t+1} = \frac{\|\mathbf{z}\|_2^2}{m}$
 - 7: **end for**
-

algorithms and the computation power of the ℓ_1 -minimization, thus conquering the large system size obstacle that often occurs in applications.

This is the first AMP algorithm that Donoho et al. proposed in their original paper [14]. In the BP-AMP iteration, ζ is the constant which can be tuned optimally before applying the algorithm based on the sampling ratio γ . In [14] the optimal ζ that achieves the maximum phase transition for sparse signal reconstruction is proved to be

$$\zeta(\gamma) = \frac{1}{\sqrt{\gamma}} \arg \max_{z \geq 0} \left\{ \frac{1 - 2/\delta[(1 + z^2)\Omega(-z) - z\phi(z)]}{1 + z^2 - 2[(1 + z^2)\Omega(-z) - z\phi(z)]} \right\} \quad (2.58)$$

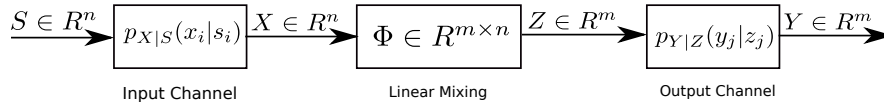
where $\phi(z) = e^{-z^2/2}/\sqrt{2\pi}$ and $\Omega(z) = \int_{-\infty}^z \phi(x) dx$.

For the LASSO problem in (2.17), its AMP counterpart is summarized as BPDN-AMP.

The construction of the graphical model for LASSO is detailed in [52, 58].

Algorithm 6 :BPDN-AMP [58]

-
- 1: **initialization:** $\hat{\mathbf{x}}^0 = \mathbf{0}, \mathbf{z}^0 = \mathbf{y}, c^0 = \sigma_x^2$
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: $\boldsymbol{\varepsilon}^t = \boldsymbol{\Phi}^T \mathbf{z}^t + \mathbf{x}^t$
 - 4: $\hat{\mathbf{x}}^{t+1} = \eta(\boldsymbol{\varepsilon}^t; \kappa + c^t)$
 - 5: $\mathbf{z}^{t+1} = \mathbf{y} - \boldsymbol{\Phi} \hat{\mathbf{x}}^{t+1} + \frac{1}{\gamma} \mathbf{z}^t \langle \eta'(\boldsymbol{\varepsilon}^t; \kappa + \zeta c^t) \rangle$
 - 6: $c^{t+1} = \frac{c^t + \lambda}{\gamma} \langle \eta'(\boldsymbol{\varepsilon}^t; c^t + \kappa) \rangle$
 - 7: **end for**
-

**Figure 2.7:** General linear mixing dealt with GAMP algorithm**2.7.3 Extensions for AMP**

For completeness, some possible extensions of the AMP framework are briefly reviewed in this section. We start with the GAMP algorithm [70]. In Fig. 2.7, a system plot featuring some general input/output channel and linear mixing is depicted. In the plot, $p_{X|S}(x_i|s_i)$ represents a signal prior with some underlying hierarchical structure, i.e. the GMD with S being the hidden states for the variables. The output channel is characterized by the conditional distribution $p_{Y|Z}(y_j|z_j)$ which is not necessarily AWGN, and generates the system output \mathbf{y} . The goal of the GAMP is to estimate \mathbf{x} and \mathbf{z} from the system input vector \mathbf{S} , the output \mathbf{y} and the linear transform $\boldsymbol{\Phi}$. As its name implies, the GAMP provides a unified methodology incorporating essentially arbitrary priors and output non-linearity. Compared to ℓ_1 -AMP and BAMP, the novelty of the GAMP lies in its ability of dealing with arbitrary output distributions $p_{Y|Z}(y_j|z_j)$. Similar to other AMP variants, its derivation is also based on approximation of the message passing algorithm over the system factor graph. The full GAMP algorithm and its derivation would be long and beyond the scope of this thesis. The interested reader is referred to [70] for more details.

Another line of AMP extension is motivated by the lack of practicality of the BAMP algorithm. Although conceptually attractive with its low complexity and optimal MMSE reconstruction performance, in practice, we rarely have the exact signal prior in advance. To overcome this limitation, one possible solution is to incorporate the EM approach to jointly estimate the unknown signal \mathbf{x} and its prior. The resulting algorithm, denoted as the EM-GM-GAMP, is introduced independently by Vila et al. [54] and Krzakala et al. [53]. In [54], the mixture of Gaussians is used as the parametric representation for $p_{\mathbf{x}}$ and the EM method is applied to

estimate the variance and weight for each Gaussian component within one AMP iteration. Extensive simulations over a wide class of distributions confirm the good performance for the EM-GM-GAMP. Later in [73], the adaptive GAMP framework is proposed with an adaptation function for prior estimation at each iteration. In this regard, the EM-GM-GAMP can be seen as a special case of the adaptive GAMP algorithm with the maximum-likelihood (ML) estimation being the adaptation function. More importantly, it is proved that in the large system limit, the adaptive GAMP with the ML parameter estimation yields asymptotically the true value for the signal prior when the distribution satisfies certain identifiability condition. It theoretically provides a rigorous justification for the EM-GM-GAMP algorithm.

2.8 State Evolution Dynamics

As we have mentioned in section 2.4.3, among many advantageous properties that AMP possesses, the state evolution formalism is indubitably the most distinguishable feature compared to all other CS reconstruction algorithms. Essentially the state evolution is a simple iteration which is proved to characterize exactly the asymptotic limit of the AMP estimates as $m, n \rightarrow \infty$ in the case of a Gaussian measurement matrix [14]. It has been stated formally in the following theorem.

Theorem 1. (*[55]*) *Let $\Phi(n)_{n \geq 0}$ be a sequence of sensing matrices $\Phi \in \mathbb{R}^{m \times n}$ indexed by n , with i.i.d entries $\Phi_{ij} \sim \mathcal{N}(0, 1/m)$, and assume $m/n \rightarrow \gamma \in (0, \infty)$. Consider further a sequence of signals $x_0(n)_{n \geq 0}$, whose empirical distributions converge weakly to a probability measure p_{x_0} on \mathbb{R} with bounded $(2k-1)^{\text{th}}$ moment, and assume $\mathbb{E}_{\hat{p}_{x_0}(n)}(\mathbf{X}_0^{2k-2}) \rightarrow \mathbb{E}_{p_{x_0}}(\mathbf{X}_0^{2k-2})$ as $n \rightarrow \infty$ for some $k \geq 2$. Also assume the noise ω has i.i.d. entries with a distribution p_W that has bounded $(2k-2)^{\text{th}}$ moment. Then, for any pseudo-Lipschitz² function $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$ of order k and all $t \geq 0$, almost surely*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \psi(x_i^{t+1}, x_{0,i}) = \mathbb{E}[\psi(\eta_t(X_0 + \tau_t Z), X_0)] \quad (2.60)$$

with $X_0 \sim p_{X_0}$ and $Z \sim \mathcal{N}(0, 1)$ independent.

²Denote the empirical distribution of a vector $x_0 \in \mathbb{R}^n$ by \hat{p}_{x_0} . For $k > 1$ we say a function $\phi : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is pseudo-Lipschitz of order k if there exists a constant $L > 0$ such that, for all $x, y \in \mathbb{R}^m$:

$$|\phi(x) - \phi(y)| \leq L(1 + \|x\|^{k-1} + \|y\|^{k-1})\|x - y\| \quad (2.59)$$

The derivation of the state evolution is inspired by the density evolution in coding theory [74]. The density evolution was first developed for analysing the low-density parity-check (LDPC) codes with iterative decoding. It is known to hold asymptotically for sparse graphs with locally tree-like structure. For CS problems, the underlying factor graph is, in contrast, a fully connected bipartite graph. With some new mathematical ideas, the state evolution is derived as the analog of density evolution in the case of dense graphs.

Another relevant asymptotic analysis for the message passing system is the replica method [47, 75–77]. As a standard statistical physics method, it has been applied successfully to study the typical compressed sensing performance in [78–81]. Although the prediction of the replica method coincides with that of the SE equations, it is not a rigorous approach. In [53], a complete replica analysis for the BAMP is provided without a proof. In this sense, the state evolution provides a theoretical foundation for the replica method based CS work.

2.8.1 State Evolution Heuristics

The detailed proof for Theorem 1 in [55] is well beyond the scope of this thesis. Nevertheless, it is useful to present the simple heuristic description to better understand the dynamics. This section summarizes the basic intuition for AMP in [55] and highlights the key role played by the “Onsager” reaction term in the update equation for \mathbf{z}^t . Recall the generic AMP algorithm is defined previously as

$$\begin{aligned}\mathbf{x}^{t+1} &= \eta_t(\mathbf{x}^t + \Phi^T \mathbf{z}^t) \\ \mathbf{z}^t &= \mathbf{y} - \Phi \mathbf{x}^t + \frac{1}{\gamma} \mathbf{z}^{t-1} \langle \eta'_{t-1}(\mathbf{x}^{t-1} + \Phi^T \mathbf{z}^{t-1}) \rangle\end{aligned}\tag{2.61}$$

The corresponding SE dynamics has the form:

$$\tau_t^2 = \sigma_\xi^2 + \frac{1}{\gamma} \mathbb{E}\{[\eta_t(X_0 + \tau_{t-1}Z) - X_0]^2\}\tag{2.62}$$

We present the argument in [14, 50, 55, 72] to explain the rationale for (2.62). Instead of directly considering the AMP in (2.61), we begin with the following modified recursion:

$$\mathbf{x}^{t+1} = \eta_t(\mathbf{x}^t + \Phi(t)^T \mathbf{z}^t)\tag{2.63}$$

$$\mathbf{z}^t = \mathbf{y}^t - \Phi(t) \mathbf{x}^t\tag{2.64}$$

Comparing to (2.61), we replace the fixed Φ with the independent copy of $\Phi(t)$ at each iteration t , where $\Phi(0), \Phi(1), \dots$ are i.i.d. Gaussian matrices of dimensions $\mathbb{R}^{m \times n}$ with $\Phi_{ij}(t) \sim \mathcal{N}(0, 1/m)$. Consequentially the observation vector at each step is $\mathbf{y}^t = \Phi(t)\mathbf{x} + \boldsymbol{\xi}$. Moreover the last term in the update for \mathbf{z}^t is removed. Eliminating \mathbf{z}^t in (2.63) by plugging in (2.64) gives us a simple recursion

$$\begin{aligned} \mathbf{x}^{t+1} &= \eta_t \{ \Phi(t)^T \mathbf{y}^t + [\mathbf{I} - \Phi(t)^T \Phi(t)] \mathbf{x}^t \} \\ &= \eta_t \{ \mathbf{x} + \Phi(t)^T \boldsymbol{\xi} + \mathbf{A}(t)(\mathbf{x}^t - \mathbf{x}) \} \end{aligned} \quad (2.65)$$

where we define the new operator $\mathbf{A}(t) = \mathbf{I} - \Phi(t)^T \Phi(t)$. Using the central limit theorem, we approximately have $\mathbf{A}_{ij}(t) \sim \mathcal{N}(0, 1/m)$. Because $\mathbf{A}(t)$ is independent of $\mathbf{x}^t - \mathbf{x}$, if we denote $\hat{\tau}_t^2 = \lim_{n \rightarrow \infty} \|\mathbf{x} - \mathbf{x}^t\|^2/n$, then $\mathbf{A}(t)(\mathbf{x}^t - \mathbf{x})$ converges to a vector of i.i.d. entries with zero mean and $\hat{\tau}_t^2/\gamma$ variance. We next consider the statistical property of $\Phi(t)^T \boldsymbol{\xi}$. It is a vector of i.i.d. Gaussian entries with zero mean and $\frac{1}{m} \|\boldsymbol{\xi}\|^2$ variance, which converges to σ_ξ^2 by the law of large numbers. Overall the sum of arguments in $\eta_t(\cdot)$ in (2.65) converges to $X_0 + \tau_t Z$ with $Z \sim \mathcal{N}(0, 1)$ independent of X_0 and

$$\begin{aligned} \tau_t^2 &= \sigma_\xi^2 + \frac{\hat{\tau}_t^2}{\gamma} \\ \hat{\tau}_t^2 &= \lim_{n \rightarrow \infty} \frac{1}{n} \|\mathbf{x} - \mathbf{x}^t\|^2 \end{aligned} \quad (2.66)$$

Given the recursion in (2.65), we have the MSE for \mathbf{x}^{t+1} calculated as

$$\begin{aligned} \hat{\tau}_{t+1}^2 &= \lim_{n \rightarrow \infty} \frac{1}{n} \|\mathbf{x} - \mathbf{x}^{t+1}\|^2 \\ &= \mathbb{E}\{[\eta_t(X_0 + \tau_t Z) - X_0]^2\} \end{aligned} \quad (2.67)$$

Combining (2.66) and (2.67) we finally have the state evolution equation (2.62) for the modified iterative algorithm (2.63).

Note that the whole argument above relies on a crucial assumption: the measurement matrix Φ is draw independently from the same Gaussian distribution at each iteration. However, for CS problems the measurement matrix is constant across iterations. In this scenario, the aforementioned heuristics for the state evolution do not hold because Φ and \mathbf{x}^t are not independent. In fact, with Φ fixed and the soft shrinkage for $\eta_t(\cdot)$, the iteration in (2.62) becomes the IST algorithm. Extensive studies have shown that IST behaves significantly different from the ℓ_1 -AMP and does not follow the state evolution prediction [14, 82].

Algorithmically, IST and ℓ_1 -AMP differ only in the last updating term. Intuitively, this Onsager term acts as the correlation cancellation for Φ and \mathbf{x}^t so that their dependence is neglectable in the large system. As a consequence, the state evolution holds for AMP irrespective of the fact that Φ is kept constant. Moreover, the Onsager term also guarantees the Gaussian behaviour of the effective noise $\hat{\tau}_t^2$.

2.8.2 State Evolution Formula

To end we formally summarize the SE formula for different AMP-based algorithms presented in Section 2.7. Given the CS system in (2.1), we have

The SE dynamics for BP-AMP ([50, 52])

$$\tau_{t+1}^2 = \sigma_\xi^2 + \frac{1}{\gamma} \mathbb{E} [\eta_S(X_0 + \tau_t Z; \zeta \tau_t) - X_0]^2 \quad (2.68)$$

The SE dynamics for BPDN-AMP ([52])

$$\begin{aligned} \tau_{t+1}^2 &= \sigma_\xi^2 + \frac{1}{\gamma} \mathbb{E} [\eta_S(X_0 + \tau_t Z; \lambda + \beta_t) - X_0]^2 \\ \beta_{t+1} &= \frac{\beta_t + \lambda}{\gamma} \mathbb{E} [\eta'_S(X_0 + \tau_t Z; \lambda + \beta_t)] \end{aligned} \quad (2.69)$$

The SE dynamics for BAMP ([52, 53, 56])

$$\tau_{t+1}^2 = \sigma_\xi^2 + \frac{1}{\gamma} \mathbb{E} [f_a(X_0 + \tau_t Z; \tau_t^2) - X_0]^2 \quad (2.70)$$

Similarly, the authors have claimed that the asymptotic behaviour of the components involved in the GAMP iteration can be described by the scalar equivalent model for large Gaussian measurement matrices as well. The parameters for the model can be tracked exactly by the state evolution dynamics. We omit the explicit formula here. Please refer to [70] for its detailed SE equations.

2.9 Summary

This chapter provides an overview of the compressed sensing problem focusing on the AMP-based algorithms. The differences and advantages of AMP over the canonical ℓ_1 -minimization

and iterative thresholding approaches are discussed. A detailed derivation of the BAMP algorithm is presented following Krzakala's procedure, which makes quadratic Gaussian approximation of the standard message passing algorithm in the large system limit. Finally the state evolution dynamics with some intuitive explanation and specific formula is reviewed for AMP-based methods. Although this chapter does not address all the aspects of the CS problem and AMP, it prepares the mathematical background and algorithms that are necessary for the rest of the thesis.

Chapter 3

Sample Distortion Framework for Compressed Sensing

In this chapter, we propose the notion of a SD function for data drawn i.i.d. from compressive distributions to fundamentally quantify the achievable reconstruction performance of compressed sensing for certain encoder-decoder pairs at a given sampling ratio. Two lower bounds on the achievable performance and the intrinsic convexity property is derived. A zeroing matrix is then introduced to improve non-convex SD functions. The SD framework is then applied to analyse compressed imaging with a multi-resolution statistical image model using both the GGD and the two-state GMD. We subsequently focus on the Gaussian encoder-BAMP decoder pair, whose theoretical SD function is provided by the rigorous SE dynamics as explained in Chapter 2. Given the image statistics, analytic bandwise sample allocation for bandwise independent model is derived as a reverse water-filling scheme. Som and Schniter’s turbo approach is further deployed to integrate the bandwise sampling with the exploitation of the hidden Markov tree (HMT) structure of wavelet coefficients. Natural image simulations confirm that with oracle image statistics, the SD function associated with the optimized sample allocation can accurately predict the possible compressed sensing gains. Finally, a general sample allocation profile based on average image statistics not only illustrates preferable performance but also makes the scheme practical.

3.1 Introduction

Traditionally in CS a lot of work has been done in improving reconstruction algorithms assuming the optimality of the homogeneous random sensing matrix. There has recently been more attention on tailoring the sensing matrix in accordance with the signal of interest. In this chapter, we focus on designing a block diagonal measurement matrix for wavelet representation of natural images, which falls under the general scope of bandwise sampling.

Donoho pioneered the use of band-wise sampling for compressed sensing in his original paper [83]. Tsaig further expanded the idea through the concept of two-gender CS, which ran-

domly samples the fine-scale wavelet coefficients while fully samples in the coarse-scale domain [84]. In [85], a specific sampling pattern is provided for the general multi-scale image model. With the key component of weighing the wavelet band importance, it achieves considerable improvement over the homogeneous measurement matrix. However, the weight for each wavelet scale is assigned empirically. Despite all the attempts to improve the measurement matrix, the prior works are algorithmic and lack a solid theoretical grounding.

Analytically optimizing the band-wise sample allocation of the sensing matrix was originally considered in [86] and [87]. The authors sought to minimize the reconstruction uncertainty in terms of the entropy of the CS approximation. However, directly quantifying the entropy is very difficult, thus the authors resorted to an ad hoc solution, which only approximately optimizes the InfoMax criterion [88].

In fact, the notion of optimized band-wise sampling dates back much further and was instrumental in Kashin’s proof of the optimal rates of approximation (n-widths) for certain classes of smooth function [89], which was a key inspiration for the theory of compressed sensing [83]. Specifically, bandwise sampling forms the basis of Maiorov’s discretization theorem which relates function n-widths to the n-widths of a sequence of finite dimensional ℓ_p balls [90].

In other recent work, the block diagonal spatially-coupled sensing matrix was used to reach the fundamental undersampling limit of compressed sensing with almost perfect reconstruction [4], [53], which we will explain in details later in Chapter 4. Unfortunately, to achieve the groundbreaking improvement, a good level of compressibility that we do not normally observe in natural images is required, which makes it impractical for compressed sensing of real images.

Main Contributions

In this chapter, we seek to better understand the nature of good sample allocation strategies for multi-resolution images. To this end, we begin by setting up the sample distortion framework for a stochastic CS model. The SD function is proposed with the purpose of assessing the performance of different encoding and decoding methods quantitatively in terms of the expected MMSE. Then an entropy based bound on the achievable MSE performance for any linear encoder (measurement matrix)-CS decoder (reconstruction algorithm) pair is derived following the classic rate distortion theory. A tighter distribution specific model based bound is further derived by leveraging the entropy based bound of the Gaussian source. We then prove that the SD function is convex in nature. It comes with a key insight: any scheme whose SD function

is concave over the sampling ratio interval $[0, \gamma_c]$ can be improved for any γ in that interval, by sensing a portion of the source at the rate γ_c and making no attempt to sense the remainder. The zeroing procedure which can convexify the SD function comes naturally as a result.

As a broad definition, the SD function is applicable to any encoder-decoder pair, i.e. the Gaussian homogeneous encoder with the linear ℓ_2 decoder or the ℓ_1 minimum CS decoder. In this work, we mainly investigate the SD function for the BAMP decoder. As shown in Chapter 2, the BAMP decoder can be tuned for optimal performance and admits a rigorous analysis in the large system-limit with a large set of sub-Gaussian encoders, which naturally provides the theoretical basis for its SD function [91], [55], [92]. Two compressible distributions: the GGD and the two-state GMD are selected as the representative examples, because they are commonly used models in the compressed imaging literature [93], [94], [95], [96].

The second part of this chapter makes a contribution to the understanding of analytically optimizing the per-band sample allocation for a band-wise independent image model. For this we use an orthogonal wavelet model to make sure our analysis is tractable. We have proved that the optimal sample arrangement with the MMSE is achieved by performing a reverse water-filling strategy, given the per-band statistics and by virtue of the convexified SD function. A similar idea was used in [86] to design the sensing matrix that is most informative about the source. A water-filling strategy is also used in [97] in the context of adaptive sensing. The reconstruction quality can be quantitatively predicted and evaluated by the SD function for the multi-resolution image model. Given the oracle image statistics, our SD function based sample allocation is the best we can achieve in terms of minimizing the MSE. In practice, when the true image statistics is not always available, the performance depends on the quality of the statistical estimation.

Finally wavelet dependencies are incorporated with the band-wise sampling by modelling the wavelet coefficients with the HMT structure [93]. Several works have exploited the local dependencies of the wavelet coefficients in the wavelet based compressed sensing literature, such as [98], [99] and [100]. In this chapter we leverage Som and Schniter's state-of-the-art turbo approach to alternate between the CS decoding and the tree structure decoding [69]. Instead of using a uniform distribution of samples across wavelet bands, we choose the optimized block diagonal sensing matrix to sample independently in the CS decoding procedure. We see that the exploitation of the wavelet tree structure enables the message propagating from coarse scale bands to fine scale bands and eventually benefiting the reconstruction. Attempts are made to find better sample allocation for the tree structure image model. Empirical results are obtained for a specific image example. However, finding the truly best sample allocation for the turbo

method is beyond the scope of this thesis.

The remainder of this chapter is organized as follows. We set up the sample distortion framework in Section 3.2. In Section 3.3 optimizing sample allocation for multi-scale band-independent wavelet image model. The combination of sample allocation with wavelet tree structure is discussed in Section 3.4. Simulation results are given in Section 3.5. Finally, conclusion and future work are discussed in Section 3.6.

3.2 Sample Distortion Framework

3.2.1 Definition

Suppose the signal of interest $\mathbf{x} \in \mathbb{R}^n$ is a random vector (source) with i.i.d. components drawn according to the prior distribution $p(\mathbf{x})$. The goal of statistical compressed sensing is to reconstruct \mathbf{x} using some Lipschitz regular mapping $\Delta : \mathbb{R}^m \rightarrow \mathbb{R}^n$ based on the knowledge of \mathbf{y} , Φ and $p(\mathbf{x})$. In our work, we are interested in the reconstruction quality for certain encoder-decoder pairs (Φ, Δ) at a sampling ratio γ , which is evaluated by the expected error distortion between the original signal \mathbf{x} and the estimation $\Delta(\Phi \mathbf{x})$:

$$D_{\{\Phi, \Delta\}}(\gamma) = \frac{1}{n} \mathbb{E} \|\mathbf{x} - \Delta(\Phi_{\gamma} \mathbf{x})\|_2^2 \quad (3.1)$$

Along the lines of the classical rate-distortion function in the communication field [101], we define a SD function for the compressed sensing setting.

Definition 3. *The SD function is defined as the infimum of sampling ratios for which there is an encoder-decoder pair, (Φ, Δ) , that can achieve an expected distortion D .*

$$D(\gamma) = \inf_{\Phi, \Delta, n} D_{\{\Phi, \Delta\}}(\gamma) \quad (3.2)$$

We will use the term operational SD function to refer to the minimum distortion level a specific encoder-decoder pair can achieve at a fixed sampling ratio for a given compressive source. In this chapter we will concentrate on the Gaussian encoder-BAMP decoder pair. As we summarized in Chapter 2, on the large-system limit assumption with i.i.d. sub-Gaussian Φ , the

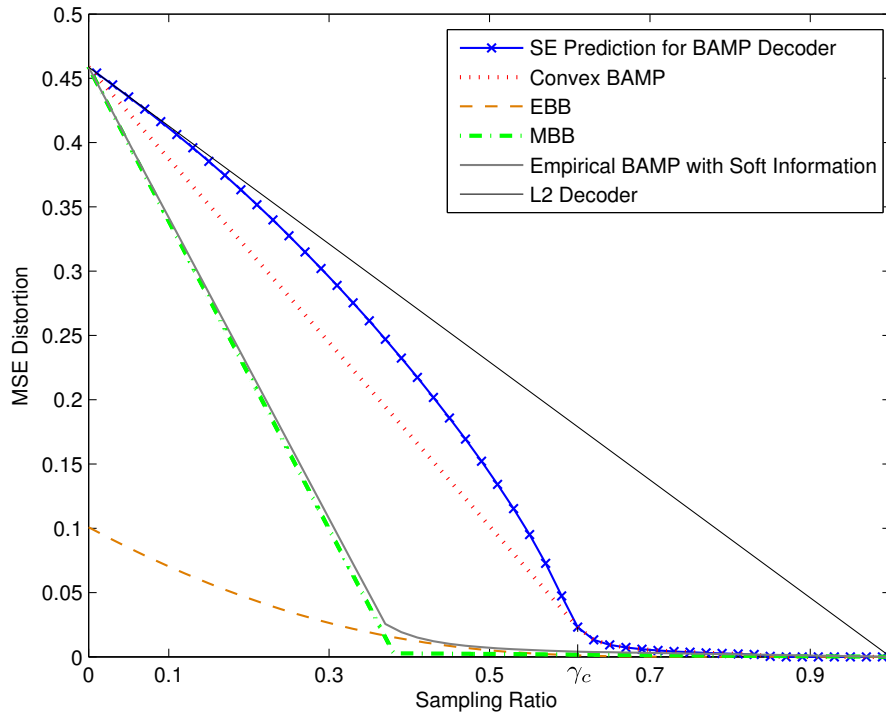


Figure 3.1: *SD functions for GMD data $p(x) = 0.38 \mathcal{N}(0, 1.198) + 0.62 \mathcal{N}(0, 0.004)$ and lower bounds. The critical sampling ratio (defined later in page 47) to convexify this SD function is $\gamma_c = 0.61$.*

distortion iteration can be derived from the SE function [52], [53]¹

$$D_{k+1} = \mathbb{E}\left\{ \left[F\left(\tilde{\mathbf{x}} + \sqrt{\frac{D_k}{\gamma}} \mathbf{z}; \frac{D_k}{\gamma}\right) - \tilde{\mathbf{x}} \right]^2 \right\} \quad (3.3)$$

where $\tilde{\mathbf{x}}$ follows the choice of the compressive distribution, $\mathbf{z} \sim \mathcal{N}(0, 1)$ is independent of $\tilde{\mathbf{x}}$, and $D_0 = \mathbb{E}(\tilde{\mathbf{x}}^2)$. The function $F(\cdot)$ is the (non-linear) scalar MMSE optimal estimator for $\tilde{\mathbf{x}}$ given $\tilde{\mathbf{x}} + \mathbf{z}$. The expectation in (3.3) is taken with respect to $\tilde{\mathbf{x}}$ and \mathbf{z} and is in general calculated numerically. The SD function for BAMP decoder $D_{\text{BAMP}}(\gamma)$ is then given by the convergence point² of (3.3).

3.2.2 Lower bounds

To understand the fundamental theoretical limits of CS for compressible distributions, we now derive two lower bounds for the SD function.

¹When the large-system limit assumption does not hold, there is no analogous results like (3.3). The finite- n case has been studied in a recent work by Rangan et al. [102].

²For the distributions considered in this chapter there is only one non-zero fixed point, i.e. BAMP exhibits no first order phase transition. We will explain this concept further in Chapter 4

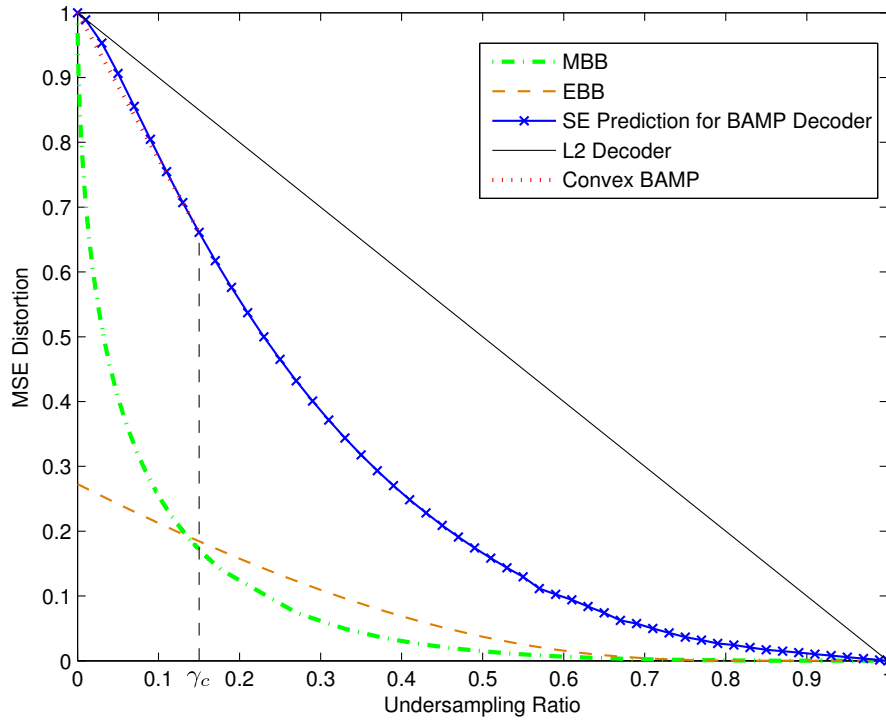


Figure 3.2: SD functions for GGD data $\alpha = 0.4$, $\sigma = 1$ and lower bounds. The critical sampling ratio (defined later in Section. 3.2.3.1) to convexify this SD function is $\gamma_c = 0.15$.

3.2.2.1 entropy based bound

We first prove the *entropy based bound* (EBB) which is a sampling analogy to the classical Shannon Rate Distortion Lower Bound.

Theorem 2. Let $\mathbf{x} \in \mathbb{R}^n$ be a realization of the random vector $\mathbf{x} = x_1, \dots, x_n$, *i.i.d.* $\sim p(x)$, $\text{Var}(x_i) = 1$ and $h(x_i) < \infty$. Let $\mathbf{y} = \Phi \mathbf{x}$, $\mathbf{y} \in \mathbb{R}^m$, $\gamma = m/n < 1$. Then for any Lipschitz reconstruction decoder $\Delta : \mathbb{R}^m \rightarrow \mathbb{R}^n$, we have:

$$D_{\Delta}(\gamma) \geq (1 - \gamma)2^{2(h(x) - h_g)/(1 - \gamma)} \quad (3.4)$$

where $h_g = \frac{1}{2} \log_2 2\pi e$ is the entropy of a unit variance Gaussian random variable.

The proof is given in Appendix D.

Remark 1. The term $h(x) - h_g$ is also known as the *negentropy* of the distribution and is a popular measure of non-Gaussianity, particularly within the field of independent component analysis [103].

Remark 2. When the source \mathbf{x} is Gaussian then the second term in the lower bound becomes

1 and $D_{EBB} = 1 - \gamma$. Here the EBB can be shown to be tight as this corresponds to the SD function for the linear estimator (optimal for Gaussian source): $\hat{\mathbf{x}} = \Phi^\dagger \mathbf{y}$, which is achievable with any full rank linear encoder.

While the EBB in Theorem 2 provides a bound on the achievable performance of CS specifically for i.i.d. sources, it is not clear how close we can expect to get to it. The EBB for two specific GMD and GGD distributions are plotted in Fig. 3.1 and Fig. 3.2. We can see that at low sampling ratios, it is unlikely to be tight. Indeed, for any sparsity promoting decoder, i.e. one for which $\text{supp}(\Delta(\mathbf{y})) \leq \text{dim}(\mathbf{y})$, we know that the MSE cannot exceed that of the best m -term approximation. For such decoders the SD function must therefore approach 1 as $\gamma \rightarrow 0$ [8].

3.2.2.2 model based bound

We next define the *model based bound* (MBB) to compensate for the disadvantage of the EBB. Inspired by the fact that the EBB is tight and achievable for Gaussian source, we resort to the hierarchical Bayesian model to approximate the target compressible distributions. By introducing the variance as a latent variable, the hierarchical representation of a compressive distribution $p(x)$ can be understood as the weighted sum of (possibly infinite) Gaussian distributions.

$$\begin{aligned} p(x) &= \int_0^\infty p(x|\tau)p(\tau) d\tau \\ &= \int_0^\infty \mathcal{N}(x; 0, \tau)p(\tau) d\tau \end{aligned} \tag{3.5}$$

where $p(\tau)$ is the weight for the Gaussian component $\mathcal{N}(x; 0, \tau)$. The MBB is then derived in the following manner: assume the source \mathbf{x} is partitioned into different groups according to the variance. For both encoder and decoder, we agree to transmit and reconstruct the source group by group in the descendant order of the variance. For each Gaussian group, the SD function is tightly bounded by its EBB. Then the lower bound for the whole procedure can be seen as the weighted combination of the EBB of Gaussian components. Thus the MBB has the form:

$$D_{\text{MBB}}(\gamma) = \int_0^c \tau p(\tau) d\tau \tag{3.6}$$

with $\gamma = \int_c^\infty p(\tau) d\tau$.

The two-state GMD model in (2.6) is intrinsically a discretized hierarchical Bayesian model

with only two Gaussian components. Thus its MBB can be seen as the discretized version of the general form given by:

$$D_{\text{MBB}}(\gamma) = \begin{cases} (1 - \lambda)\sigma_s^2 + (\lambda - \gamma)\sigma_l^2 & 0 \leq \gamma \leq \lambda \\ (1 - \gamma)\sigma_s^2 & \lambda < \gamma \leq 1 \end{cases} \quad (3.7)$$

For the GGD model, the detailed procedure for inferring its hierarchical Bayesian prior $p(\tau)$ is relegated to Appendix E. As we can see in both Fig. 3.1 and Fig. 3.2, the MBB is much tighter than the EBB for small sampling ratios, although neither the MBB nor the EBB dominates for the whole range of the sampling ratios. The supremum of the two therefore yields a better lower bound for the SD function.

3.2.3 Convex property

Inspired by the convex property of the rate distortion function, we first prove that the SD function defined in (3.2) is necessarily convex in this section. A direct application of this property is then illustrated to effectively improve the reconstruction quality of the Gaussian encoder-BAMP decoder in the low sample ratio regime.

Theorem 3. *The SD function $D(\gamma)$ in (3.2) is convex.*

Proof. Consider two achievable SD points $(\gamma_1, D(\gamma_1))$ and $(\gamma_2, D(\gamma_2))$. To prove the SD function is convex, we only need to show the convex combination of the two points is also achievable. Let $\gamma_t = t\gamma_1 + (1 - t)\gamma_2$, $0 \leq t \leq 1$. To sample the source $\mathbf{x} \in \mathbb{R}^n$ at the sampling ratio γ_t , we could split \mathbf{x} into two parts $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2]^T$, where $\mathbf{x}_1 \in \mathbb{R}^{tn}$, $\mathbf{x}_2 \in \mathbb{R}^{(1-t)n}$, and apply encoders with sampling ratio γ_1, γ_2 to $\mathbf{x}_1, \mathbf{x}_2$, respectively. Then the reconstruction of \mathbf{x}_1 and \mathbf{x}_2 has achievable MSE: $tnD(\gamma_1)$ and $(1 - t)nD(\gamma_2)$. So the MSE of the reconstruction of X is:

$$nD(\gamma_t) \leq tnD(\gamma_1) + (1 - t)nD(\gamma_2) \quad (3.8)$$

Therefore

$$D(t\gamma_1 + (1 - t)\gamma_2) \leq tD(\gamma_1) + (1 - t)D(\gamma_2) \quad (3.9)$$

□

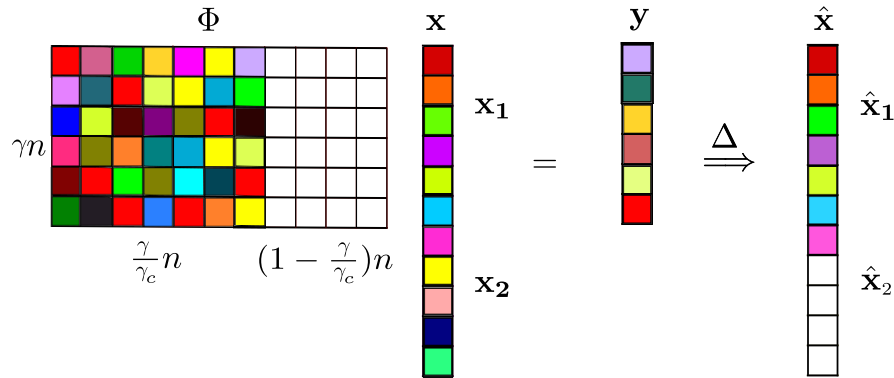


Figure 3.3: Hybrid zeroing Gaussian matrix as the convex combination of a trivial decoder $\hat{\mathbf{x}} = 0$ and a BAMP decoder Δ . Elements equal to 0 are represented with white blocks.

3.2.3.1 hybrid zeroing matrix

The convexity property is applicable to the operational SD function for any specific encoder-decoder pair in the large-system limit. The application of Theorem 3 is that for a given encoder-decoder pair with a concave operational SD function between γ_1 and γ_2 ($\gamma_1 < \gamma_2$), there exists a hybrid system with better SD performance: it can be easily achieved by applying the two encoder-decoders to different portions of the source to get the convex combination of $D(\gamma_1)$ and $D(\gamma_2)$. A special case is when $\gamma_1 = 0$ with the corresponding trivial decoder ($\hat{\mathbf{x}} = 0$) and $\gamma_2 = \gamma_c$ with γ_c being the crucial sampling ratio. In this case, instead of sampling the source \mathbf{x} with a full Gaussian matrix, $\Phi \in \mathbb{R}^{\gamma n \times n}$, we split \mathbf{x} as before with $\mathbf{x}_1 \in \mathbb{R}^{tn}$ and $\mathbf{x}_2 \in \mathbb{R}^{(1-t)n}$, $t = \gamma/\gamma_c$. We then sample \mathbf{x}_1 with the Gaussian matrix, $\tilde{\Phi} \in \mathbb{R}^{\gamma n \times tn}$ and reconstruct, while the remaining \mathbf{x}_2 we reconstruct as zero. Since this is equivalent to setting part of the encoder to zero, $\Phi = [\tilde{\Phi}, \mathbf{0}]$, we call this the *zeroing procedure*, as illustrated in Fig. 3.3.

Close observation of the operational SD functions for the Gaussian encoder-BAMP decoder system in Fig. 3.1 and Fig. 3.2 reveals that the curves are convex for large sampling ratios but concave for small sampling ratios. By applying the hybrid zeroing Gaussian matrix, we convexify the SD function for γ below the critical sampling ratio γ_c .

Definition 4. To best improve the SD performance, γ_c is chosen as the largest sampling ratio below which the SD function is concave.

The Gaussian sensing matrix has been widely assumed within the CS community to be optimal in terms of CS performance. Indeed this has been proved to be the case for the distributions that exhibit exact sparsity [104]. However, under the assumption that the BAMP achieves the

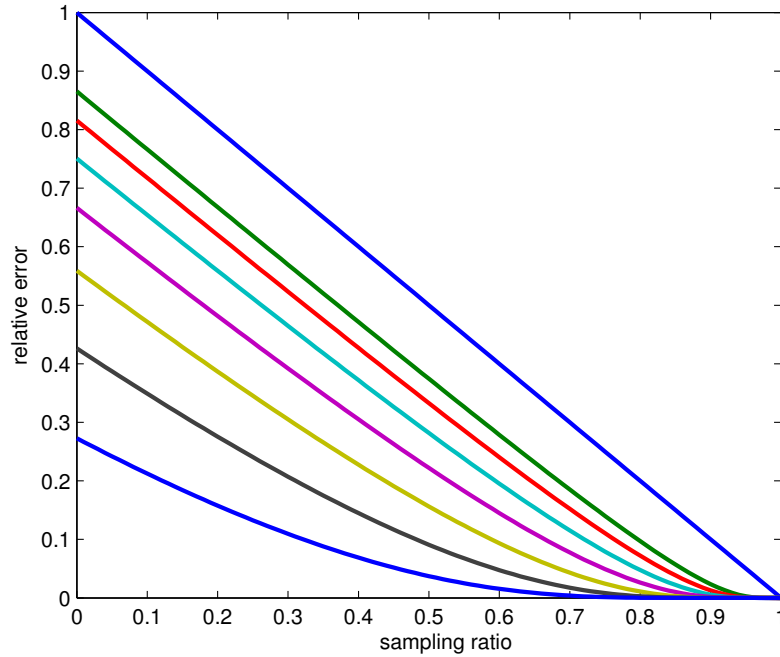


Figure 3.4: EBB for GGD model with $\alpha = 0.4$ (left most curve), $0.5, \dots, 1.0$ and $\alpha = 2$.

Bayes optimal reconstruction - this would follow, for example, if the replica method could be proved to be rigorous [53] - then the *zeroing procedure* resulting from Theorem 3 indicates this assumption to be false.

3.3 Measurement Matrix for Multi-resolution Image Model

In this section we build upon the aforementioned SD framework and study the SD behaviour of the compressive imaging. We investigate the optimal band-wise sampling strategy with a fixed sample budget, in a similar manner to [86], but in terms of minimizing the expected MSE. We begin by introducing the bandwise independent multi-resolution statistical model for natural images.

3.3.1 Compressible Distributions

Given that we consider the CS problem in the stochastic setting, the probability density function to characterize the compressible signals is required. For this section, we consider two specific non-Gaussian distributions introduced in Chapter 2, the two-state GMD (2.6) and the GGD (2.5), to model the wavelet coefficients of natural images.



Figure 3.5: Ten 512×1024 HDR Images. From left to right, top to bottom: Chapel, Dog, Pine, Sea, Man, Wedding, Hill, Penguin, Room, Sign.

Fig. 3.4 shows a plot of the EBB for GGD distribution with different shape parameter α . While the right most curve, $D_{\text{EBB}} = 1 - \gamma$, is always achievable, in CS we are mainly interested in distributions that can be well approximated at significant undersampling ratios, i.e. we want $D \ll 1$ and $\gamma \ll 1$ simultaneously. From Fig. 3.4 we can conclude that the Laplace distribution ($\alpha = 1$) cannot admit such a low distortion at significant undersampling ratios. Indeed, low SD functions appear to require very small values of $\alpha \sim 0.4$. For images we are typically interested in the GGD with $\alpha \sim [0.3, 1]$ since these distributions provide a good approximation for the distribution of the wavelet coefficients in a given band for natural images. This is illustrated in Fig. 3.6 (a) and Fig. 3.9 (b).

Examples of the theoretical prediction for the SD function of GMD and GGD data using BAMP decoder can be found in Fig. 3.1 and Fig. 3.2 respectively. The function $F(\cdot)$ has a close-form expression for the GMD [53], [69] and can be solved numerically for the GGD.

3.3.2 Band-wise Independent Image Model

Natural images are typically transform compressible: they have a more concise representation in the wavelet domain. The wavelet decomposition of an image $f(\mathbf{X})$ has the form [105]:

$$f(\mathbf{X}) = \sum_k \mu_{i,k} \phi_{i,k}(\mathbf{X}) + \sum_{j \geq i, k} \omega_{j,k} \psi_{j,k}(\mathbf{X}) \quad (3.10)$$

where $\phi_{i,k}(\mathbf{X}) = 2^{\frac{i}{2}} \phi(2^i \mathbf{X} - k)$ are the scaling functions, $\psi_{j,k}(\mathbf{X}) = 2^{\frac{j}{2}} \psi(2^j \mathbf{X} - k)$ are the prototype bandpass functions such that together they form an orthonormal basis. The variables $\mu_{i,k}$ are in turn the scaling coefficients at scale i and $\omega_{j,k}$ are the wavelet coefficients at scale j . We can group the coefficients into a single vector according to the scale and assign each a band index: denote the scaling coefficients as band 0, the coarsest wavelet coefficients as band 1 and so on. In this manner we obtain the vector $\boldsymbol{\theta} = [\mu_0, \mu_1, \mu_2, \dots]$. Next we follow [9], [10]

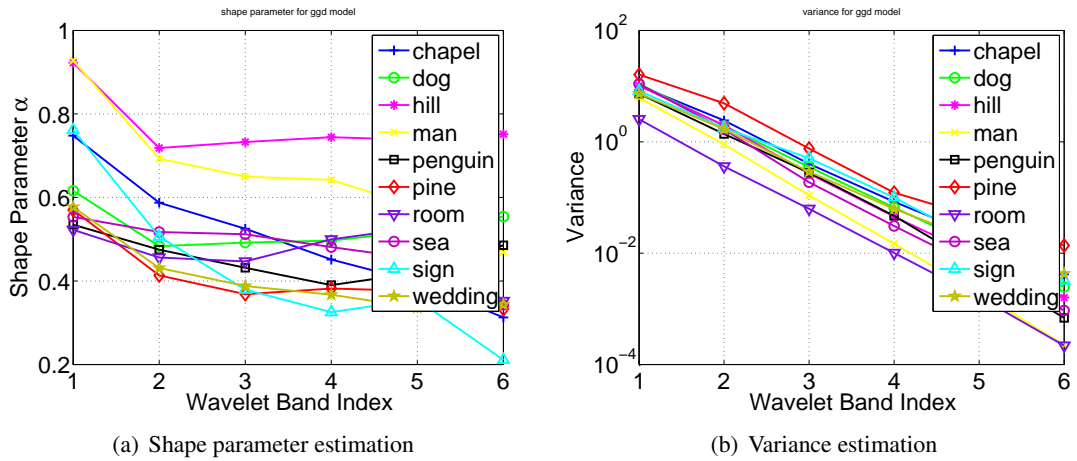


Figure 3.6: GGD parameters for six wavelet bands of HDR images in Fig. 3.5.

and consider a simple statistical model defined directly on the wavelet coefficients. The band 0 is always treated as Gaussian since these coefficients typically exhibit no sparsity. This can be seen as a worse case assumption in terms of its SD function. For the other bands, we model the wavelet coefficients within each band as mutually independent and impose a compressive distribution for each wavelet band. To be specific, $\omega_{j,k}$ at scale j can be modelled as

$$\omega_{j,k} \sim \text{GGD}(0, \sigma_j^2, \alpha_j) \quad (3.11)$$

or

$$\omega_{j,k} \sim \text{GMD}(\lambda_j, \sigma_{L,j}^2, \sigma_{S,j}^2), \quad (3.12)$$

where typically for natural images the distributions exhibit a self-similar structure with an exponential decay across scale, i.e. $\sigma_j^2 = 2^{-j\beta} \sigma_0^2$ for the GGD and $\sigma_{a,j}^2 = 2^{-j\beta} \sigma_{a,0}^2$, $a = S, L$ for the two-state GMD for some $\beta > 0$. For the bandwise independent image model, we assume an uniform activity rate λ_j for each wavelet band in spite of the coefficient index. In particular, we define $\lambda_j := \Pr\{s_{j,k} = 1\}$.

Extensive statistic studies for both natural images (Fig. 3.8) and high resolution high-dynamic range (HDR) images (Fig. 3.5) are conducted and presented in Fig. 3.6, Fig. 3.7 and Fig. 3.9. The parameters for GGD and GMD model are derived through moment matching. The log-log scale plots of the variance confirm the power law decay assumption. And the shape parameter estimation agrees with our previous analysis for the GGD model that α normally has the value between 0.3 and 1 for images.

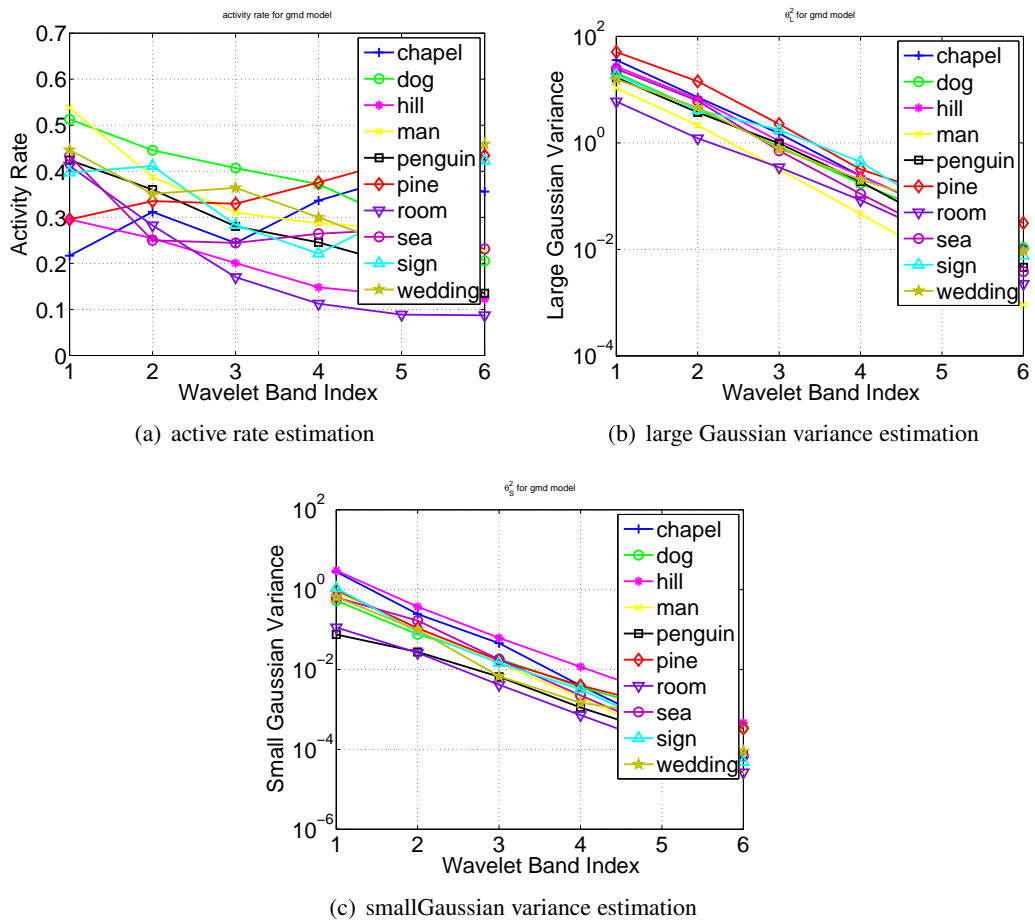


Figure 3.7: Two-state GMD parameters for 6 wavelet bands of HDR images in Fig. 3.5.

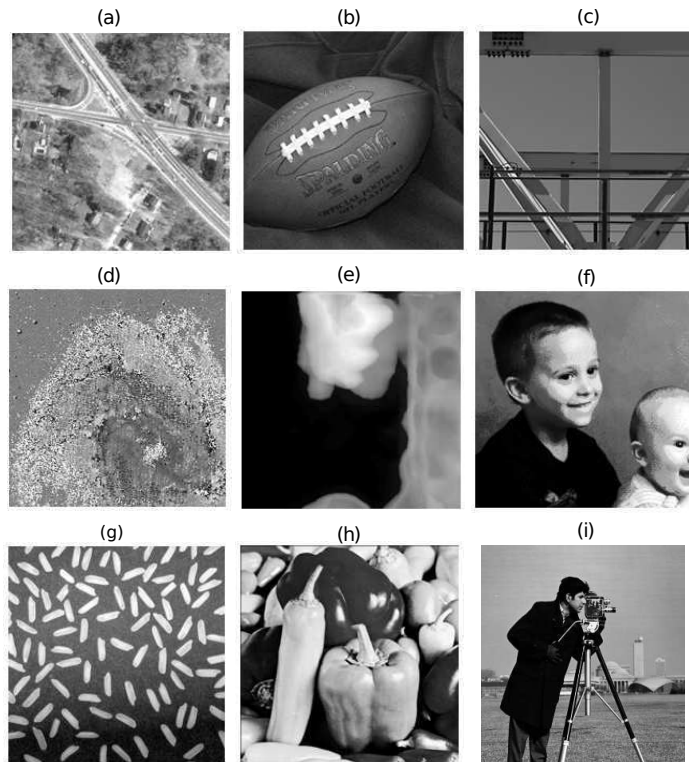


Figure 3.8: *Nine 256×256 Natural Images: (a)Concordant (b) football (c) Gantry Crane (d) M83 (e) Spine (f) Kids (g) Rice (h) Peppers (i) Cameraman*

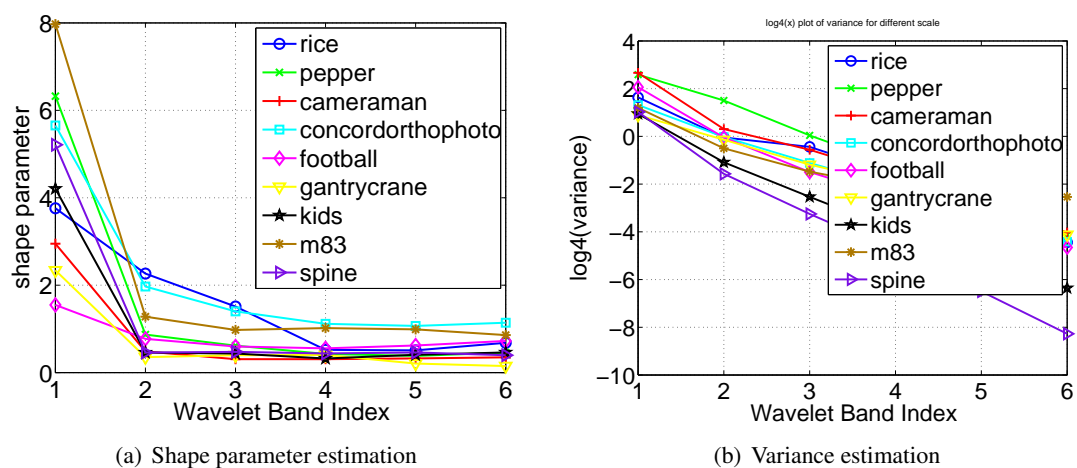


Figure 3.9: *GGD parameters for 6 wavelet bands of natural images in Fig. 3.8.*

3.3.3 Band-wise Sampling

3.3.3.1 Sample Allocation Strategy

To keep things tractable we restrict ourselves to the class of linear encoders, $\mathbf{y} = \Phi\theta$, that are block diagonal and sample the different wavelet bands separately with the following form:

$$\Phi = \begin{pmatrix} \Phi_0 & & & \\ & \Phi_1 & & \\ & & \ddots & \\ & & & \Phi_L \end{pmatrix} \quad (3.13)$$

where $\Phi_i \in \mathbb{R}^{m_i \times n_i}$, $m_i \leq n_i$ puts m_i measurements to sample the i th band. The equality holds when the i th band is fully sampled with Φ_i being an identity matrix. Otherwise Φ_i is a possibly zero padded (for convexity) Gaussian random matrix. And $\underline{\mathbf{y}}_i = \Phi_i \underline{\omega}_i$ is the CS observation for each block. To derive the SD function for the multi-resolution images, we first consider the L wavelet bands as independent and parallel. The question then is how to allocate a fixed number of samples to the various bands, with the aim of minimizing the total reconstruction distortion. First let us assume for now that m_i, n_i be continuous and $\gamma_i = m_i/n_i \in [0, 1]$. The problem is reduced to the following optimization

$$\begin{aligned} \min_{m_i} & \sum_{i=1}^L \sigma_i^2 n_i D_i(m_i/n_i) \\ \text{s.t.} & \sum_{i=1}^L m_i = m \text{ and } 0 \leq m_i \leq n_i, \quad i = 1, \dots, L. \end{aligned} \quad (3.14)$$

where D_i is the (convex) SD function for band i normalized to have unit variance. Using Lagrange multipliers, we construct the objective function

$$\begin{aligned} L = & - \sum_i \sigma_i^2 n_i D_i(m_i/n_i) - \lambda (\sum_i m_i - m) \\ & - \sum_i \mu_i (m_i - n_i) + \sum_i \nu_i m_i \end{aligned} \quad (3.15)$$

Differentiating with respect to m_i and setting equal to 0 we have

$$\frac{\partial L}{\partial m_i} = -\sigma_i^2 n_i \frac{\partial D_i}{\partial \gamma_i} \cdot \frac{\partial \gamma_i}{\partial m_i} - \lambda - \mu_i + \nu_i = 0 \quad (3.16)$$

or

$$-\sigma_i^2 \frac{\partial D_i}{\partial \gamma_i} - \lambda - \mu_i + \nu_i = 0 \quad (3.17)$$

Define the distortion reduction function as

$$\eta_i(\gamma_i) = -\sigma_i^2 \frac{\partial D_i}{\partial \gamma_i}, \quad (3.18)$$

noting that this function is non-increasing in terms of γ_i . Now applying the Kuhn-Tucker (KT) conditions we arrive at:

$$\eta_i(\gamma_i) - \lambda - \mu_i + \nu_i = 0, \quad (3.19)$$

with

$$\mu_i(n_i - m_i) = 0, \quad \mu_i \geq 0, \quad (3.20)$$

and

$$\nu_i m_i = 0, \quad \nu_i \geq 0. \quad (3.21)$$

We therefore have three cases for the distortion reduction function. First, if $0 < m_i < n_i$ then $\mu_i = \nu_i = 0$ and the sampling ratio, γ_i , is set so that $\eta_i(\gamma_i) = \lambda$. Next suppose that $m_i = n_i$ so that $\gamma_i = 1$. In this case, the KT conditions imply that

$$\eta_i(\gamma_i) \geq \lambda, \quad \forall \gamma_i \quad (3.22)$$

In the final case we have $m_i = 0$ and $\gamma_i = 0$. Here the KT conditions imply:

$$\eta_i(\gamma_i) \leq \lambda, \quad \forall \gamma_i \quad (3.23)$$

This gives us an optimal sample allocation strategy which is similar to the reverse water-filling idea in rate distortion theory [106]. We allocate samples to the band with the greatest distortion reduction value until another band has a greater one or that band has been fully sampled. The procedure is stopped when the total distortion reaches the desired level.

To apply this idea to natural images we need to take account of the fact that m_i , n_i and L are all discrete and finite. Thus we define a discretized distortion reduction (DR) function for each wavelet band.

$$\eta_i(m_i) = \sigma_i^2 [D_i(m_i/n_i) - D_i((m_i + 1)/n_i)] \quad (3.24)$$

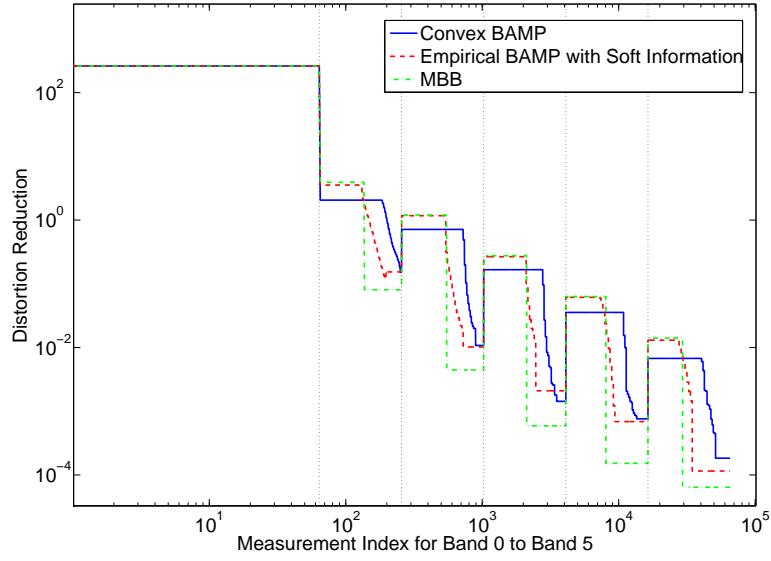


Figure 3.10: *Distortion reduction function of six bands Daubechies 2 wavelet decomposition of cameraman image using GMD model (including the low-pass band). The statistics is reported in Table 3.1 in page 63.*

Suppose that m_i samples have been allocated to the i th band. The DR function gives the amount of distortion decreased by adding one more sample to that band. Then the number of samples allocated to the band i is

$$m_i = \begin{cases} 0 & \text{if } \max \eta_i(m_i) < \kappa \\ n_i & \text{if } \min \eta_i(m_i) > \kappa \\ \hat{m}_i \text{ s.t. } \eta_i(\hat{m}_i) = \kappa & \text{otherwise} \end{cases} \quad (3.25)$$

where κ is chosen so that $\sum_i m_i = m$. With a convex SD function, the optimal allocation is again achieved by performing a greedy sample allocation strategy. The DR function for a six-band Daubechies 2 decomposition of the “cameraman” using the two-state GMD model is illustrated in Fig. 3.10. One thing worth noting is that neither the convexity property nor the resulting greedy sample allocation method is restricted to the form of the decoder. For example the optimized bandwise sensing matrix can be designed in the same manner for the CS ℓ_1 and ℓ_2 decoder. The consequential SD lower bounds can be obtained as well as demonstrated in Fig. 3.10.

3.3.3.2 Comparison to the Theory of Widths

In [9], parallels are drawn between the statistical wavelet model we have considered here and the family of Besov function spaces. In particular, the authors argue that under appropriate conditions realizations drawn from the GMD or GGD based wavelet model almost surely lie in an associated Besov space. It is therefore interesting to explore the similarities and differences between the achievable distortion rates derived here and those known in the deterministic setting for Besov spaces.

n-widths of Besov spaces Consider the Lipschitz class of r -smooth functions on the interval $[0, 1]$ and the unit ball, B_p^r , defined as:

$$B_p^r := \{f : \|f^{(r)}\|_p \leq 1\} \quad (3.26)$$

where $f^{(r)}$ denotes the r th derivative of f and the L_p ball acts as the deterministic counterpart to the coefficient prior above.

The ℓ_2 error of the best n -dimensional linear approximation for these spaces is known to scale as $\sim n^{-r+1/p-1/2}$ for $1 \leq p \leq 2$ [107, Chapter 14, Theorem 1.1]. In contrast, the ℓ_2 error for the best CS reconstruction is characterized by the Gelfand width of B_p^r which can be written as:

$$d^n(B_p^r) := \inf_{\Phi} \sup_h \{\|h\|_2, h \in \mathcal{N}(\Phi) \cap B_p^r\}. \quad (3.27)$$

and measures the uncertainty in B_p^r within the null space of Φ . Here, for $1 \leq p \leq 2$ the best CS approximation error decays at the faster rate of $\sim n^{-r}$, i.e. inversely proportional to the smoothness [107, Chapter 14, Theorem 1.1]. This result was derived in Kashin's seminal paper [89], which is better known in the CS community for accurate bounds for the n -widths of l_p balls in \mathbb{R}^n .

Similarities and differences Interestingly Kashin's result relied on a discretization theory of Maiorov [90] that uses a similar bandwise sampling to our own. Specifically Maiorov uses a subband decomposition of spline spaces to bound the n -width of B_p^r in terms of a weighted sum of finite dimensional n -widths for the individual subbands - effectively performing a bandwise sampling. Furthermore in both the deterministic and stochastic settings the allocation scheme is broadly the same: fully sample the first few low resolution subbands; then partially sample a number of intermediate subbands; and finally set coefficients of all the higher resolution sub-

bands to zero. However, in Kashin’s theory, the number of partially sampled subbands grows as the distortion decreases and, indeed, it is this that accounts for the different rate of approximation compared with the best linear approximation. In contrast, in the sample allocation framework, the number of partially sampled bands, P , is bounded by the range of the distortion reduction function:

$$P < \beta \log_2(\eta(0)/\eta(1)). \quad (3.28)$$

For the two-state GMD model this bound is finite since from the MBB we can deduce that:

$$\frac{\eta(0)}{\eta(1)} < \frac{\sigma_{L,0}^2}{\sigma_{S,0}^2} \quad (3.29)$$

Note the same bound applies to the SD function for the MBB oracle decoder where the band-wise sampling is optimal. Hence, the fact that we do not get a growing number of partially sampled subbands implies that in the large system limit the CS approximation error will decay at the same rate as for the best linear approximation. We can therefore conclude that the gains in CS solutions over optimal linear approximation for such a model are fundamentally limited. We can see this, for example, in Fig. 3.10 where we would only ever partially sample at most 3 subbands for the convexified BAMP decoder.

3.4 Sample Allocation with Tree Structure

Until now we have developed an analytic sample allocation method for a multi-resolution image model by assuming the independence of the wavelet band. In this subsection we look beyond the signal sparsity and incorporate the wavelet dependencies with the aim of getting closer to the model based bound. We will start with the review of the wavelet quad-tree structure. Then the HMT-based compressed imaging is introduced and the turbo inference scheme is presented as the tool. Finally the combination of the sample allocation and the turbo reconstruction is discussed.

3.4.1 Hidden Markov Tree Model

Beside the primary properties, i.e. multi-resolution and compressibility, the wavelet coefficients are well known for some secondary properties, one of which is known as *persistence across scale* (PAS) [93]. When modelled with the two-state GMD, the wavelet coefficients for 2D images naturally form a quad-tree structure, as illustrated in Fig. 3.11(a). Except the “root”

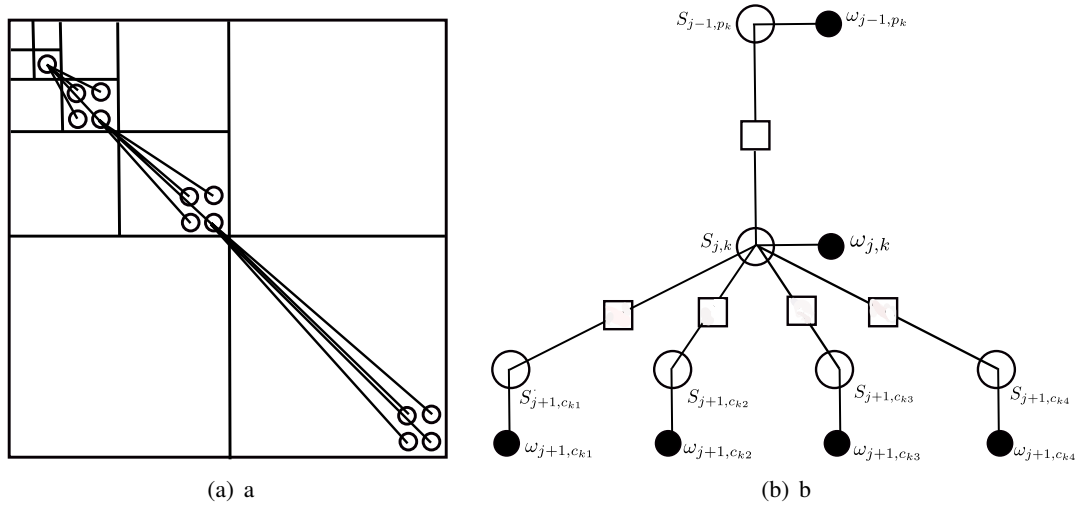


Figure 3.11: (a) An illustration of the image quad-tree structure. (b) A zoomed in factor graph of the HMT structure featuring a typical variable node (the hidden state) $s_{j,k}$ connected with its four children $\{s_{j+1,c_{ki}}\}_{i=1}^4$ and parent node s_{j-1,p_k} by the factor nodes (the transition matrix).

(wavelet coefficients at the coarsest scale, also noted as band 1), each wavelet coefficient has a “parent” in the upper wavelet scale and serves as the parent for four “children” in the next scale. The HMT connects the hidden states across scale and readily models the PAS [93]. Here we repeat the GMD model in (2.6) for the wavelet coefficient $\omega_{j,k}$.

$$p_{\text{GMD}}(\omega_{j,k}) = p(s_{j,k} = 1)\mathcal{N}(\omega_{j,k}; 0, \sigma_{j,L}^2) + p(s_{j,k} = 0)\mathcal{N}(\omega_{j,k}; 0, \sigma_{j,S}^2) \quad (3.30)$$

$$= \lambda_{j,k}\mathcal{N}(\omega_{j,k}; 0, \sigma_{j,L}^2) + (1 - \lambda_{j,k})\mathcal{N}(\omega_{j,k}; 0, \sigma_{j,S}^2) \quad (3.31)$$

The PAS property states that if the parent is large, some of its children are likely to be large; if the parent is small, all of its children tend to be small. In other words, the activity rate $\lambda_{j,k}$ for $\omega_{j,k}$ depends on the activity rate of its parent on scale $j - 1$, λ_{j-1,p_k} . To take the Bayesian approach, we model the parent-children relationship across scale by the transition matrix T_j .

$$T_j = \begin{bmatrix} p(s_{j+1} = 0 | s_j = 0) & p(s_{j+1} = 0 | s_j = 1) \\ p(s_{j+1} = 1 | s_j = 0) & p(s_{j+1} = 1 | s_j = 1) \end{bmatrix} = \begin{bmatrix} t_j^{00} & 1 - t_j^{11} \\ 1 - t_j^{00} & t_j^{11} \end{bmatrix} \quad (3.32)$$

where t_j^{00} is the probability of a child state at scale $j + 1$ being 0 given its parent’s state at scale j is 0. Similarly, t_j^{11} is the probability of a child state being 1 if its parent’s state is 1. For the HMT structure, the activity rate for a child wavelet $\omega_{j+1,c_{ki}}$ can be calculated with the

transition matrix and the activity rate of its parent

$$\begin{bmatrix} 1 - \lambda_{j+1, c_{ki}} \\ \lambda_{j+1, c_{ki}} \end{bmatrix} = T_j \begin{bmatrix} 1 - \lambda_{j, k} \\ \lambda_{j, k} \end{bmatrix} \quad (3.33)$$

In summary, for a L level wavelet decomposition, we can specify the HMT structure with the following set of parameters

$$\Theta = \left[\{\lambda_{1,k}\}_{k=1}^{n_1}, \{t_j^{00}\}_{j=1}^{L-1}, \{t_j^{11}\}_{j=1}^{L-1}, \{\sigma_{j,L}^2\}_{j=1}^L, \{\sigma_{j,S}^2\}_{j=1}^L \right]^T \quad (3.34)$$

where $\{\lambda_{1,k}\}_{k=1}^{n_1}$ are the activity rates for the root coefficients. In our work, we assume the variances for the Gaussian components are known. We can either estimate the variances from the wavelet coefficients or adopt a general image model. Then the parameter set is reduced to

$$\hat{\Theta} = \left[\{\lambda_{1,k}\}_{k=1}^{n_1}, \{t_j^{00}\}_{j=1}^{L-1}, \{t_j^{11}\}_{j=1}^{L-1} \right]^T \quad (3.35)$$

We can further treat the activity rate and the transition probability as random variables and impose some pdf to complete the statistical model for the HMT structure. In [69], Beta and Gamma hyperpriors are assumed.

One of the important applications of HMT is the signal estimation from noisy observation. The denoising algorithm was introduced in [93] and can be summarized as a two-step procedure: Given the noisy data, we first fit an HMT model Θ to the data. Then we use the model as the prior to compute the conditional mean as the denoised estimate. Since the factor graph of the HMT model has a loop-free structure, the exact calculation of the posteriors can be obtained using two passes of the sum-product algorithm [62]. In [93], the upward-downward algorithm was introduced for the model fitting. The denoising power of the HMT model can serve as an assistant for the compressed imaging. It will help further enhance the confidence about the activity rate for each coefficients through the dependency across wavelet levels, and thus should improve the image reconstruction.

3.4.2 Turbo Decoding

Several authors have investigated the HMT-aided compressed sensing reconstruction: In [98], the Markov-chain Monte-Carlo (MCMC) techniques are exploited; In [108], the Variational Bayes based approach is introduced. In this work, we focus on a HMT-based compressive imag-

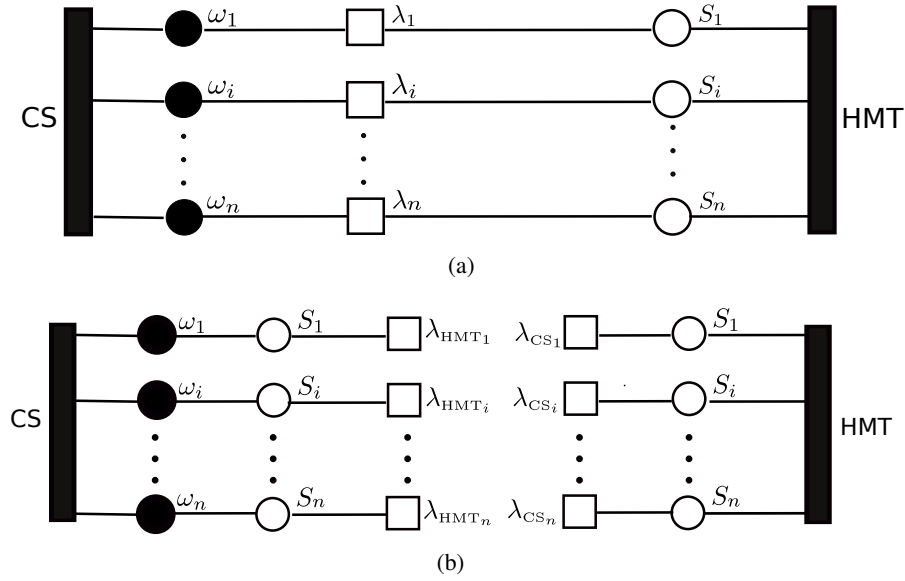


Figure 3.12: Top: factor graph for the compressive imaging with HMT structure. Bottom: two sub-graphs for the turbo decoding. The likelihood from one sub-graph is used as the prior for the other sub-graph.

ing scheme based on the loopy belief propagation (LBP), first proposed by Schniter in [109]. It has been shown to have the state-of-art reconstruction performance with a low complexity.

In the Bayesian compressed sensing setting, the reconstruction of the wavelet coefficients ω from the CS observation \mathbf{y} is interpreted as approximating the posterior mean of the density $p(\omega|\mathbf{y})$. When introducing dependencies across wavelet scales, the factor graph for the whole reconstruction system has a loopy structure as illustrated in 3.12(a). Although exact inference of $p(\omega|\mathbf{y})$ is known to be NP hard, the marginal posterior $p(\omega_i|\mathbf{y})$ can be approximated using the LBP. In [109] and [69], the LBP is conducted through the “turbo” decoding approach, which we summarize as follows.

To perform the turbo decoding, we first split the factor graph for the whole system as two decoupled sub-graphs, with one representing the compressed sensing mixing and the other exploiting the HMT structure, as shown in Fig. 3.12(b). The essence of turbo decoding is to exchange the local belief of the hidden state $s_{j,k}$ between the CS decoding and HMT decoding alternately. To be specific, the likelihood on $s_{j,k}$ derived from the HMT decoding is treated as the prior for the active rate in the CS decoding. When the CS reconstruction terminates, the resulting likelihood on the active rate is used for the next round HMT decoding. The turbo reconstruction converges when both decoding procedures terminate. The CS decoding is performed through the AMP based algorithms. For the HMT decoding, the aforementioned upward-downward algo-

rithm can be applied. Although there is no rigorous convergence analysis for the turbo scheme, Schniter et al. have demonstrated some promising results for compressed imaging in [69].

3.4.3 Sample Allocation with Tree Structure

In this subsection, we discuss the application of the turbo decoding approach with the sample allocation strategy. Let $\boldsymbol{\omega} = [\omega_1, \omega_2, \dots, \omega_L]^T$ denote the collection of the wavelet coefficients of different bands and $\mathbf{s} = [s_1, s_2, \dots, s_L]^T$ be the corresponding hidden states vector. Assume $\mathbf{y} = [y_1, y_2, \dots, y_L]^T$ is the CS observation vector using the block diagonal sensing matrix. Then the posterior $p(\boldsymbol{\omega}|\mathbf{y})$ has the form:

$$\begin{aligned} p(\boldsymbol{\omega}|\mathbf{y}) &= Z^{-1} p(\mathbf{y}|\boldsymbol{\omega}) \sum_{\mathbf{s}} p(\mathbf{s}) p(\boldsymbol{\omega}|\mathbf{s}) \\ &= Z^{-1} \sum_{\mathbf{s}} p(\mathbf{s}) \prod_j [\prod_t p(y_{j,t}|\underline{\boldsymbol{\omega}}_j)] [\prod_k p(\omega_{j,k}|s_{j,k})] \end{aligned} \quad (3.36)$$

where $Z = p(\mathbf{y})$. The factor graph plotted in Fig. 3.13 visualizes this global function [58], [60]. Here, unlike [69], the AMP decoder is bandwise independent due to the block diagonal form of Φ . The interaction across different wavelet bands only comes from the HMT decoding.

The SD function for the bandwise independent image model is unlikely to be optimal for the turbo decoding scenario since it does not take the HMT decoding into consideration. The role of the HMT decoding is to better provide estimation of the activity rate $\lambda_{j,k}$ for the scalar MMSE estimator of each wavelet coefficient, instead of using an identical λ_j over the coefficient index k , thus improving the reconstruction quality. To see the impact of the HMT decoding, we feed the BAMP decoder with the *soft information*, $\hat{\lambda}_{j,k}$, defined as follows:

$$\begin{aligned} \hat{\lambda}_{j,k} &= \frac{p(\omega_{j,k}|s_{j,k} = 1)}{p(\omega_{j,k}|s_{j,k} = 1) + p(\omega_{j,k}|s_{j,k} = 0)} \\ &= \frac{\mathcal{N}(\omega_{j,k}; 0, \sigma_{j,L}^2)}{\mathcal{N}(\omega_{j,k}; 0, \sigma_{j,L}^2) + \mathcal{N}(\omega_{j,k}; 0, \sigma_{j,S}^2)} \end{aligned} \quad (3.37)$$

This provides a soft estimate of the state of the GMD and thereby gives a better prediction of individual coefficient variances. The empirical SD curve for the BAMP decoder with soft information is generated from the Monte Carlo simulations with synthetic GMD data and illustrated in Fig. 3.1 page 43. To be specific, we use the $\hat{\lambda}_{j,k}$ in (3.37) instead of λ_j for the scalar MMSE estimator of each synthetic GMD component. Fig. 3.1 demonstrates that providing the BAMP decoder with good estimation of activity rate information dramatically improves the

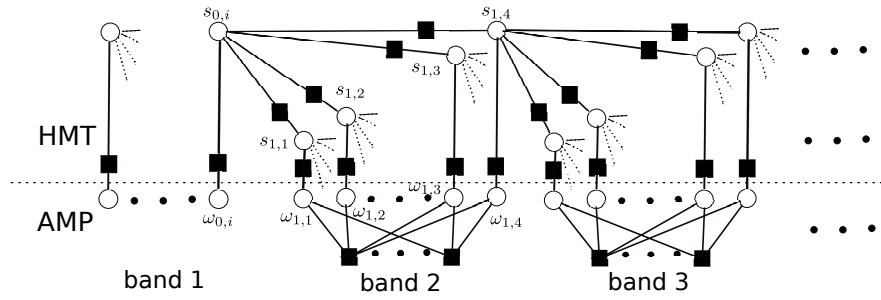


Figure 3.13: Factor graph for band-wise sampling with HMT decoding. The upper graph illustrates a quad-tree structure of the wavelet hidden states. The lower graph is the band-wise independent random mixing.

reconstruction quality, with the SD function lying very close to the lower bound.

Based on the per-band image statistics, the SD function for BAMP decoder with soft information can be obtained empirically for each wavelet band in the same fashion. Then the DR function with soft information for the multi-resolution image model can be established following the aforementioned definition, as shown in Fig. 3.10 page 55. To clarify the terminology, we denote the corresponding sample allocation profile as the HMT based sample allocation, or HSA. And we use the term SA to denote the sample allocation derived from the bandwise independent wavelet model. We should note here that neither SA nor HSA is optimal for turbo decoding. The problem with SA is that it tends to undersample the fine scale bands since they contain less energy than the coarse bands when treated independently and we are less confident on the activity rates. While HSA is served as the benchmark by assuming we have the accurate activity rate information for each wavelet coefficient. The optimal sample allocation for turbo decoding should combine the merits of both SA and HSA.

3.5 Simulations

Reconstruction performance for natural images with the band wise sampling matrices introduced in Section 3.3 and Section 3.4 is demonstrated and compared with several existing sensing matrices in this section. We start with the 256×256 cameraman image as an introductory example. With the knowledge of the image statistics, we show that the bandwise independent image model based SD function can accurately predict the reconstruction quality for the proposed sample allocation scheme. It also confirms the theoretical optimality of our band-wise sensing matrix. We then extend the scheme to practical compressive imaging by designing the general sample allocation with the average image statistics estimated from the training set of

the Berkeley dataset [3]. Simulation with ten images from the test set further confirms that with good statistical estimation, the proposed SD sample allocation exhibits state-of-the-art performance.

3.5.1 Sample Allocation with Oracle Image Statistics

The cameraman image is decomposed into six bands using the Daubechies 2 wavelet. GGD and GMD model parameters estimated directly from the wavelet coefficients are reported in Table 3.1 as the oracle image statistics, using moment matching [10] and the EM algorithm [110] respectively. Given the parameter estimation, we are able to generate the image SD function and the subsequent band-wise sample allocation using the aforementioned method.

subband		b_0	b_1	b_2	b_3	b_4	b_5
GGD	α	2	0.7	0.4	0.3	0.3	0.4
	σ^2	261.4383	2.0822	0.4559	0.0902	0.0167	0.0033
GMD	λ	1	0.4155	0.5309	0.4842	0.3664	0.2792
	σ_L^2	261.4383	4.4215	0.8542	0.1856	0.0453	0.0115
	σ_S^2		0.3331	0.0038	0.0004	0.0002	0.0001

Table 3.1: Statistics for Daubechies 2 wavelet coefficients of cameraman

To show the sample allocation method is not restricted to the form of the decoders, we consider three reconstruction options: the linear ℓ_2 decoder, the CS ℓ_1 decoder, and the BAMP decoder. The SPGL1 toolbox ³ is used to implement the ℓ_1 decoder. Its SD function can also be derived using the SE formalism [55]. Both the ℓ_2 and the ℓ_1 decoder are considered for the GGD and the GMD model. Although in [54] the authors show that the BAMP decoder is applicable to the GGD data by approximating it with the finite term of Gaussian mixture distribution, the approximation error may contribute to the final reconstruction distortion. Thus the BAMP decoder results are only reported for the GMD model here. The detailed algorithm can be found in [69], [53].

For quantitative comparison, the peak signal-to-noise ratio (PSNR) is used for both theoretical prediction and simulations. We examined the cameraman image at four different sampling ratios: 10%, 15.26%, 25% and 30% associated with $m = 6554, 10000, 16384, 19661$ noiseless measurements. Two different wavelet image models are considered. First, the wavelet bands are assumed to be mutually independent. The proposed SA matrix is compared with five sensing

³<http://www.cs.ubc.ca/labs/scl/spgl1/index.html>

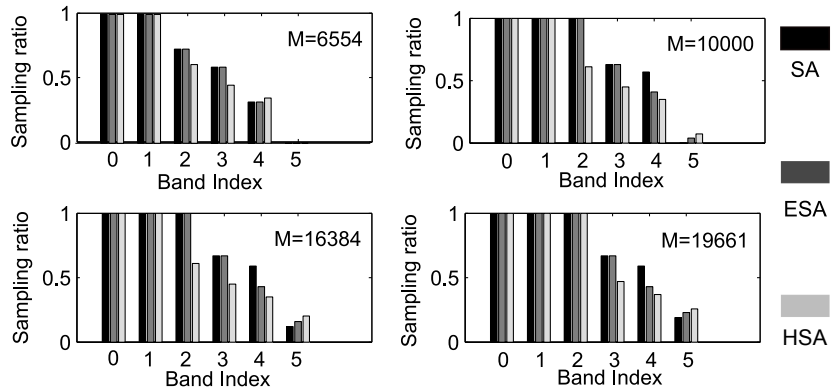


Figure 3.14: *Sample allocation per band for Daubechies 2 wavelet with the GMD model. SA: sample allocation based on the bandwise independent model. HSA: sample allocation based on the empirical SD functions for BAMP decoder with soft information. ESA: empirically optimized sample allocation for turbo decoding.*

matrices: the homogeneous Gaussian matrix (Uniform), the two-gender matrix (2 Gender) [84], the informative sensing matrix (InforSA) [87] and the multi-scale sensing matrix (MBSA) in [85]. The 2 Gender matrix is implemented as fully sampling the scaling band and uniformly allocating the remaining samples to all the wavelet bands. As a statistic-dependent sample allocation scheme, InforSA is also generated based on Table 3.1.

The corresponding PSNR results are shown in Fig. 3.15 and Fig. 3.16 for GGD and GMD model, respectively. The SD function predicts the expected distortion quite accurately for all three choices of the decoder with SA. For both image models, SA achieves the best performance among the five sensing matrices. The advantages of SA over the Uniform matrix and the 2 Gender matrix is significant in spite of the sample ratio. MBSA has a relatively good performance since it has the essence of putting more samples to the coarse bands. Provided with the same image statistics, InforSA tends to allocate more samples to the fine wavelet bands compared with SA. Thus it is not as effective as SA in the low sampling ratio regime. Interestingly the CS scheme, even with an optimized sample allocation, only provides modest reconstruction gains over the classical linear approximation with similarly optimized sample allocation. This suggests the discussion in Section 3.3.3.2: the rate of decay of error is the same for both the BAMP and ℓ_2 decoder (though the constants are different). Thus we do not observe overwhelmingly better performance for the BAMP decoder even when SA is performed.

Secondly, the quad-tree structure is exploited with the GMD model. Within the turbo decoding regime, simulations are reported for four different sensing matrices: Uniform (amounts to the algorithm proposed in [69]), SA, HSA, and the empirically optimized sample allocation, or

ESA. As analysed in section 3.4, ESA should be the balance between SA and HSA. For the cameraman image, the ESA is obtained by adaptively reallocating samples from band four to band five based on SA, with the step size of 100 samples, until the PSNR does not increase. The sample allocation per band under four specific sampling ratios are reported in Fig. 3.14. We see that the scaling band and the coarsest wavelet band always have priority over the fine wavelet bands. For this particular image, around 2000 samples are reallocated to the finest scale band to achieve the ESA. For the turbo decoding, the soft information in (3.37) is used. It is fixed if band j is fully sampled during the HMT decoding. For partially sampled bands, activity rates λ_j in Table 3.1 are used to initialize the turbo decoding and updated by the HMT decoding for each turbo iteration. Other hyperparameters to initialize the HMT decoding are set in accordance with the recommendation in [69]. For various choices of sample allocations, we ran 20 turbo iterations, within which 500 BAMP iterations are performed.

As evident in Fig. 3.16, adding the HMT decoding ingredient indeed improves the reconstruction quality. It is the joint use of optimized bandwise sampling and the tree structure that delivers by far the best PSNR performance. Again, sample allocation shows its importance when there is a tight budget of samples: even without the turbo decoding procedure, SA+BAMP is 1.5 dB better at $\gamma = 0.1, 0.15$ than Uniform+TurboAMP. In the large sampling ratio regime $\gamma = 0.3$, the effectiveness of the sample allocation is not as obvious and the HMT alone is responsible for the excellent performance: SA+TurboAMP is 0.5 dB better than the Uniform+TurboAMP. It shows that both sample allocation and the HMT play a role in improving the performance of compressive imaging, and which matters more depends on several factors, including the sampling ratio. We also observe that the ESA is only slightly better than the SA, which means that even when we have the luxury of manipulating samples, the benefit is limited because of the exponential energy decay of the multi-resolution model.

The 256×256 cameraman image along with the reconstructed images by different encoder-decoder pairs are visualized in Fig. 3.17 at the sampling ratio $\gamma = 15\%$. It further confirms that given accurate image statistics, our proposed SA is the optimal distribution of samples.

3.5.2 Sample Allocation with General Image Statistics: The GSA

In practice, we may not have access to the accurate image statistics. In this section, reconstruction results for a general sample allocation (GSA) which is not tuned to a specific image distribution are presented. The GSA is designed based on the fixed per-band natural image

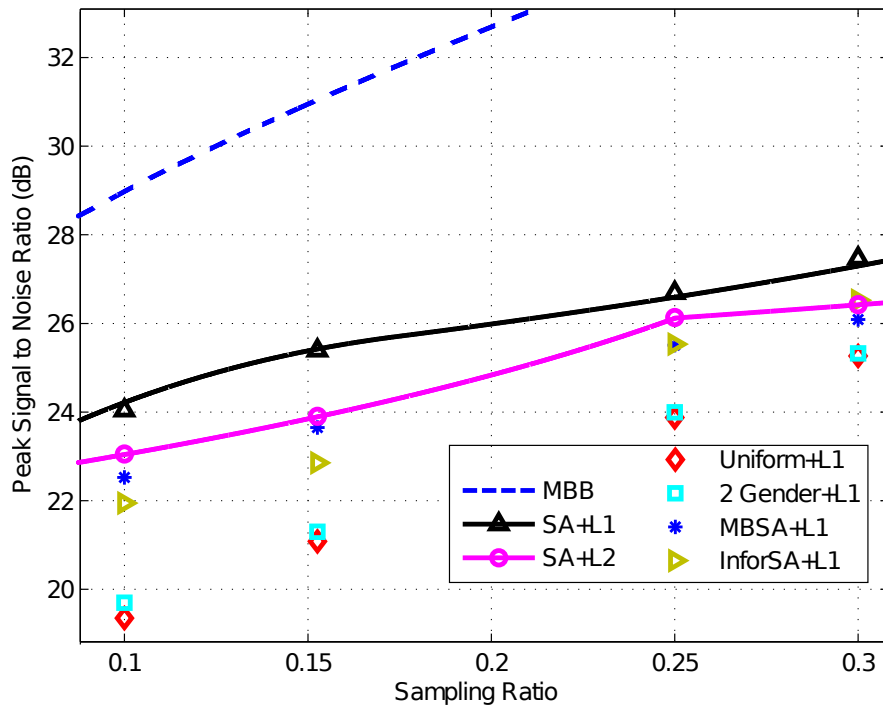


Figure 3.15: PSNR comparison of different encoder-decoder pairs for cameraman Daubechies 2 wavelet with the GGD model. The lines are theoretical predictions with the SD function. While dots represent simulations with the cameraman image.

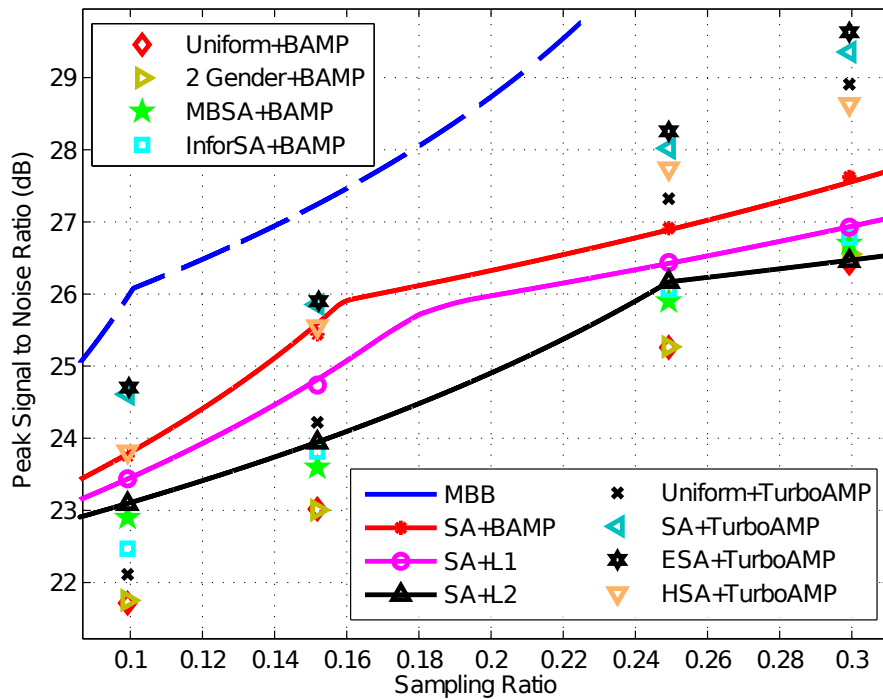


Figure 3.16: PSNR comparison of different encoder-decoder pairs for cameraman Daubechies 2 wavelet with the GMD model. The lines are theoretical predictions with the SD function. While dots represent simulations with the cameraman image.

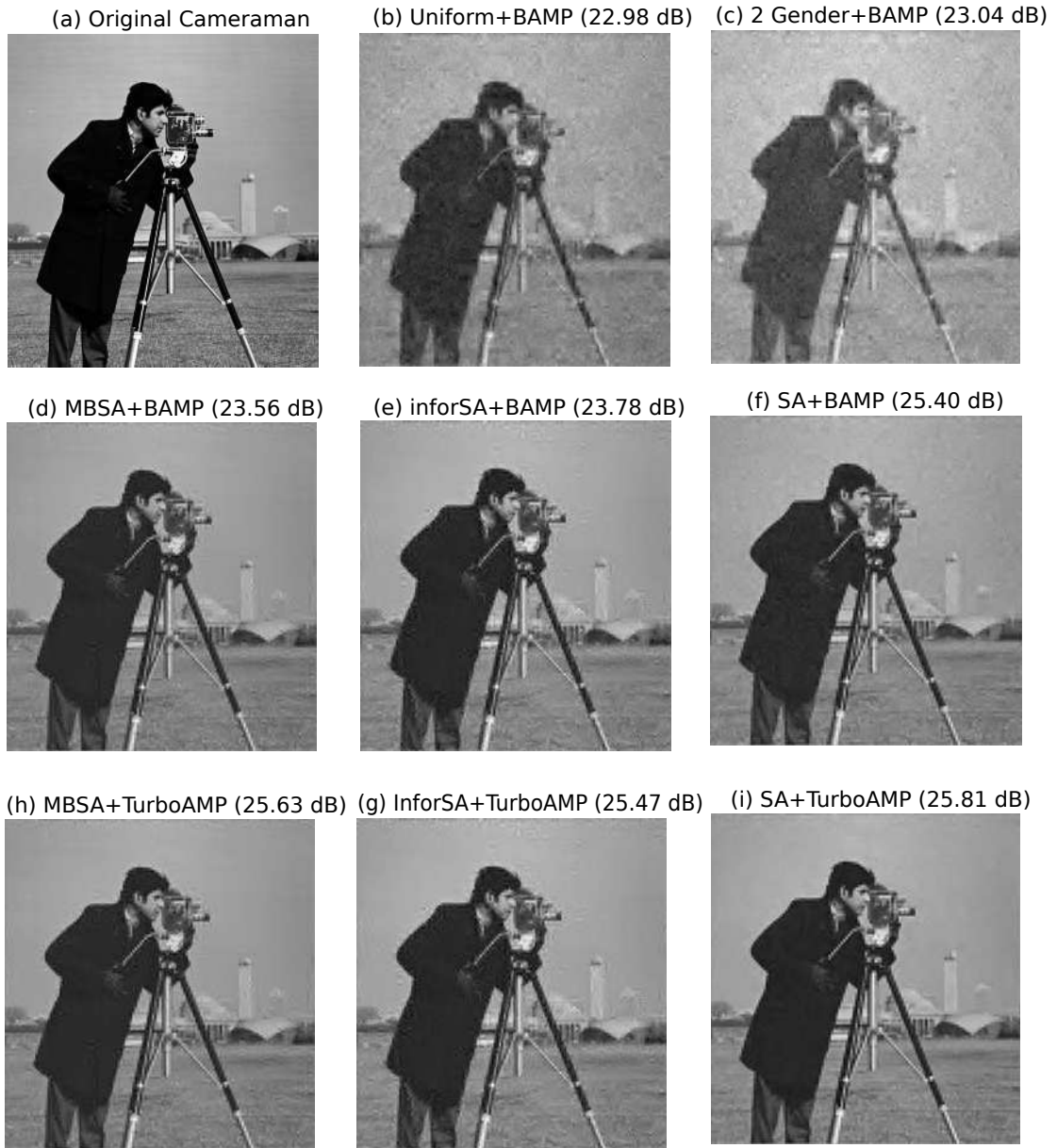


Figure 3.17: Reconstruction using 10000 (15%) samples of the 256×256 cameraman image with different encoder-decoder pairs. The GMD is used to model the Daubechies 2 wavelet coefficients statistics. The encoding matrices for the cameraman simulations are explained in details in Sec. 3.5.1.



Figure 3.18: Ten test images from the Berkeley dataset [3]. From left to right, top to bottom are: car, plane, eagle, sculpture, surfer, tourists, building, castle, man and fish.

statistics. We estimated the GMD statistics for the six-band Daubechies 2 wavelet decomposition of 200 training images from the *Berkeley Segmentation Dataset* [3]. Each training image is cropped to the size of 256×256 . The pixel intensity value is normalized between 0 and 1. The average per-band GMD parameters are reported in Table 3.2 and used to generate the general (albeit dictionary and algorithm dependent) sample allocation profile.

subband	b_1	b_2	b_3	b_4	b_5
λ	0.5108	0.4374	0.4076	0.3616	0.3137
σ_L^2	3.6910	0.7506	0.1595	0.0385	0.0081
σ_S^2	0.4596	0.0490	0.0075	0.0015	0.0003

Table 3.2: Average Statistics for Daubechies 2 wavelet coefficients of 200 test images from the Berkeley dataset [3]

The resulting GSA is then applied to ten test images outside the training set, as shown in Fig. 3.18, and again compared with the Uniform matrix, the 2 Gender matrix, MBSA and InforSA. Table 3.2 is also used to generate InforSA. The BAMP decoder is used as the reconstruction algorithm. The PSNR performance for sampling ratio $\gamma = 0.1, 0.2, 0.3$ are reported in Table 3.3, Table 3.4 and Table 3.5, respectively.

The reconstruction quality of GSA depends on the accuracy of the image statistics. We see that with reasonable image statistics estimation, GSA outperforms the Uniform matrix and the 2 Gender matrix with roughly 2 dB gain consistently for all cases. The MBSA and InforSA have comparable yet slightly worse performance except three images at sampling ratio $\gamma = 0.3$. It is due to the actual image deviates from the average image statistics. Not surprisingly, adding the HMT decoding component can only improve the reconstruction quality, if not significantly.

Image	GSA	InforSA	MBSA	Uniform	2 Gender	SA+TurboAMP
car	22.52	21.67	22.28	20.61	20.65	23.12
plane	25.87	25.27	25.63	24.16	24.26	26.57
eagle	25.23	24.53	24.88	23.39	23.44	26.30
sculpture	22.42	21.72	22.36	20.75	20.81	22.68
surfer	22.37	21.58	22.11	20.42	20.59	23.14
tourists	22.17	21.35	22.08	20.41	20.50	22.52
building	22.01	21.42	21.84	20.39	20.41	22.73
castle	21.40	20.93	21.26	19.82	19.78	21.74
man	26.86	26.02	26.42	24.84	24.89	28.52
fish	24.60	23.52	24.43	22.57	22.63	24.85
average	23.55	22.80	23.33	21.74	21.80	24.22

Table 3.3: Image reconstruction results for ten 256×256 test images from the berkeley image database [3] with $\gamma = 0.1$. Entries are the peak signal-to-noise ratio (PSNR) in decibels, $PSNR := 10 \log_{10}(N/\|\hat{x} - x\|_2^2)$. All results use the average image statistics reported in Table 3.2 and the BAMP decoder.

3.6 Summary

The main contribution of this chapter is to understand the nature of the sampling for multi-resolution images. For this, the complete sample distortion framework with the definition, lower bounds and the convex property is presented. Given the image statistics, we have derived a tractable sample allocation method for minimizing the reconstruction distortion and shown that it provides an accurate prediction of the achievable SD performance. We have also shown that when the optimized sample allocation is performed, the reconstruction gain of the CS decoder is limited over the linear reconstruction techniques. To get closer to the model based bound, we have deployed the tree structured sparsity within the optimized band-wise sampling framework by the turbo decoding approach. Various encoder-decoder combinations examined with the cameraman image illustrate the merit of band-wise sampling, especially in the regime of very low sampling ratios. For practical sample allocation, a general sampling profile is constructed based on average image statistics and demonstrates competitive performance.

Image	GSA	InforSA	MBSA	Uniform	2 Gender	SA+TurboAMP
car	25.56	24.11	25.29	22.92	22.98	25.92
plane	28.28	27.32	28.13	26.19	26.25	28.52
eagle	28.66	27.84	28.59	26.31	26.44	28.95
sculpture	23.81	22.89	23.54	22.05	22.61	24.58
surfer	25.37	24.00	25.13	22.81	22.95	25.65
tourists	24.15	22.93	23.75	22.08	22.37	24.53
building	24.84	23.59	24.66	22.48	22.55	25.37
castle	23.65	22.76	23.41	21.02	21.42	23.96
man	30.32	29.33	30.08	28.05	28.49	30.80
fish	27.26	27.57	26.76	24.62	24.83	27.76
average	26.10	25.23	25.93	23.85	24.09	26.60

Table 3.4: Reconstruction PSNR for test images with $\gamma = 0.2$

Image	GSA	InforSA	MBSA	Uniform	2 Gender	SA+TurboAMP
car	26.21	26.24	26.15	25.00	25.22	26.97
plane	28.96	29.20	28.89	28.21	28.54	29.82
eagle	29.97	29.22	29.17	28.61	28.94	30.25
sculpture	24.94	23.93	25.02	23.00	23.11	25.72
surfer	26.04	25.96	25.85	24.91	25.05	26.85
tourists	25.35	24.22	25.15	23.35	23.57	25.79
building	25.50	25.28	25.32	24.28	24.42	26.17
castle	24.32	24.21	24.16	23.04	23.06	24.75
man	31.56	30.85	30.77	30.05	30.29	33.09
fish	28.76	27.97	28.26	26.31	26.53	29.31
average	27.16	26.71	26.87	25.68	25.87	27.87

Table 3.5: Reconstruction PSNR for test images with $\gamma = 0.3$

Chapter 4

Modulated Matrix Design

In this Chapter, the modulated matrix design is proposed as an extension of the seeded matrix in the CS literature. The structure of the modulated matrix can be seen as the product of a homogeneous Gaussian matrix and a rescaling matrix. A 1-D SE equation is derived for the modulated matrix by modifying the block SE function for the seeded matrix. Thus, the corresponding SD performance can be accurately predicted. The relatively simple form of the modulated matrix potentially reduces the complexity of the parameter optimization procedure while retains the ability to perform as well as the seeded matrix. The two block matrix is then presented as an exemplary realization of the modulated matrix design. Interestingly, the zeroing matrix introduced in Chapter 3 falls into the two block matrix framework. Since the performance of the proposed measurement matrix depends on the first order phase transition (FOPT), we analyse this signal property using the SE equation. We have shown that for sparse and dense signals with a FOPT, exact reconstruction can be achieved in the region where the homogeneous Gaussian matrix is not optimal. For compressible signals without a FOPT, the two block matrix can effectively improve the SD function, with the zeroing matrix being the empirically optimal choice.

4.1 Introduction

One of the major focuses in compressed sensing is the optimal configuration for recovery, i.e. the optimal measurement matrix and reconstruction algorithm. In [104], an extensive study of the optimal CS reconstruction for sparse signals is reported with some rigorous proof. As we have shown already in Chapter 3, there is no such thing as the universally optimal measurement matrix. The optimality of the measurement matrix is highly related to both the signal prior and the recovery scheme. Despite the general advantages of the homogeneous Gaussian matrix, there have recently been a number of studies on tailoring the measurement matrix Φ with the signal distribution and the reconstruction algorithm, aiming for better CS performance. In Chapter 3, a hybrid zeroing matrix was introduced by exploring the convex property of the SD function, which successfully convexifies the SD function in the low sampling regime, thus

improving the reconstruction quality. Another measurement matrix design attracting a lot of attention has the spatially coupled structure. The spatial coupling concept was first developed and implemented in the coding theory [111–113]. Kudekar et al. first presented the effectiveness of the spatial coupling in CS in [111]. Krzakala and colleagues further promoted its usage in CS and denoted the corresponding matrix as the seeded matrix [4]. Designed as the spatially coupled block diagonal matrix, it has been shown heuristically that exact recovery of the sparse signal can be obtained under a sampling ratio approaching the sparsity level. Rigorous proof for its success is given in [114]. Asymptotic analysis and state evolution prediction for the block matrix structure are derived in [53] using the replica method. The application for compressible signals is considered in [115].

In this chapter, we propose a new block measurement matrix structure, by introducing a random rescale distribution to modify the homogeneous Gaussian matrix. The proposed matrix is denoted as the modulated matrix. Different from the block diagonal structure of the seeded matrix, the modulated matrix $\Phi_M \in \mathbb{R}^{m \times n}$ consists of several m -row Gaussian matrices with different variances. The variances admit a probability density function specified by the rescale distribution. By varying the variance for the sub-matrices, we are essentially reweighting the signal prior. Another key difference between the modulated matrix and the seeded matrix is the complexity of the SE analysis. For the seeded matrix, the dimension of the SE dynamical system is on the order of the number of the blocks. For the modulated matrix, we derive a 1-D SE equation to track the performance when used with the AMP based reconstruction algorithm in the large system limit. The 1-D equation makes the analysis and the optimization of the measurement matrix relatively easy. Inspired by the zeroing matrix in Chapter 3, we then consider a rescale distribution consisting of two Dirac delta functions. The rest of the chapter will focus on the special form of the modulated matrix, known as the two block matrix.

In [115], Barbier et al. pointed out the sub-optimal sampling region for BAMP with the homogeneous Gaussian matrix and the associated first order phase transition phenomenon for the signal. It is shown that with the presence of a FOPT, the seeded matrix is empirically able to improve the reconstruction and reach the optimal achievable MSE. As an important property, the FOPT is explained using the replica method in [115]. In this chapter, it is interpreted from the state evolution point of view. Three different types of SE behaviour are presented and the cause of the FOPT is analysed. We further derive the necessary and sufficient condition for signals without a FOPT.

The work in [4, 115] emphasizes on achieving the perfect reconstruction for sparse signals

with a sampling ratio that near the sparsity level with the seeded matrix. Here, we focus on improving the reconstruction quality with the full range of the sampling ratio. The performance of the modulated matrix is evaluated by the SD function. We show that in the SD framework, the FOPT is related to the position of the critical sampling ratio defined in the Chapter 3. For signals with a FOPT, the two block matrix is able to reduce the critical sampling ratio. In the case of the sparse and dense signal priors, this means the exact reconstruction is available in the region that the homogeneous Gaussian matrix fails. For compressible signals without a FOPT, theoretically we prove that the two block matrix will retain this non-FOPT property. Empirical results indicate applying the two block matrix will not change the position of the critical sampling ratio. Regardless of the FOPT, the two block matrix delivers a significantly improved SD performance. Finally, the theoretical and empirical SD functions for sparse, compressible and dense signals confirm the 1-D SE analysis and demonstrate the power of the two block matrix design.

4.2 Seeded Matrix Review

The ultimate principle of compressed sensing is acquiring only enough information necessary to restore the original signal. For a sparse signal with k non-zeros elements in n dimension $k < n$, essentially the knowledge of the $k + 1$ components is enough to represent the signal. In principle, exact recovery of the signal is possible with m measurements with $m > k$. In compressed sensing, ideally we would like to have a measurement matrix and a reconstruction algorithm that achieves exact recovery with a sampling rate reaching the optimal limit, i.e. $m/k \rightarrow 1$, in the large system limit.

In [4], the fundamental reconstruction limit is achieved in the limit of large systems for sparse signals. Krzakala et al. proposed the seeded belief propagation (s-BP) procedure, which is essentially the TAP-AMP summarized in Chapter 2 and presented the carefully designed measurement matrix, designated as the seeded matrix.

It was demonstrated both numerically and analytically that the seeded matrix together with s-BP is able to exactly reconstruct sparse signals with m very close to the sparsity level k . In [115], the compressible signal reconstruction with the seeded matrix is investigated. The authors pointed out that for compressible signals with a FOPT, the homogeneous Gaussian measurement matrix together with BAMP does not truly achieve the optimal Bayes inference in the small sampling regime. The important contribution of [115] is the explanation of the

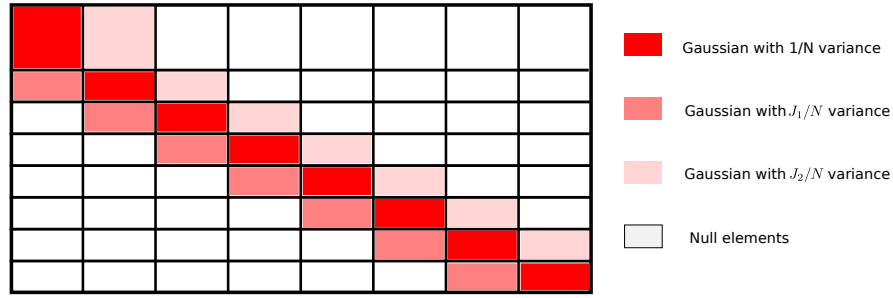


Figure 4.1: Construction of the seeded matrix for compressed sensing [4].

FOPT phenomenon and the demonstration of how the seeded matrix is able to aid the BAMP algorithm to achieve the optimal Bayes inference.

This section presents a brief summary of the seeded matrix, including its structure, heuristic explanation for its working principle and the theoretical state evolution dynamics, as the background information for the modulated matrix design.

4.2.1 Seeded Matrix Structure

The seeded matrix has a spatial coupling structure. The measurement matrix Φ is divided into $L_r \times L_c$ blocks with each being $\Phi_{qp} \in \mathbb{R}^{m_q \times n_p}$, $q = 1, \dots, L_r$, $p = 1, \dots, L_c$. Consequently, the signal of interest can be seen as divided into L_c equal-sized blocks. For non-zero blocks, the components for each block are drawn i.i.d. from the Gaussian distribution with zero mean and variance J_{qp}/n . The standard seeded matrix has the block diagonal structure, as illustrated in Fig. 4.1. The principle diagonal blocks have Gaussian elements with variance $1/n$. The coupling blocks sitting above and below have variance J_2/n and J_1/n , respectively. Empirical experiments suggest good reconstruction is obtained with large J_1 and small J_2 . The seeded matrix in Fig. 4.1 is not the only valid structure to achieve the reconstruction optimality. More designs can be found in [53], all of which have the general spatially coupling structure.

There is a heuristic explanation for the working principle of the seeded matrix. To approach the theoretical limit for the perfect reconstruction of sparse signals, the first block of the seeded matrix Φ_{11} is chosen to be near square shaped to achieve almost exact reconstruction for the first signal block. Then the reconstruction propagates through the coupling matrices as a wave for the following blocks, making good reconstruction possible even for blocks with very small sampling ratio $\gamma_q = m_q/n_p$. Reconstruction with the seeded matrix has an analogy to the crystal nucleation. For the supercooled liquid trapped in a glassy state, a large enough seed

of crystal will enable it to escape from the metastable state and let the crystal grow from its seed. For the seeded matrix, the first square shaped block acts as the nuclear seed to embark the perfect reconstruction.

4.2.2 When Does Seeding Work?

In crystallization, the reason that the nuclear seed is able to trigger the procedure is essentially because the liquid is in the unstable glassy state. In physics, the way liquid changes into solid is a typical example of a system undergoing a first order phase transition. For CS reconstruction, we borrow this physics term to describe the scenario when the signal reconstruction is trapped at a sub-optimal solution. Equivalently, the seeded matrix triggers better reconstruction for signals exhibiting an unstable state, or undergoing a first order phase transition. In the context of the sample distortion framework, the FOPT is a discontinuous drop of the MSE with the increasing sampling ratio. More precisely, for a fixed sparsity level, there exists a sampling ratio γ_{BP} that separates a phase with a small MSE, obtained at $\gamma > \gamma_{BP}$, from the phase with a large MSE for $\gamma < \gamma_{BP}$. The MSE discontinuity happens at $\gamma = \gamma_{BP}$.

It was argued heuristically in [4] and shown empirically in [115] that only if the FOPT is present, the optimal Bayes inference can be restored by the spatially coupling measurement matrix. For sparse signals with a FOPT, BAMP is able to obtain exact recovery for $\gamma > \gamma_{BP}$. Below γ_{BP} an unstable state appears. Instead of finding the original signal, the BAMP reconstruction gets stuck in a sub-optimal solution. In this case, the seeded matrix achieves the optimal Bayes inference in the sense that the exact reconstruction is obtained. For compressible signals with a FOPT, the optimal Bayes inference corresponds to a solution with a smaller MSE, while the BAMP with the homogeneous Gaussian measurement matrix terminates at one with a larger MSE. In this scenario, the seeded matrix improves the BAMP recovery by leading the algorithm to converge to the optimal Bayes inference result.

As an important factor effecting the signal reconstruction, we will provide more detailed explanation for the FOPT later in Section 4.4.

4.2.3 State Evolution for Seeded Matrix

As we have shown in Chapter 2, the behaviour of the AMP-based algorithm with the homogeneous Gaussian matrix can be characterized by the SE dynamical system in the large limit. For the s-BP algorithm with the joint use of seeded matrix, the same analysis can be applied to

track its MSE behaviour. In [4], the authors presented the detailed derivation of the SE equation for the seeded matrix from the replica method perspective. Here we summarize the conclusion and present the resulting $2L_c$ -D dynamical system:

$$E_p^{(t+1)} = \mathbb{E} \left\{ \left[F\left(x + \frac{z}{\sqrt{\tau_p}}; \frac{1}{\tau_p}\right) - x \right]^2 \right\} \quad (4.1)$$

$$\tau_p = \sum_{q=1}^{L_r} \frac{m_q J_{qp}}{\sum_{r=1}^{L_c} J_{qr} n_r E_r^t} \quad (4.2)$$

The SE prediction for s-BP with the seeded matrix has the same interpretation as for the homogeneous matrix: $z \sim \mathcal{N}(0, 1)$ is the Gaussian noise which is independent of x . The function $F(\cdot)$ is the non-linear denoising estimator of x given $x + z$. Given the signal prior, $F(\cdot)$ can be the MMSE estimator. The expectation in (4.1) is taken with respect to both x and z . $E_p^{(t+1)}$ represents the reconstruction MSE for the p th signal block. This dynamical system allows us to obtain the theoretical performance for the seeded matrix. Moreover, it can be used to optimize the matrix parameters, L_c , L_r , J_1 and J_2 for good reconstruction quality. However, the optimization may not be trivial as the number of blocks grows.

4.3 Modulated Matrix Framework

As stated in [53], the seeded matrix is not the only choice for improving the CS reconstruction. In this section, we introduce the *modulated matrix* design, as a general measurement matrix framework, which possesses a much simpler SE dynamics. The two block matrix is then presented as a special realization of the modulated matrix design.

4.3.1 Modulated Matrix Structure

Instead of dividing both columns and rows of the measurement matrix into blocks, the modulated matrix Φ_M is composed of L_c m -row sub-matrices $\Phi_i \in \mathbb{R}^{m \times n_i}$, $i = 1, \dots, L_c$, and $\sum_i n_i = n$. Each consists of i.i.d. random elements drawn from the Gaussian distribution with zero mean and J_i/n variance.

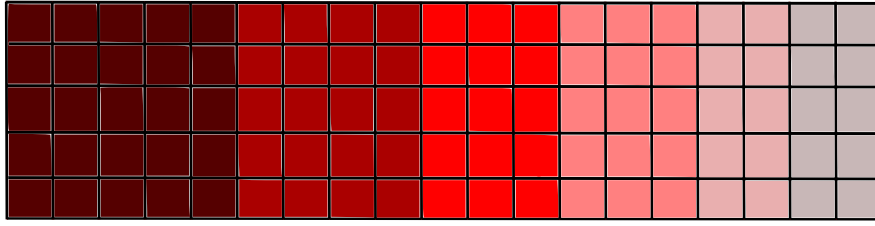


Figure 4.2: Construction of the modulated matrix. Gaussian random elements with different variances are indicated by different shade.

Let us define the rescaling matrix $\mathbf{R} \in \mathbb{R}^{n \times n}$ as:

$$\mathbf{R} = \begin{pmatrix} \sqrt{J_1} \mathbf{I}_1 & 0 & \cdots & 0 \\ 0 & \sqrt{J_2} \mathbf{I}_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{J_{L_c}} \mathbf{I}_{L_c} \end{pmatrix} \quad (4.3)$$

where $\mathbf{I}_i \in \mathbb{R}^{n_i \times n_i}$ is the identity matrix. The modulated matrix is then the product of the homogeneous Gaussian matrix \mathbf{G} and the rescaling matrix:

$$\Phi_M = \mathbf{G}\mathbf{R} \quad (4.4)$$

A plot illustrating the modulated matrix structure is shown in Fig. 4.2. As opposed to the standard seeded matrix, the modulated matrix has a dense structure. The different shade in Fig. 4.2 corresponds to difference rescaling parameter J_i sorted in a decreasing order.

4.3.2 1-D State Evolution

The state evolution equations in (4.1) and (4.2) are not exclusive for the seeded matrix. They actually provide a general formulation to track the MSE evolution for any block structured measurement matrix. Thus they can also be used to derive the SE dynamics for the modulated matrix as a special case. To be specific, for the modulated matrix we set $L_r = 1$ in eq. (4.2).

Then for each signal block, we have:

$$E_i^t = S(\tau_i^t) \quad (4.5)$$

$$= \mathbb{E} \left\{ \left[F(x + z\sqrt{\tau_i^t}; \tau_i^t) - x \right]^2 \right\} \quad (4.6)$$

$$\tau_i^{(t+1)} = \frac{\sum_k J_k \varepsilon_k S(\tau_k^t)}{\gamma J_i} \quad (4.7)$$

where $\varepsilon_k = n_k/n$. The total reconstruction MSE at each iteration is the average over all blocks:

$$\bar{E}^t = \frac{1}{L_c} \sum_{i=1}^{L_c} E_i^t \quad (4.8)$$

Eq. (4.5) to (4.7) are a straightforward implementation of the general SE dynamics for the block sensing matrix. The distinctive feature of the modulated matrix though, is that its SE dynamical system has a 1-D form. To show this we define a rescaled variable $\hat{\tau} = J_i \tau_i$, which is independent of the block index i . Then the update rule for $\hat{\tau}$ becomes:

$$\hat{\tau}^{(t+1)} = \frac{\sum_k J_k \varepsilon_k S(\hat{\tau}^t / J_k)}{\gamma} \quad (4.9)$$

Unlike the update of τ_p in (4.2), the evolution of $\hat{\tau}$ involves only its previous state and thus forms a 1-D SE equation. When the iteration of $\hat{\tau}$ converges to $\hat{\tau}^*$, the CS reconstruction MSE can be accurately predicted by

$$\bar{E} = \frac{1}{L_c} \sum_k S\left(\frac{\hat{\tau}^*}{J_k}\right) \quad (4.10)$$

We can also extend the aforementioned modulated matrix design to the stochastic setting by introducing a random rescaling parameter J for each column. That is, set $L_c = n$ and J with the distribution $p(J)$. In the limit of large systems, the SE equation (4.9) and the distortion prediction (4.10) become

$$\hat{\tau}^{(t+1)} = \frac{1}{\gamma} \mathbb{E}_J \left\{ JS \left(\frac{\hat{\tau}^t}{J} \right) \right\} \quad (4.11)$$

$$\bar{E} = \mathbb{E}_J \left\{ S \left(\frac{\hat{\tau}^*}{J} \right) \right\} \quad (4.12)$$

where the expectation is calculated with respect to J .

Both the deterministic (4.9) and stochastic (4.11) dynamics are described by a 1-D SE equation.

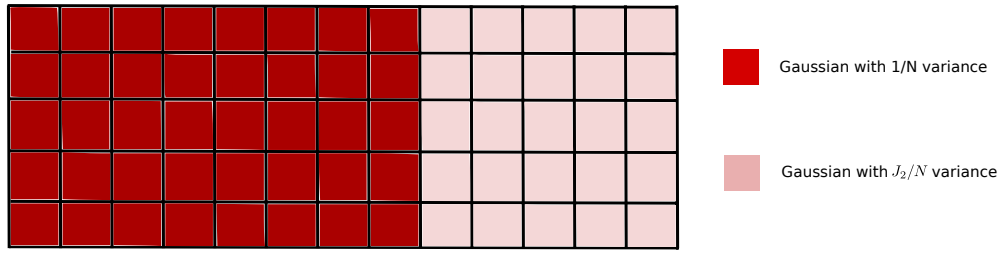


Figure 4.3: Construction of the two block matrix.

It makes the analysis and the optimization of the modulated matrix easier than the general seeded matrices of [4].

4.3.3 Two Block Matrix

In Chapter 3, we have proved the convexity of the SD function and showed the hybrid zeroing matrix can effectively convexify the concave SD function. It has been illustrated analytically that better performance can be achieved in the concave region by setting a portion of the measurement matrix to zero. Motivated by this design, we consider a special form of the rescaling matrix

$$\hat{\mathbf{R}} = \begin{pmatrix} \mathbf{I}_{n_1} & \mathbf{0} \\ \mathbf{0} & \sqrt{J_2} \mathbf{I}_{n_2} \end{pmatrix} \quad (4.13)$$

where $n_1 + n_2 = n$ and $\mathbf{I}_{n_i} \in \mathbb{R}^{n_i \times n_i}$. We denote the corresponding $\hat{\Phi}_M$ as the *two block matrix*. Fig. 4.3 illustrates the structure of the two block matrix.

There is a strong link between the two block matrix and the hybrid zeroing matrix: Setting $J_2 = 0$ and $\gamma_1 = \gamma/\gamma_c$ with γ_c being the critical sampling ratio results in the hybrid zeroing matrix and the convexified SD function. Here, we consider J_2 being non-zero and without loss of generality assume $0 < J_2 < 1$. The SE equation and the MSE prediction for $\hat{\Phi}_M$ become:

$$\hat{\tau}^{(t+1)} = \frac{1}{\alpha} M(\hat{\tau}^{(t+1)}) \quad (4.14)$$

$$= \frac{1}{\alpha} \left[\gamma_1 S(\hat{\tau}^t) + (1 - \gamma_1) J_2 S\left(\frac{\hat{\tau}^t}{J_2}\right) \right] \quad (4.15)$$

$$\bar{E} = \gamma_1 S(\hat{\tau}^*) + (1 - \gamma_1) S\left(\frac{\hat{\tau}^*}{J_2}\right) \quad (4.16)$$

The two block matrix design is closely related to the seeded matrix with four sub-matrices.

According to [4], the four block seeded matrix Φ_s takes the form:

$$\Phi_s = \begin{pmatrix} \mathbf{G}_1 & \sqrt{J_2}\mathbf{G}_2 \\ \sqrt{J_1}\mathbf{G}_3 & \mathbf{G}_4 \end{pmatrix} \quad (4.17)$$

where \mathbf{G}_i is the homogeneous Gaussian matrix. For the seeded matrix to work it requires $J_1 \gg J_2$. If we set $J_1 = 1/J_2$, which is close to what was found to be optimal in [4], the two block matrix $\hat{\Phi}_M$ turns out to be the rescaled seeded matrix.

$$\hat{\Phi}_M = \begin{pmatrix} \mathbf{G}_1 & \sqrt{J_2}\mathbf{G}_2 \\ \mathbf{G}_3 & \sqrt{J_2}\mathbf{G}_4 \end{pmatrix} \quad (4.18)$$

$$= \begin{pmatrix} \mathbf{I}_1 & 0 \\ 0 & \sqrt{J_2}\mathbf{I}_4 \end{pmatrix} \Phi_s \quad (4.19)$$

where the index matrix \mathbf{I}_i has the same number of rows as \mathbf{G}_i .

The intuitive idea behind the two block matrix design is that it simply shrinks a fraction of the signal to be very small. This leaves fewer large coefficients which can consequently be recovered through the reconstruction algorithm. In the SE dynamics, when the uncertainty for large coefficients is small enough, it acts as noise for the rescaled signal so that the two parts denoise together. Compared to the seeded matrix, the two block matrix has a relatively simple SE dynamics and fewer parameters to be selected, which makes analytical optimization possible. The potential downside maybe a reduced robustness to noise.

4.4 First Order Phase Transition

As previously mentioned, the FOPT is a crucial phenomenon indicating the possible improvement for the spatially coupling matrices. The cause of FOPT is first explained in [115] using the statistical physics tool. In the section, we will first summarize their analysis. Then we provide our own explanation for the FOPT phenomenon from the state evolution perspective and derive the necessary and sufficient condition for signals without a FOPT. Finally we analyse how the two block matrix effects the FOPT thus the reconstruction dynamics.

4.4.1 Analysis via Statistical Physics

To study the typical performance of CS in the large system limit, many works are published using statistical physics methods, whose principle goal is to study the macroscopic properties of physical systems from the principle of microscopic interactions. As randomness plays a key role in CS, it falls under the general scope of statistical physics. Although non-rigorous, empirical works in many research fields, i.e. compressed sensing [78–81], multi-user detection [47, 77], have shown promising results with the statistical physics analysis.

To explain the FOPT for BAMP reconstruction, the potential function (also known as the free energy function) is leveraged to characterize the CS system. The potential function is a statistical physics concept, which is originally used to characterize the thermodynamic properties of a disordered system. It has been shown in [53] that the fixed points of the message passing of a CS system are the stationary points of the corresponding potential function. In [4], the authors provide the detailed procedure to derive the potential function for a given sparse/compressible signal prior. Here we use the potential function for the two-state GMD in [115] to explain the cause of the FOPT.

Given the probability

$$p(x) = w_1 \mathcal{N}(x; 0, \sigma_1^2) + w_2 \mathcal{N}(x; 0, \sigma_2^2) \quad (4.20)$$

where $w_1 + w_2 = 1$, the corresponding potential function for BAMP with the homogeneous Gaussian matrix is formed as [115]:

$$\Lambda(E) = -\frac{\gamma}{2} \left(\log E + \frac{w_1 \sigma_1^2 + w_2 \sigma_2^2}{E} \right) + \sum_{a=1}^2 w_a \int \frac{e^{-z^2/2}}{\sqrt{2\pi}} \log \left[\sum_{b=1}^2 w_b \frac{e^{\frac{(t^2 \sigma_a^2 + t) z^2}{2(t+1/\sigma_b^2)}}}{\sqrt{t\sigma_b^2 + 1}} \right] dz \quad (4.21)$$

where $t = \frac{\gamma}{E}$ and E is the reconstruction MSE. The BAMP iterations correspond to the gradient ascent of $\Lambda(E)$. Initialized with $\hat{\mathbf{x}}^0 = 0$, BAMP obtains a better signal estimate thus a smaller E at each iteration. Consequently, $\Lambda(E)$ is in general increasing as the BAMP proceeds. When the BAMP terminates, the potential function $\Lambda(E)$ arrives at one of its maximas.

In Fig. 4.4, we plot potential functions associated with two GMD priors and one BG model under various sampling ratios. The potential function $\Lambda(E)$ for BG priors can be obtained by setting $\sigma_2^2 = 0$ in (4.21). For the GMD without a FOPT in Fig. 4.4.(c), all potential functions have only one global maximum irrespective of γ . Moreover, the value of the global maximum evolves smoothly, in the sense that it is a continuous function of γ . In Fig. 4.4.(c), we only

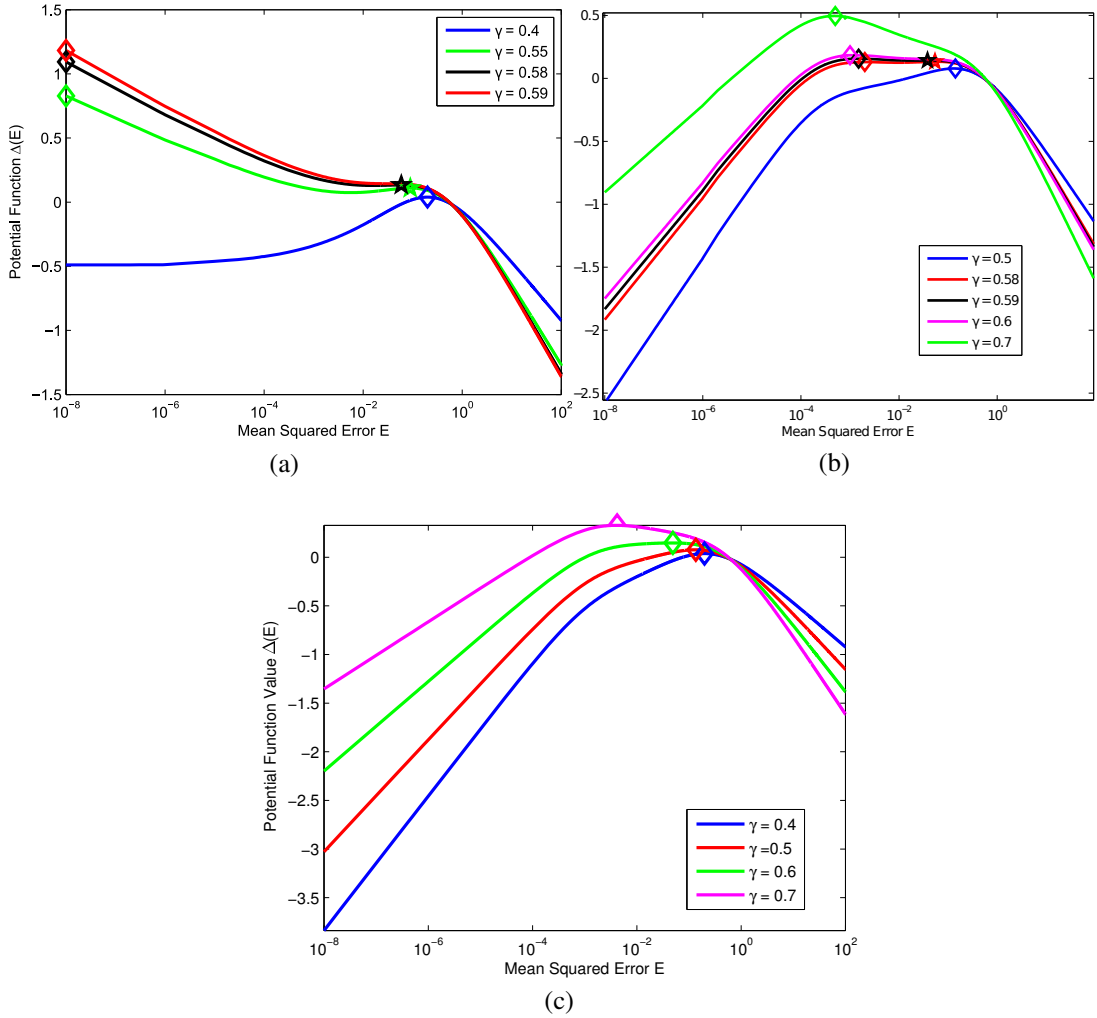


Figure 4.4: The potential function $\Lambda(E)$ for different signals at various sampling ratios with the homogeneous Gaussian matrix: (a) Bernoulli-Gaussian data with FOPT, $p_{BG}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\delta(x)$; (b) Gaussian-mixture data with FOPT, $p_{GM_1} = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 5e - 4)$ (c) Gaussian-mixture data without FOPT, $p_{GM_2} = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 0.003)$. The diamond-shaped dots represent the global maximums, while the star-shaped dots are the secondary local maximums.

plotted the potential functions for $\gamma = 0.4, 0.5, 0.6, 0.7$. One can imagine that if we apply a finer sampling ratio grid, connecting all the global maximum dots will form a continuous curve. As we increase the sampling ratio $\gamma \rightarrow 1$, the global maximum will appear at $E \rightarrow 0$, achieving the exact recovery. Since the iteration of BAMP acts like the steepest ascent of the potential function, we are expecting a continuous SD function, thus a SD function without FOPT for this particular prior. Later, the simulation results in Fig. 4.9 confirm this analysis.

Things are different for priors that induce a FOPT for BAMP. In Fig. 4.4.(a) we plot the evolution of the potential functions for the BG prior. In the small sampling ratio regime, i.e. $\gamma = 0.4$, there is only one maxima for $\Lambda(E)$ and the MSE of the BAMP recovery is the corresponding $E = 0.2007$. As we increase γ , a secondary local maxima shows its presence with the global maxima remaining at $E \rightarrow 0$. The reconstruction quality will then highly depend on the initialization of BAMP. If we could somehow initialize the algorithm from any point beyond the local maxima, in principle, BAMP is still able to reach the global maxima of $\Lambda(E)$, thus the exact recovery. However, such initialization barely happens in practice. With $\mathbf{x}^0 = 0$, one can accurately estimate the starting variance E^0 , which is normally larger than the one associated with the local maxima for $\Lambda(E)$. In such a scenario, BAMP with the homogeneous Gaussian measurement matrix will eventually get stuck at the local maximum instead of finding the optimal Bayes inference with the global maximum. The FOPT occurs when this spurious local maximum starts to vanish as we keep increasing the sampling ratio. In Fig. 4.4.(b), BAMP converges at $E = 0.0546$ for $\gamma = 0.58$. For $\gamma = 0.59$, the disappearance of the local maximum for the potential function leads to the significant change of MSE. Instead of having a MSE in the vicinity of $E = 0.0546$, BAMP converges to the global and the only maxima at $E \rightarrow 0$. This is how the sudden drop of the MSE, or the FOPT, happens in the SD function, as we will see later in Fig. 4.7.

The same phenomenon can be observed for GMD priors when the small Gaussian variance is neglectable, for example in Fig. 4.4.(b). As γ increases from 0.59 to 0.6, the local maximum $\sim 10^{-2}$ vanishes. In the SD function, we are expecting to see a sudden drop of the MSE from $\sim 10^{-2}$ to $\sim 10^{-4}$ in the sampling region between 0.59 and 0.6. In [115], the authors define γ_{BP} as the sampling ratio at which the MSE discontinuity happens (the largest γ that the potential function has two maximas) and γ_{opt} as the sampling ratio for which the two maximas has the same height. According to the previous analysis, it is in the region $\gamma_{\text{opt}} < \gamma < \gamma_{\text{BP}}$ that BAMP with a homogeneous Gaussian measurement matrix is sub-optimal, in the sense that it does not reach the global maximum. Also this region is where the spatial coupling measurement matrix

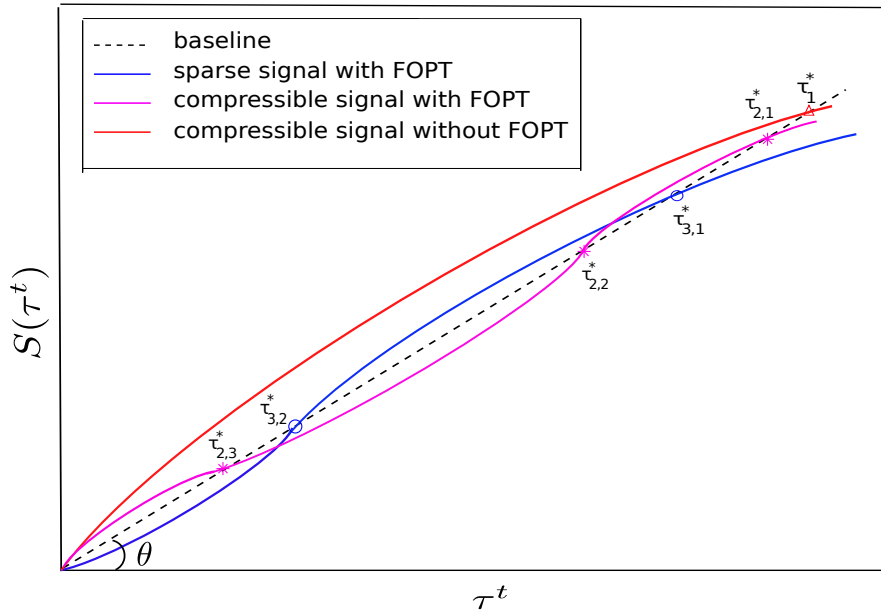


Figure 4.5: The schematic plot of three types of SE behaviour to explain FOPT. The dash line is the baseline $\gamma\tau^t$. The solid lines are $S(\tau^t)$ for BAMP with the homogeneous Gaussian matrix. The number of non-zero intersection points with the baseline varies for different types of signals.

could improve the recovery quality. In other words, the FOPT needs to be present for the seeded matrix to restore optimality.

4.4.2 Analysis via State Evolution

Although the analysis using the statistical physics tool coincides with the asymptotic behaviour of the AMP reconstruction, there is no rigorous proof for the connection. In contrast, SE dynamics is theoretically valid for characterizing the AMP behaviour. In this section, we study the FOPT phenomenon by analysing the SE equation for BAMP with the homogeneous Gaussian measurement matrix.

To better illustrate the dynamics, a schematic plot for three typical types of SE behaviour is presented in Fig. 4.5, which corresponds to the three types of the potential functions illustrated in Fig. 4.4. As summarized in Chapter 2 eq. (2.70), the SE equation for BAMP with a homogeneous Gaussian measurement matrix has the 1-D form:

$$\begin{aligned}\tau^{t+1} &= \frac{1}{\gamma} S(\tau^t) \\ &= \frac{1}{\gamma} \mathbb{E}\{[F(x + \sqrt{\tau^t}z; \tau^t) - x]^2\}\end{aligned}\tag{4.22}$$

Let us define τ^* as the convergence point of (4.22). In the SE dynamics, the intersection points of the baseline function $\gamma\tau^t$ and the function $S(\tau^t)$ are the fix points of the SE equation. They are also the possible values for τ^* . For any arbitrary prior, zero is always a fix point representing the exact recovery. As previously stated, whether BAMP can achieve the perfect reconstruction depends on the signal prior as well as the algorithm initialization. The curvature of $S(\tau^t)$ is determined by the signal prior. Changing the sampling ratio γ varies the baseline angle θ in Fig. 4.5, and the position of the fixed points.

As demonstrated in Fig. 4.5, for the compressible signal without a FOPT, apart from zero, there is only one non-zero fix point τ_1^* . Since the BAMP iteration starts with a large τ^0 , it will always converge to τ_1^* with $\tau^* = \tau_1^*$. As we gradually increase γ , the non-zero fix point decreases continuously to zero with γ approaching 1. It therefore leads to a smooth transition of τ^* with respect to γ , thus a continuous SD function. It corresponds to the potential functions with only one maximum in Fig. 4.4.(c).

For compressible signals with a FOPT, $S(\tau^t)$ consists of three smooth arcs. For small γ , the baseline intersects with $S(\tau^t)$ at three non-zero fix points. Since BAMP initializes with a large τ^0 , the BAMP iteration always terminates at the largest fix point $\tau_{2,1}^*$ associated with the first concave arc. As we gradually increase γ , $\tau_{2,1}^*$ and $\tau_{2,2}^*$ will move closer to each other and merge as one eventually. The FOPT happens when we keep increasing γ beyond this point. The baseline will surpass the first (concave) and the second (convex) curvature, resulting in only one non-zero fix point at $\tau_{2,3}^*$. Because of the existence of the convex curve between the two concave arcs, we cannot obtain a continuous transition of τ^* between $\tau_{2,3}^*$ and the merged $\tau_{2,1}^*/\tau_{2,2}^*$. The sudden vanishing of $\tau_{2,1}^*$ is the FOPT. We observe the same dynamical change in the potential function analysis in Fig. 4.4.(b).

When using the homogeneous Gaussian encoder and the BAMP decoder, sparse signals may also belong to the category of signals with a FOPT. However, their SE function behaves slightly different. As illustrated in Fig. 4.5, its $S(\tau^t)$ consists of a convex and a concave arc. The evolution of the convergence point is similar to the one for a compressible signal with a FOPT. For small γ , the baseline intersects with the concave curve of $S(\tau^t)$ at fix points $\tau_{3,1}^*$, $\tau_{3,2}^*$ and AMP converges at $\tau_{3,1}^*$. Once γ is large enough for the two points to merge, the convergence point τ^* will suddenly drop to zero as γ keeps increasing. Thus a discontinuity of the MSE to zero is expected in the SD function. For sparse signals, the local gradient of their SE dynamics at zero is such that the local stability is lost only when $m \leq k$. The corresponding type of potential function is shown in Fig. 4.4.(a).

4.4.3 FOPT Condition

Based on the FOPT analysis from the state evolution perspective, the baseline $\gamma\tau^t$ must always lie below $S(\tau^t)$ for any $\tau^t \in [0, \infty)$ if we want a smooth transition for the convergence point τ^* . Mathematically speaking, the slope of the baseline must be less than the gradient of the SE function for any $\tau^* > 0$. Here we present the necessary and sufficient condition for signals without a FOPT.

$$\frac{f(\tau^*)}{\tau^*} < \eta(\tau^*) \text{ for all } \tau^* > 0, \quad \text{and} \quad \eta(\tau) = \frac{df(\tau)}{d\tau} \quad (4.23)$$

where $\gamma\tau^{t+1} = f(\tau^t)$ is the general form of the SE equation.

4.4.4 Two Block Matrix Effect on FOPT

The way that the seeded matrix enables BAMP achieving the optimal Bayes inference, looking from the SE perspective, is that it changes the SE dynamics so that the sub-optimal fixed point vanishes. It is interesting to ask how the proposed two block matrix alters the SE dynamics. To better illustrate the two block matrix effect on the FOPT, $S(\tau^t) - \gamma\tau^t$ is plotted against τ^t for three types of SE behaviours in Fig. 4.6 with both homogeneous Gaussian and two block measurement matrix. The fixed points for SE iterations are the ones with $S(\tau^t) = \gamma\tau^t$.

For sparse signals with a FOPT, the two block matrix is capable of fundamentally altering the shape of the SE function to remove the spurious fixed points. As shown in Fig. 4.6.(b), with proper choice of the rescaling parameters, both non-zero fix points are eliminated so that the exact recovery is achievable at $\gamma = 0.58$ for the BG prior. This improvement can also be seen in the simulation in Fig. 4.7.

For the compressible prior with a FOPT as in Fig. 4.6.(c), the two block matrix is also capable of changing the structure of the SE function and accelerating the FOPT. For the homogeneous Gaussian measurement matrix with $\gamma = 0.58$, there are three non-zero fixed points and BAMP terminates at $\tau_{2,1}^*$, associated with some relatively large MSE. When the two block matrix is applied, there is only one non-zero fixed point left, associated with some very small MSE. It is as if the baseline angle θ is increased so that the intersection with the SE curve happens only at $\tau_{2,3}^*$ in Fig. 4.5. The actual SD function for is signal prior is shown in Fig. 4.8.

Finally, for signals which have no FOPT with the homogeneous Gaussian matrix, the dynamics

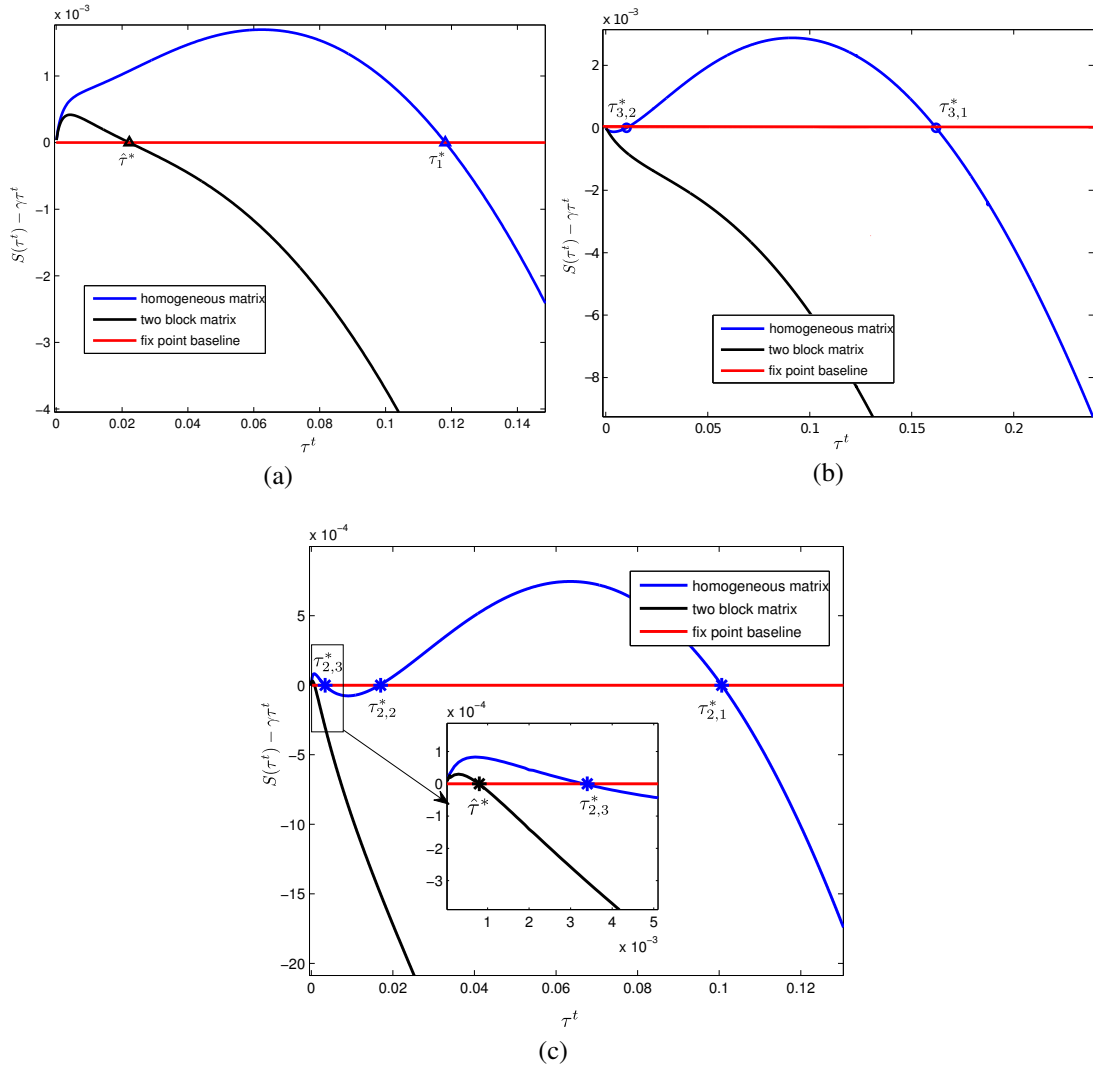


Figure 4.6: Fixed points of the SE evolution for both homogeneous Gaussian matrix and the two block matrix. (a) The compressible prior $p_{GM_2}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 3 \times 10^{-3})$ with $\gamma = 0.58$. For the homogeneous Gaussian matrix, the SE function has only one non-zero fixed point at $\tau_1^* = 0.1181$. Applying the two block matrix $J_2 = 1e - 3$, $\gamma_1 = 0.9206$ will not alter the shape of the SE function, it only shrinks the function so that the fix point is moved to $\hat{\tau}^* = 0.0222$. (b) The sparse prior $p_{BG}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\delta(x)$ with $\gamma = 0.55$. For the homogeneous Gaussian matrix, the SE function has two non-zero fixed points at $\tau_{3,1}^* = 0.1619$ and $\tau_{3,2}^* = 0.01020$. With the two block matrix $\gamma_1 = 0.847$, $J_2 = 10^{-3}$, the SE evolution successfully removes the spurious fixed points and leads to perfect reconstruction. (c) The compressible signal is $p_{GM_1}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 5 \times 10^{-4})$ with $\gamma = 0.58$. For the homogeneous Gaussian matrix, the SE equation has three non-zero fixed points at $\tau_{2,1}^* = 0.1006$, $\tau_{2,2}^* = 0.017$ and $\tau_{2,3}^* = 3.4 \times 10^{-3}$, With the two block matrix $\gamma_1 = 0.847$, $J_2 = 10^{-3}$, the fix point is moved to $\hat{\tau}^* = 0.0008$.

of the two block matrix keeps this property. This can be seen in Fig. 4.6.(a). Although the two block matrix successfully reduced the value of the convergence point of AMP, the general shape of the SE equation is unaltered. The following theorem confirms such observation.

Theorem 4. *If the SE equation for signals with the homogeneous Gaussian matrix $S(\tau)$ satisfies the no FOPT condition, then the SE equation for using the two block matrix $M(\tau)$ also satisfies the no FOPT condition.*

Proof. To prove the signal does not have FOPT with the two block matrix, we only need to check the gradient of $M(\tau)$.

$$\kappa(\tau) = \frac{dM(\tau)}{d\tau} \tag{4.24}$$

$$= \gamma_1 \eta(\tau) + (1 - \gamma_1) \eta\left(\frac{\tau}{J_2}\right) \tag{4.25}$$

$$< \gamma_1 \frac{S(\tau)}{\tau} + (1 - \gamma_1) \frac{J_2}{\tau} S\left(\frac{\tau}{J_2}\right) \tag{4.26}$$

$$= \frac{M(\tau)}{\tau} \tag{4.27}$$

where $\eta(\tau) = \frac{dS(\tau)}{d\tau}$ and the inequality is based on the no FOPT condition for the homogeneous matrix. □

Theorem 4 is true for any modulated matrix. Later in Section 4.5.2 we will also see that as the two block matrix delivers an improved SD performance, the critical sampling ratio for the SD function as defined in Chapter 3 remains the same as the one for the homogeneous Gaussian matrix. .

4.5 Simulations

In this section, we investigate the SD performance using the two block matrix for three different types of signals: BG, two state GMD and the k -dense prior. We also demonstrate the SD function of the homogeneous Gaussian measurement matrix for three priors as the performance benchmark. The seeded matrix performance is only demonstrated for the BG data since there are many parameters involved and the optimal configuration is only suggested for BG data in [4]. Throughout, we assume a noiseless scenario for all simulations. The theoretical SD functions with the two block matrix are calculated according to (4.14), (4.15) and (4.16). The empirical curves are obtained through Monte Carlo simulations with signals of length $N =$

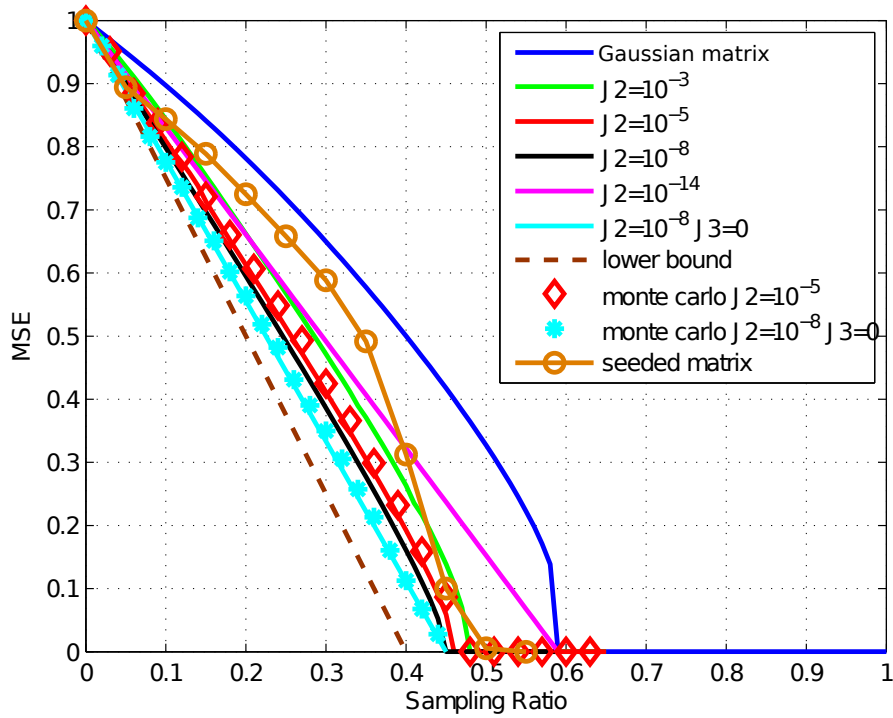


Figure 4.7: The normalized SD function for the sparse signal $p_{BG}(x)$ with different measurement matrix configuration. For two-block matrix, $\gamma_1 = \gamma/\gamma_c$ for $\gamma < \gamma_c$, where $\gamma_{c1} = 0.59$ is the perfect reconstruction ratio for the homogeneous Gaussian matrix. The three-block matrix is achieved by convexify the SD function of the two block matrix with $\gamma_1 = \frac{\gamma}{\gamma_{c1}}$, $\gamma_2 = \frac{\gamma}{\gamma_{c2}} - \frac{\gamma}{\gamma_{c1}}$, $\gamma_3 = 1 - \gamma_2 - \gamma_3$, where $\gamma_{c2} = 0.45$ is the perfect reconstruction ratio achieved by the two block matrix.

5000. The distortion performance for each sampling ratio is an average over 100 problem realizations. For reconstruction, the TAP-BAMP algorithm summarized in Chapter 2 is used, which does not assume the normalized columns for the measurement matrix. It is realized with the MATLAB toolbox provided in [70].

4.5.1 Two Block Matrix for Sparse Signal

First, we demonstrate the results for the sparse signal generated from the BG prior.

$$p_{BG}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\delta(x) \quad (4.28)$$

In Fig. 4.7, we plot the average MSE against the sampling ratio, under various choices of the rescaling parameter, J_2 . With the two block matrix we can reduce the sampling ratio for perfect reconstruction by decreasing J_2 : the perfect reconstruction ratio is moved from 0.59 to 0.45 with $J_2 = 10^{-8}$. However, further shrinking of J_2 does not improve the reconstruction to

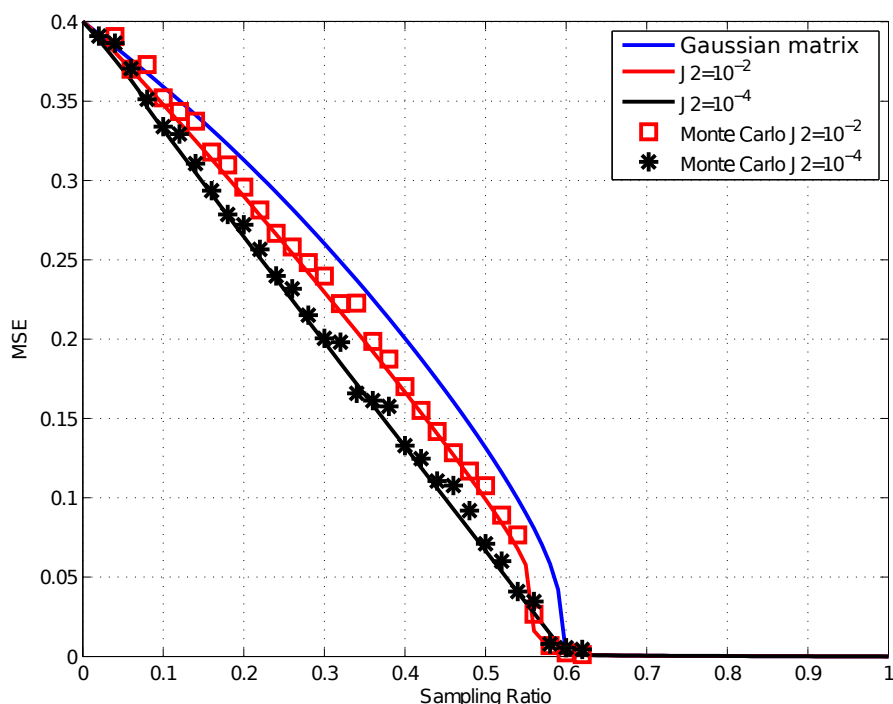


Figure 4.8: The normalized SD function for the compressible signal $p_{GM_1}(x)$ with different measurement matrix configuration. For two-block matrix, $\gamma_1 = \gamma/\gamma_c$ for $\gamma < \gamma_c$, where $\gamma_c = 0.6$ is the critical sampling ratio for the homogeneous Gaussian matrix.

the optimal limit (the sparsity level). Setting $J_2 = 10^{-14}$ moves the perfect reconstruction ratio back to 0.59. In fact, its SD function is very close to the hybrid zeroing matrix performance with $J_2 = 0$. Comparing to the seeded matrix, the two block matrices exhibit an improved reconstruction for $\gamma < 0.4$. With the suggested configuration, the seeded matrix achieves the perfect reconstruction at $\gamma = 0.5$.

One thing worth noting is that even with the improved performance, the two block matrix still has a concave SD function up to a new critical sampling ratio. A further convexifying procedure with a three-block structure can then be easily applied to achieve a slightly better reconstruction. In fact, if we introduce multiple J_i , we conjecture that this approach will tend to the optimal recovery as with the seeded matrix. Note again that for the multi-block matrix structure, the SE equation is still 1-D. The Monte Carlo simulation implies that for the finite size problem the SE prediction is accurate.

4.5.2 Two Block Matrix for a Compressible Signal

In this section, we present the SD functions for two different two-state Gaussian mixture priors. In Fig. 4.8, we have the SD functions for

$$p_{\text{GM}_1}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 5 \times 10^{-4}) \quad (4.29)$$

As demonstrated in Section 4.4, this signal model also exhibits a FOPT. It is not surprising to see a very similar SD curve as in Fig. 4.7 for BAMP with the homogeneous Gaussian matrix, since the small Gaussian variance is very close to zero. For the homogeneous Gaussian matrix with BAMP, the FOPT happens at $\gamma_c = 0.6$. For the two block matrix with $J_2 = 10^{-2}$, we can successfully remove the discontinuity point of the SD function to $\gamma_c = 0.56$. Decreasing the rescaling parameter J_2 to 10^{-4} further improves the SD performance.

Next we show some results for a compressible signal without FOPT. The signal is drawn i.i.d. from

$$p_{\text{GM}_2}(x) = 0.4\mathcal{N}(x; 0, 1) + 0.6\mathcal{N}(x; 0, 0.003) \quad (4.30)$$

which is motivated from the statistics of natural images. The SD functions, as well as the achievable model based bound for the prior are shown in Fig. 4.9. Similarly to the sparse signal case, the two block matrices outperform the homogeneous Gaussian matrix up to the same critical sampling ratio γ_c . Also, the SD performance is better as we decrease J_2 . We obtained an excellent agreement between the SD prediction and the Monte Carlo simulation. Empirically, we observed that the optimum weighting for J_2 is zero for the compressible signal without FOPT. This suggests that without FOPT, the only gains come from the convexification of the SD function. However, the proof remains an open question.

4.5.3 Two Block Matrix for Dense Signals

In this section, the two block matrix design is applied to the k -dense signal model, which has been defined in Chapter 2.

$$p_{\text{KD}}(x) = 0.45\delta(x + 1) + 0.45\delta(x - 1) + 0.1\mathcal{U}(-1, 1) \quad (4.31)$$

It has been observed in [116] and proved in [117] (Proposition 3.12) that the k -dense signal can

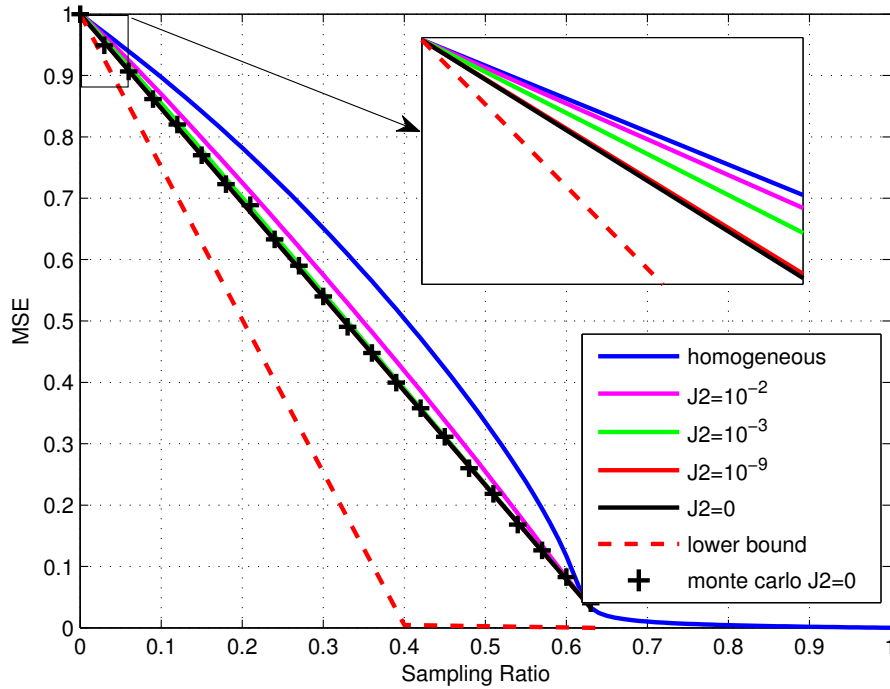


Figure 4.9: The normalized SD function for the compressible signal $p_{GM_2}(x)$ with different measurement matrix configuration. For the two-block matrix, $\gamma_1 = \gamma/\gamma_c$ for $\gamma < \gamma_c$, where $\gamma_c = 0.63$ is the critical sampling ratio for the homogeneous Gaussian matrix.

be reconstructed with high probability by solving the following convex optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\tilde{\mathbf{x}}} \|\tilde{\mathbf{x}}\|_{\ell_\infty} \text{ s.t. } \mathbf{y} = \Phi \tilde{\mathbf{x}} \quad (4.32)$$

In both [13] and [14], the authors concluded that the sampling ratio is required to be more than 0.5 to ensure successful recovery for the k -dense signal with convex optimization. In [116], an iterative, fast method is proposed to solve (4.32). We plot the empirical SD function for the convex optimization with the homogeneous Gaussian measurement matrix as the benchmark in Fig. 4.10.

The resulting SD functions for different measurement matrices and reconstruction algorithms are presented in Fig. 4.10. As expected, the convex recovery algorithm with the homogeneous Gaussian matrix will not achieve perfect reconstruction until $\gamma = 0.62$. While BAMP with the homogeneous Gaussian matrix pushes the perfect sampling ratio down to 0.49. Similar to the sparse signal with a FOPT, there is also a sudden drop of MSE at the perfect recovery ratio. Thus we can expect a similar fixed point change as in Fig. 4.6.(b) when the two block matrix is applied. Unsurprisingly, the two block matrix manages to largely reduce the perfect

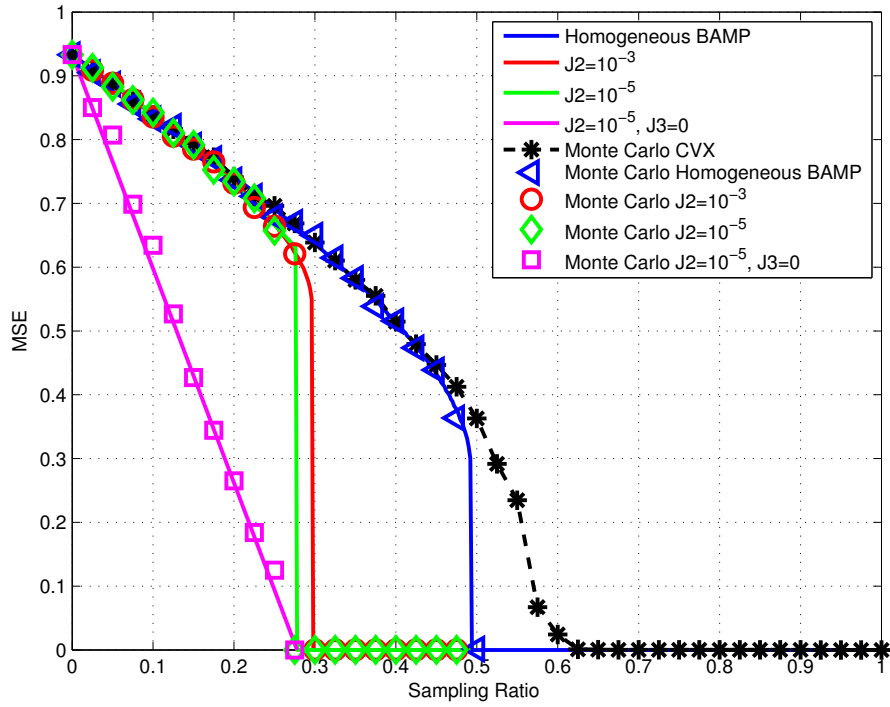


Figure 4.10: The SD function for the dense signal with different measurement matrix configuration. With $\gamma_1 = 0.6$, $J_2 = 10^{-5}$, the two block matrix is able to move the perfect reconstruction sampling ratio from 0.49 to 0.29.

recovery to $\gamma_c = 0.29$. The combination of the two block matrix and TAP-AMP demonstrates a dramatic improvement over both convex optimization and BAMP with homogeneous Gaussian measurement matrix. Moreover, the reconstruction performance can be further improved by convexifying the SD function of the two block matrix. Monte Carlo results also confirm the 1-D SE prediction.

4.6 Summary

In the chapter, a novel measurement matrix, the modulated matrix, is introduced. With the simple 1-D dynamics and the flexible rescaling matrix, it provides us a whole range of measurement matrix design. As a special case, we understand the advantage and limitation of the two block matrix based on the analysis of the first order phase transition. Extensive simulations with sparse, compressible and dense signals demonstrate that with the two block matrix, better recovery quality is achievable in the least squared sense. For the sparse signal, the two block matrix does not push the perfect reconstruction sampling ratio down to the sparsity level but close. Part of the reason is that the two block design is still crude. A further research direction

involves examination of multi-block modulated matrices and parameter optimization. Different rescaling distributions should also be considered. As discussed in [53], the additive noise level has a significant impact on the reconstruction performance for the seeded matrices. As a relative measurement matrix design, we believe the modulated matrix is also likely to be sensitive to the noise level. Understanding the relationship between the noise sensitivity of the modulated matrices in comparison with the seeded matrices is the topic worthy of further investigation.

Chapter 5

Bayesian optimal reconstruction without priors: parametric SURE-AMP algorithm

The generic AMP algorithm reviewed in Chapter 2 is revisited here to present some motivations for the enhancement of the algorithm. As we can see from the previous chapters, both theoretical analysis and empirical evidence confirm that the AMP algorithm can be interpreted as recursively solving a signal denoising problem: at each AMP iteration, one observes a Gaussian noise perturbed original signal. Retrieving the signal amounts to a successive noise cancellation until the noise variance decreases to a satisfactory level. In this chapter we incorporate the SURE based parametric denoiser with the AMP framework and propose the novel parametric SURE-AMP algorithm. At each parametric SURE-AMP iteration, the denoiser is adaptively optimized within the parametric class by minimizing SURE, which depends purely on the noisy observation. In this manner, the parametric SURE-AMP is guaranteed with the best-in-class recovery and convergence rate. If the parameter family includes the family of the MMSE estimators, we are able to achieve the BAMP performance without knowing the signal prior. In the chapter, we resort to the linear parameterization of the SURE based denoiser and propose three different kernel families as the base functions. Numerical simulations with the BG, k -dense and Student's-t signals demonstrate that the parametric SURE-AMP does not only achieve the state-of-the-art recovery but also runs more than 20 times faster than the EM-GM-GAMP algorithm.

5.1 Introduction

Recall that the generic AMP algorithm for the CS system (2.1) takes the simple iterative form:

$$\mathbf{r}^t = \hat{\mathbf{x}}^t + \Phi^T \mathbf{z}^t \quad (5.1)$$

$$\hat{\mathbf{x}}^{t+1} = \eta_t(\mathbf{r}^t) \quad (5.2)$$

$$\mathbf{z}^{t+1} = \mathbf{y} - \Phi \hat{\mathbf{x}}^{t+1} + \frac{1}{\gamma} \mathbf{z}^t < \eta'_t(\mathbf{r}^t) > \quad (5.3)$$

Initialized with $\hat{\mathbf{x}}^0 = \mathbf{0}$ and $\mathbf{z}^0 = \mathbf{y}$, AMP iteratively produces an estimation of the original signal $\hat{\mathbf{x}}^t$ with a scalar non-linear function $\eta_t(\cdot)$, which is applied elementwise to \mathbf{r}^t . Throughout this chapter we assume the elements of Φ are drawn i.i.d from $\mathcal{N}(\Phi_{i,j}, 0, m^{-1})$. As discussed in Chapter 2, the key feature of AMP is the Onsager reaction term $\frac{1}{\gamma} \mathbf{z}^t < \eta'_t(\mathbf{r}^t) >$, which guaranteed the Gaussian behaviour of the AMP residual $\mathbf{r}^t - \mathbf{x}^t$ at each iteration. We display the Gaussianity of the residual for AMP iterations again in Fig. 5.1. The QQ plot of the empirical pdf of $\mathbf{r}^t - \mathbf{x}$ against the normal distribution at various iteration conform its Gaussian behaviour. In other words, we approximately have $\mathbf{r}^t \approx \mathbf{x} + \sqrt{\tau_t} \mathbf{z}^t$, $z_i \in \mathcal{N}(z_i; 0, 1)$, where τ_t is the effective noise variance [14, 72] at each AMP iteration. Then the non-linearity $\eta_t(\cdot)$ essentially acts as a denoising function to remove the Gaussian noise $\sqrt{\tau_t} \mathbf{z}^t$.

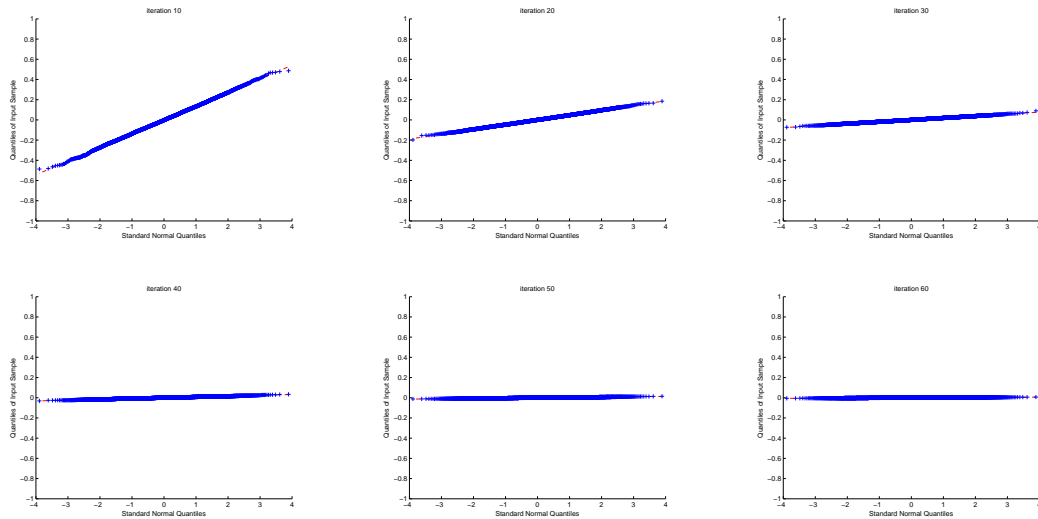


Figure 5.1: QQ plots tracking the effective noise of the AMP algorithm under various iterations while reconstruction a 40% sampled Bernoulli-Gaussian data with pdf in (5.35). The residual of AMP remains Gaussian because of the Onsager reaction term. Decreasing slope as the iteration increasing indicates the decreasing standard deviation.

Treating AMP reconstruction as an iterative denoising procedure, we can reconsider the ℓ_1 -AMP and BAMP algorithm introduced in Chapter 2 from the signal denoising prospective. In the original AMP paper [14, 118], the denoising is achieved with the simple soft thresholding function. Despite the fact that the noisy vector \mathbf{r}^t has multiple i.i.d. distributed elements, the ℓ_1 -AMP treats the denoising as a 1-D problem. However, since the true signal pdf is visible in the noisy estimate in the large system limit and the effective noise variance is estimated at each AMP iteration, we should be able to exploit such information to achieve better recovery than the ℓ_1 -AMP. The BAMP algorithm deploys the MMSE estimator for denoising and achieves the best reconstruction in the least square sense. However, the requirement of $p(\mathbf{x})$ to be known in advance can be restrictive in practice. The advantages and limitation of BAMP also motive us to find an alternative approach which is able to fill the gap between the ℓ_1 -AMP and the BAMP, or even performs as well as BAMP without knowing the signal distribution a priori.

Main contributions

In the large system limit, the true prior for \mathbf{x} at each AMP iteration is essentially embedded in the data \mathbf{r}^t , which is the convolution of the original signal with the Gaussian noise kernel. To improve the recovery, we could either estimate the pdf and then deduce the associated MMSE estimator, or directly optimize the denoising. In this chapter, we adopt the latter approach and propose the parametric SURE-AMP algorithm. Realizing the recursive denoising nature of the AMP iteration, we introduce a class of parameterized denoising functions to the generic AMP framework. At each iteration, the denoiser with the least MSE is selected within the class by optimizing the free parameters. In this manner, the parametric SURE-AMP algorithm adaptively chooses the best-in-class denoiser and achieves the best possible denoising within the parametric family at each iteration. When the denoiser class contains all possible MMSE estimators for a specific signal, the parametric SURE-AMP is expected to achieve the BAMP recovery without knowing the signal prior in the large system limit.

The key feature of the parametric SURE-AMP algorithm is that the denoiser optimization does not require prior knowledge of $p(\mathbf{x})$. To make this possible, we resort to the SURE based parametric least squarer denoiser construction. There exists a rich literature on signal denoising with SURE [119–125]. Since SURE is the unbiased estimate of MSE, the pursuit of the best denoiser with the least MSE is nothing more than minimizing the corresponding SURE. More importantly, for a Gaussian noise corrupted signal, the calculation of SURE depends purely on the sampled average of the noisy data [126]. By leveraging the large system limit, the best-in-

class denoiser can be determined without the prior information [125].

The success of the parametric SURE-AMP relies heavily on the parameterization of the denoiser class. The number of parameters as well as the linearity determine the optimization complexity. In this chapter, we restrict ourselves to the linear combination of non-linear kernel functions as the denoiser structure. The non-linear parameters of the kernel functions are set to have a fixed ratio with the effective noise variance. The linear weights for the kernels are optimized by solving a linear system of equations. We presented two types of piecewise linear kernel family and one exponential kernel family for both sparse and heavy-tailed signal reconstruction. The numerical simulation with the BG, k -dense and Student's-t signals show that with a limited number of kernel functions, we are able to adaptively capture the evolving shape of the MMSE estimator and achieves the state-of-art performance in the sense of reconstruction quality and computational complexity.

Related literature

The pre-requisite of the signal prior to implement BAMP has been noticed by several research groups. To tackle this limitation, the prior estimation step was proposed to be incorporated within the AMP framework. The corresponding EM-GM-GAMP has been summarized in Chapter 2, page 33. The key difference between the EM-GM-GAMP and the parametric SURE-AMP is that fitting the signal prior is an indirect adaptation for minimizing the reconstruction MSE while we directly tackle the problem by adaptively selecting the best-in-class denoiser with the least MSE. When the signal distribution can be well approximated by a GM model, fitting the prior and minimizing MSE lead to subtle difference. However, for distributions that are difficult to be approximated as the finite sum of Gaussians, as we demonstrate later in section 5.4, the parametric SURE-AMP algorithm provides a better solution. In terms of computational complexity, the parametric SURE-AMP significantly outperforms the EM-GM-GAMP with the linear parameterization of the denoisers.

In [73], the authors generalized the EM step with an adaptive prior selection function. The proposed adaptive GAMP algorithm includes the EM-GM-GAMP as a special case. Although the general form of the prior adaptation also enables other learning methods, i.e. ML, to be deployed in the AMP framework, in principle the adaptive GAMP still focuses on fitting the signal prior rather than directly minimizing the reconstruction MSE.

Another relevant work is the denoising-based AMP (D-AMP) algorithm [51]. The intrinsic

denoising problem within AMP iterations has also been noticed by the authors. The intuition for D-AMP is to take advantage of the rich existing literature on signal denoising to enhance the AMP algorithm. In the paper, the existing image denoising algorithm BM3D has been utilized as the denoiser in D-AMP and produced the state-of-art recovery for natural images. The authors essentially share the same understanding as us for the AMP algorithm and point out the possibility of using the SURE based estimator for denoising.

Structure of the Chapter

The remainder of the chapter is organized as follows: The parametric SURE-AMP algorithm is presented in Section 5.2. Section 5.3 is devoted to introducing the construction of the SURE-based parametric denoiser class. Three types of kernel families as well as the parameter optimization scheme are discussed herein. The simulation results are summarized in Section 5.4. It compares both the reconstruction performance and the computational complexity of the parametric SURE-AMP algorithm with other CS algorithms. We conclude the chapter in Section 5.5.

5.2 Parametric SURE-AMP Framework

5.2.1 Parametric SURE-AMP algorithm

Algorithm 7 : Parametric SURE-AMP

```

1: initialization:  $\hat{\mathbf{x}}^0 = \mathbf{0}, \mathbf{z}^0 = \mathbf{y}, c^0 = \langle \|\mathbf{z}^0\|^2 \rangle$ 
2: for  $t = 0, 1, 2, \dots$  do
3:    $\mathbf{r}^t = \hat{\mathbf{x}}^t + \Phi^T \mathbf{z}^t$ 
4:    $\boldsymbol{\theta}^t = H_t(\mathbf{r}^t, c^t)$ 
5:    $\hat{\mathbf{x}}^{t+1} = f_t(\mathbf{r}^t, c^t | \boldsymbol{\theta}^t)$ 
6:    $\nu^{t+1} = \langle f'_t(\mathbf{r}^t, c^t | \boldsymbol{\theta}^t) \rangle$ 
7:    $\mathbf{z}^{t+1} = \mathbf{y} - \Phi \hat{\mathbf{x}}^{t+1} + \frac{1}{\gamma} \nu^{t+1} \mathbf{z}^t$ 
8:    $c^{t+1} = \langle \|\mathbf{z}^{t+1}\|^2 \rangle$ 
9: end for

```

We begin with a description of the parametric SURE-AMP algorithm, which extends the generic AMP iteration defined in (5.1), (5.2) and (5.3) with an adaptive signal denoising module. The implementation of the parametric SURE-AMP algorithm is summarized in Algorithm 7.

Most of the entities have the same interpretation as in AMP: \mathbf{r}^t is the noisy version of the original signal, which can be effectively approximated as $\mathbf{r}^t \approx \mathbf{x} + \sqrt{c^t} \mathbf{z}^t$, $z_i \sim \mathcal{N}(z_i; 0, 1)$.

Here c^t is the estimation of the effective noise variance. A new signal estimate $\hat{\mathbf{x}}^{t+1}$ is obtained by denoising \mathbf{r}^t at each iteration. The key modification to AMP is the introduction of the parametric denoising function $f_t(\cdot|\boldsymbol{\theta}^t)$ and the parameter selection function $H_t(\cdot)$. Consider a class of denoising functions $\mathbb{F}(\cdot|\mathbf{Q})$ characterized by the parameter set \mathbf{Q} . At each iteration, the best-in-class denoiser $f_t(\cdot|\boldsymbol{\theta}_t) \in \mathbb{F}(\cdot|\mathbf{Q})$ is chosen by selecting the parameter $\boldsymbol{\theta}^t$ via the parameter selection function $H_t(\cdot)$. We design $H_t(\cdot)$ as a function of the noisy data \mathbf{r}^t and the effective noise variance c^t to close the parametric SURE-AMP iteration.

The next question is what should be the parameter selection criteria for the parametric SURE-AMP algorithm. Our fundamental reconstruction goal is to obtain a signal estimate $\hat{\mathbf{x}}$ with the MMSE. Theoretically speaking, we want to jointly select the denoisers across all iterations. However, solving the joint optimization is not trivial. Based on the state evolution analysis in the subsequent section, we propose to break the joint selection into separate independent steps. Specifically, the parameter vector $\boldsymbol{\theta}^t$ at iteration t is selected by solving

$$\begin{aligned} \boldsymbol{\theta}^t &= \arg \min_{\boldsymbol{\theta}} \mathbb{E}[(\hat{\mathbf{x}}^{t+1} - \mathbf{x})^2] \\ &= \arg \min_{\boldsymbol{\theta}} \mathbb{E}\{[f_t(\mathbf{r}^t, c^t|\boldsymbol{\theta}) - \mathbf{x}]^2\} \end{aligned} \tag{5.4}$$

which achieves the MMSE among the parametric family. As the signal estimate $\hat{\mathbf{x}}^t$ is optimized within the denoiser class at each step, one would expect to obtain a "global" optimal reconstruction as the algorithm converges.

5.2.2 SURE based denoiser selection

The Stein's unbiased estimate is an unbiased estimate for MSE. It becomes more accurate as more data is available, which is particularly apt for AMP since it is designed with the large system limit in mind. It has been widely used as the surrogate for the MSE to tune the free parameters of estimation functions for signal denoising. In [126], it has been proved that for the Gaussian noise corrupted signal, the calculation of SURE can be performed entirely in terms of the noisy observation. This property is summarized in the following theorem.

Theorem 5. [126] *Let x be the signal of interest and $r = x + \sqrt{c}z$ be noisy observation with $z \sim \mathcal{N}(z; 0, 1)$. Without loss of generality, we assume the denoising function $f(r, c|\boldsymbol{\theta})$ is parameterized by $\boldsymbol{\theta}$ and has the form*

$$f(r, c|\boldsymbol{\theta}) = r + g(r, c|\boldsymbol{\theta}) \tag{5.5}$$

The denoised signal is obtained through $\hat{x} = f(r, c|\boldsymbol{\theta})$. Additionally, we assume as γ goes to $\pm\infty$, $p(\gamma)$ dies off faster than $g(\gamma, c|\boldsymbol{\theta})$ grows. That is, $p(\gamma)g(\gamma, c|\boldsymbol{\theta})|_{\pm\infty} = 0$. Then SURE is defined as the expected value over the noisy data alone and is the unbiased estimate of the MSE. That is,

$$\begin{aligned}\mathbb{E}_{\hat{x},x}\{(\hat{x} - x)^2\} &= \mathbb{E}_{r,x}\{[f(r, c|\boldsymbol{\theta}) - x]^2\} \\ &= c + \mathbb{E}_r\{g^2(r, c|\boldsymbol{\theta}) + 2cg'(r, c|\boldsymbol{\theta})\}\end{aligned}\tag{5.6}$$

Proof. Given the parametric form of the estimator $f(\cdot)$ we have

$$\mathbb{E}_{r,x}\{(f(r, c|\boldsymbol{\theta}) - x)^2\}\tag{5.7}$$

$$= \mathbb{E}_{r,x}\{(r + g(r, c|\boldsymbol{\theta}) - x)^2\}\tag{5.8}$$

$$= \mathbb{E}_r\{g^2(r, c|\boldsymbol{\theta})\} + 2\mathbb{E}_{r,x}\{g(r, c|\boldsymbol{\theta})(r - x)\} + \mathbb{E}_{r,x}\{(r - x_o)^2\}\tag{5.9}$$

$$= \mathbb{E}_r\{g^2(r, c|\boldsymbol{\theta})\} + 2\mathbb{E}_{r,x}\{g(r, c|\boldsymbol{\theta})(r - x)\} + c\tag{5.10}$$

The middle term in (5.10) can be further written as

$$\begin{aligned}\mathbb{E}_{r,x}\{g(r, c|\boldsymbol{\theta})(r - x)\} &= \mathbb{E}_r\{g(r, c|\boldsymbol{\theta})[r - \mathbb{E}_{x|r}(x|r)]\} \\ &\stackrel{(a)}{=} \mathbb{E}_r\{g(r, c|\boldsymbol{\theta})[r - r - c\frac{p'(r)}{p(r)}]\} \\ &= -c\mathbb{E}(g(r, c|\boldsymbol{\theta})\frac{p'(r)}{p(r)}) \\ &= -c\int g(r, c|\boldsymbol{\theta})\frac{p'(r)}{p(r)}p(r)dr \\ &= -c\int g(r, c|\boldsymbol{\theta})p'(r)dr \\ &\stackrel{(b)}{=} c\int g'(r, c|\boldsymbol{\theta})p(r)dr \\ &= c\mathbb{E}_r\{g'(r, c|\boldsymbol{\theta})\}\end{aligned}\tag{5.11}$$

where we use the following observation

- (a) The MMSE estimator for the Gaussian noise corrupted data can be written entirely in terms of the measurement density [127]

$$\mathbb{E}_{x|r}(x|r) = r + c\frac{p'(r)}{p(r)}\tag{5.12}$$

The proof is in Appendix C.

- (b) Apply integration by parts and the assumption $p(r)g(r, c|\boldsymbol{\theta})|_{-\infty}^{\infty} = 0$

Combine (5.11) with (5.10) completes the proof. \square

According to Theorem 5, the parameter selection for the parametric SURE-AMP algorithm can thus be conducted via the minimization of SURE. By the law of large numbers, the expectation in (5.6) can be approximated as the average over multiple realizations of the noisy data r . For parametric SURE-AMP, we naturally have a vector \mathbf{r}^t at each iteration. Since the term c will disappear in the minimization of (5.6), the corresponding parameter selection function is thus defined as

$$\begin{aligned}\boldsymbol{\theta}^t &= H_t(\mathbf{r}^t, c^t) \\ &= \arg \min_{\boldsymbol{\theta}} \langle g^2(\mathbf{r}^t, c^t|\boldsymbol{\theta}) + 2c^t g'(\mathbf{r}^t, c^t|\boldsymbol{\theta}) \rangle\end{aligned}\tag{5.13}$$

It fundamentally eliminates the dependency on the original signal for selecting the denoisers with the minimum MSE. Applying (5.13) into line 4 of Algorithm 7 we have a complete parametric SURE-AMP algorithm.

5.2.3 State evolution

As reviewed in Chapter 2, the asymptotic behaviour of the AMP algorithm can be accurately characterized by the SE formalism in the large system limit. As an extension of the AMP algorithm, one expects the parametric SURE-AMP would also follow the SE analysis incorporating the denoising adaptation. We hereby formally summarize our finding:

Finding 1. *Starting with $\tau^0 = \frac{\|\mathbf{y}\|^2}{m}$, the state evolution equation for the parametric SURE-AMP algorithm has the following iterative form*

$$\bar{\boldsymbol{\theta}}^t = H_t(x + \sqrt{\tau^t}z, \tau^t)\tag{5.14}$$

$$\tau^{t+1} = \sigma_w^2 + \frac{1}{\gamma} \mathbb{E}\{\tau^t f'_t(x + \sqrt{\tau^t}z, \tau^t|\bar{\boldsymbol{\theta}}^t)\}\tag{5.15}$$

where $x \sim p(x)$ has the same distribution as the original signal, $z \sim \mathcal{N}(z; 0, 1)$ is the white Gaussian noise. In the large system limit, i.e. $m \rightarrow \infty$, $n \rightarrow \infty$ with $\gamma = m/n$ fixed, the MSE of parametric SURE-AMP estimate at iteration t can be predicted as

$$\mathbb{E}\{(\mathbf{x} - \hat{\mathbf{x}}^t)^2\} = \sigma_w^2 + \frac{1}{\gamma} \mathbb{E}\left\{\left[x - f_t(x + \sqrt{\tau^t}z, \tau^t|\bar{\boldsymbol{\theta}}^t)\right]^2\right\}\tag{5.16}$$

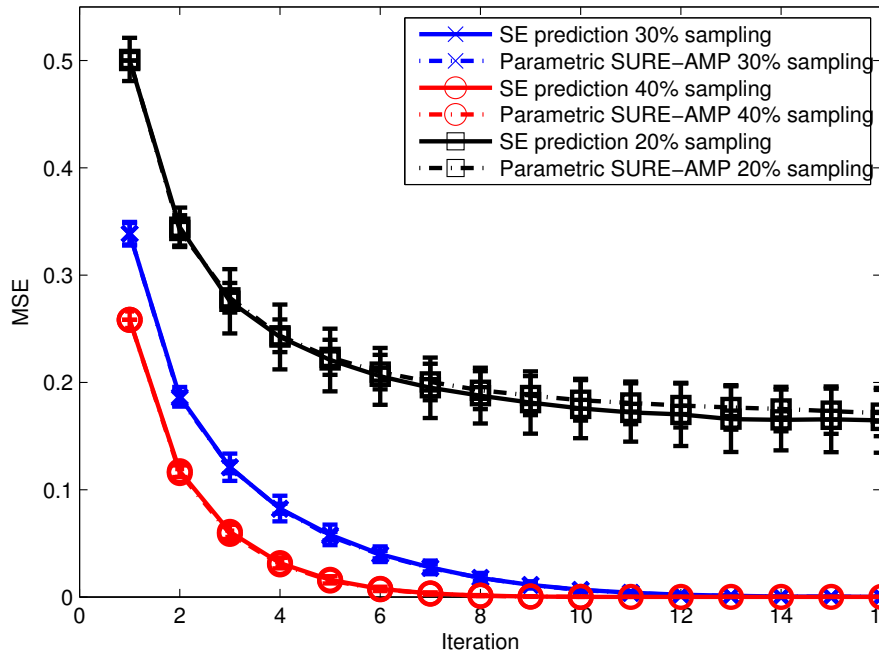


Figure 5.2: The actual MSE for the noiseless Bernoulli-Gaussian data reconstruction at each parametric SURE-AMP iteration versus the state evolution prediction. The signal is generated i.i.d. according to eq. (5.35). The first piecewise linear kernel family is utilized within the parametric SURE-AMP algorithm, which will be discussed in section 5.3.1.1. The reconstruction MSE is an average over 100 Monte Carlo realizations.

We use the term Finding here to emphasize the lack of rigorous proof. However, the empirical simulation supports our finding. In Fig. 5.2, the state evolution prediction for the noiseless BG signal reconstruction with the parametric SURE-AMP algorithm is compared against the Monte Carlo average at multiple iterations. It is clear from the figure that at various sampling ratios, SE accurately predicts the MSE of the parametric SURE-AMP reconstruction.

Finding 1 coincides with the SE analysis for the adaptive GAMP algorithm in [73] when the output channel is assumed to be Gaussian white noise and $H_t(\cdot)$ is the prior fitting function. The authors have proved that when $H_t(\cdot)$ has the weak pseudo-Lipschitz continuous property and the denoising function $f_t(\cdot|\theta^t)$ is Lipschitz continuous, the adaptive GAMP can be asymptotically characterized by the corresponding state evolution equations in the large system limit. Unfortunately, their analysis does not apply directly to the parametric SURE-AMP algorithm since our $H_t(\cdot)$ and $f_t(\cdot|\theta^t)$ do not satisfy the required pseudo-Lipschitz continuous properties. The theoretical proof of Finding 1 is beyond the scope of this work and remains an open question for further study.

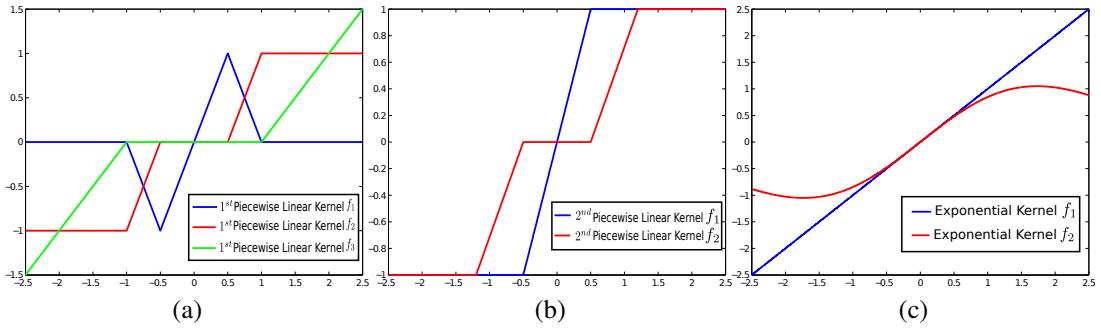


Figure 5.3: Kernel families used for linear parameterization of the SURE based denoiser: (a) the first piecewise linear kernel family (b) The second piecewise linear kernel family. (c) The exponential kernel family.

5.3 Construction of the Parametric Denoiser

The reconstruction quality of the parametric SURE-AMP algorithm primarily counts on the construction of the adaptive denoiser class and the tuning of free parameters. Inspired by the SURE-LET algorithm for image denoising in [124], we choose to form the denoiser $f_t(\cdot|\boldsymbol{\theta}^t)$ as a weighted sum of some kernel functions to give an adaptive non-linearity. To be specific,

$$f_t(\mathbf{r}^t, c^t|\boldsymbol{\theta}^t) = \sum_{i=1}^k a_{t,i} f_{t,i}(\mathbf{r}^t|\boldsymbol{\vartheta}_{t,i}(c^t)) \quad (5.17)$$

where $f_{t,i}(\mathbf{r}^t|\boldsymbol{\vartheta}_{t,i}(c^t))$ is the non-linear kernel function with $\boldsymbol{\vartheta}_{t,i}(c^t)$ summarizes all non-linear parameters that depend on the effective noise variance c^t . The linear weight for the kernel function is represented with $a_{t,i}$. At each parametric SURE-AMP iteration, we need to select the parameter set $\boldsymbol{\theta}^t = [a_{t,i}, \boldsymbol{\vartheta}_{t,i}]_{i=1}^k$ to obtain the best denoising function in the class. For the rest of this section, we drop the iteration index t to simplify the notation.

This parameterization method for denoisers has been used before. In [125], the ‘‘bump’’ kernel family is designed to approximate the MMSE estimator of the generalized Gaussian signal. In [122–124], the exponential kernels are specifically designed for natural image denoising in the transformed domain. In this section, we start by presenting three types of kernel families for both sparse and heavy-tailed signal denoising. Then we will explain the parameter selection rule for both linear and non-linear parameters of the kernels. Finally the constructed denoiser is applied to three different signal priors to validate the design.

5.3.1 Kernel families

5.3.1.1 First piecewise linear kernel family

The underlying principle for the kernel function design is to keep it simple and flexible at the same time. One way to do this is to use the piecewise linear function as the kernel format and proposed the first piecewise linear kernel family which consists three kernel functions:

$$f_1(r|\alpha_1) = \begin{cases} 0 & r \leq -2\alpha_1, r \geq 2\alpha_1 \\ -\frac{r}{\alpha_1} - 2 & -2\alpha_1 < r < -\alpha_1 \\ \frac{r}{\alpha_1} & -\alpha_1 < r \leq \alpha_1 \\ -\frac{r}{\alpha_1} + 2 & \alpha_1 < r < 2\alpha_1 \end{cases} \quad (5.18)$$

$$f_2(r|\alpha_1, \alpha_2) = \begin{cases} -1 & r \leq -\alpha_2 \\ \frac{r+\alpha_1}{\alpha_2-\alpha_1} & -\alpha_2 < r < -\alpha_1 \\ 0 & -\alpha_1 \leq r \leq \alpha_1 \\ \frac{r-\alpha_1}{\alpha_2-\alpha_1} & \alpha_1 < r < \alpha_2 \\ 1 & r \geq \alpha_2 \end{cases} \quad (5.19)$$

$$f_3(r|\alpha_2) = \begin{cases} r + \alpha_2 & r \leq -\alpha_2 \\ 0 & -\alpha_2 < r < \alpha_2 \\ r - \alpha_2 & r \geq \alpha_2 \end{cases} \quad (5.20)$$

where $\alpha_1 > 0$ and $\alpha_2 > 0$ are hinge points closely related to the effective noise level c . The three kernels are plotted in Fig. 5.3(a). Eq. (5.20) is the soft thresholding function to promote sparsity. It sets all vector elements whose magnitude smaller than α_2 to zero and keeps the linear behaviour of large elements. The linear part with positive gradient in (5.18) aims to soften the “brutal” correction of the soft thresholding function on the small elements. It is designed for removing the Gaussian perturbation for small but non-zero elements of compressible signals. Eq. (5.19) is constructed to add a denoising transition between the small and large elements to increase the denoiser flexibility. With proper rescaling of the three kernels and appropriate setting for the hinge points, we expect the denoiser class constructed with the first piecewise linear kernels to be flexible and accurate enough to capture the evolving shape of the MMSE estimators for various CS signals at different noise levels.

5.3.1.2 Second piecewise linear kernel family

As we presented in Chapter 2, the k -dense model is a non-conventional CS prior. In [14], the soft thresholding function with the adaptive thresholding level is suggested as a generic AMP algorithm for such signals. For the k -dense signals, we propose the second piecewise linear kernel functions to construct the denoiser.

$$f_1(r|\beta_1) = \begin{cases} -1 & r \leq -\beta_1 \\ \frac{r}{\beta_1} & -\beta_1 < r < \beta_1 \\ 1 & r \geq \beta_1 \end{cases} \quad (5.21)$$

$$f_2(r|\beta_1, \beta_2) = \begin{cases} -1 & r \leq -\beta_2 \\ \frac{r+\beta_1}{\beta_2-\beta_1} & -\beta_2 < r < -\beta_1 \\ 0 & -\beta_1 \leq r \leq \beta_1 \\ \frac{r-\beta_1}{\beta_2-\beta_1} & \beta_1 < r < \beta_2 \\ 1 & r \geq \beta_2 \end{cases} \quad (5.22)$$

Similar to the first piecewise linear kernel family, the hinge points β_1 and β_2 depend on the effective Gaussian noise level. With proper scaling of the second piecewise linear kernels, the constructed denoiser is able to mimic the MMSE estimator behaviour for the k -dense signal under different noise levels.

5.3.1.3 Exponential kernel family

For the third type of kernel family, we resort to more sophisticated exponential functions.

$$f_1(r) = r \quad (5.23)$$

$$f_2(r|T) = r e^{-\frac{r^2}{2T^2}} \quad (5.24)$$

This kernel family is motivated from the derivatives of Gaussians (DOG) and has been used for natural image denoising in the transformed domain in [122–124]. The virtue of DOGs is that they decay rapidly and ensure a linear behaviour close to the identity for large elements [123]. It has been demonstrated that with kernels defined in (5.23) and (5.24), the constructed denoiser delivers the near-optimal performance regarding both quality and computational cost. The parameter T in (5.24) has the same functionality as the hinge points for the piecewise linear

kernels. It controls the transition between small and large elements and is linked tightly with the effective noise variance. Given the fact that most natural images are compressible in the wavelet or DCT domain, we believe that the exponential kernel family used for image denoising can also be applied in the parametric SURE-AMP algorithm to recovery compressible signals.

One thing worth noting is that the proposed kernel families are not designed to fit any specific signal prior, but are motivated from the general sparse or compressible pattern. Thus they are, to some extent, suitable for many CS signal reconstructions. It is also straightforward to construct new kernel functions to increase the sophistication of the constructed denoiser. For the exponential kernel family, high order DOGs can be used. For the piecewise linear kernel families, more functions with various hinge points could be added. In our work, we find that with just three kernel functions, the constructed denoiser is able to deliver a near Bayesian optimal performance. Moreover, we do not necessarily require the denoiser class to contain the true MMSE estimator to achieve good reconstruction performance. As proved in [128], denoisers constructed by the piecewise linear kernels are not eligible for the true MMSE estimator since they are not in $C^\infty(\mathbb{R}^n)$. Nevertheless, they exhibit excellent performance for the CS signal denoising and integrate well with the parametric SURE-AMP reconstruction as we will see later. When the parametric denoiser class includes all possible MMSE estimators for a specific prior, the parametric SURE-AMP algorithm is guaranteed to obtain the BAMP recovery in the large system limit.

5.3.2 Non-linear parameter tuning for kernel functions

To cope with the developing noise level during the parametric SURE-AMP iteration, the aforementioned kernel functions all have some non-linear dependency, i.e. the hinge point and the variance for the exponential kernel. While the non-linearity is necessary, finding the global optimizer for the non-linear parameter can be computationally expensive. To mitigate this problem, we propose a fixed linear relationship between the non-linear parameters and the effective noise level. Since at each parametric SURE-AMP iteration we obtain an estimated effective noise variance c^t , the non-linear parameters are consequently selected. In this section, we will explain the non-linear parameter tuning for all three kernel families.

5.3.2.1 First piecewise linear kernel family

In [82], the authors discussed the thresholding choice for iterative reconstruction algorithms for compressed sensing. For iterative soft thresholding, they proposed to set the threshold as a fixed multiple of the standard derivation of the effective noise variance. The rule of thumb for the multiple is between 2 and 4. This threshold choice has been tested with the StOMP algorithm [36] and the underlying rationale has been explained therein. For the first piecewise linear kernel family which has the soft thresholding element, we take their recommendation and set the hinge points as

$$\alpha_1 = 2\sqrt{c}, \quad \alpha_2 = 4\sqrt{c} \quad (5.25)$$

5.3.2.2 Second piecewise linear kernel family

In [129], a novel iterative dense recovery (IDR) algorithm is proposed to replace the MMSE estimator for the k -dense signal with an adaptive denoiser within the AMP iteration. The essence of the IDR is the employment of a piecewise linear function with one flexible hinge point to approximate the MMSE estimator class. Inspired by the selection of hinge point in [129], we choose to fix the linear ratio for the second piecewise linear kernels as following

$$\beta_1 = \frac{1}{1 + 6\sqrt{c}}, \quad \beta_2 = \frac{1}{1 + 2\sqrt{c}} \quad (5.26)$$

The ratio in (5.26) is based on the empirical denoising experiments with k -dense signals under different noise levels. Although not very critical, we find it to be a good choice for implementing the parametric SURE-AMP algorithm to recover the k -dense signal.

5.3.2.3 Exponential kernel family

For the non-linear parameter of the exponential kernel, we adopt the recommendation in [123] and set T as

$$T = 6\sqrt{c} \quad (5.27)$$

It has been demonstrated through extensive simulations in [123] that the image denoising quality is not very sensitive to the ratio between T and \sqrt{c} . Eq. (5.27) is shown to be a practical setting for removing various noise perturbation irrespective of the images. The denoising and reconstruction simulations in the subsequent sections will also confirm that it is a plausible

choice for both sparse and heavy-tailed signals.

5.3.3 Linear parameter optimization

With the non-linear parameters fixed with the effective noise variance, the only parameters left to be optimized are the kernel weights a_i . Denote ε as the MSE of the denoised signal using the parametric function $f(\mathbf{r}, c|\boldsymbol{\theta})$. With Theorem 5, we have

$$\varepsilon = c + \langle g^2(\mathbf{r}, c|\boldsymbol{\theta}) + 2cg'(\mathbf{r}, c|\boldsymbol{\theta}) \rangle \quad (5.28)$$

where

$$\begin{aligned} g(\mathbf{r}, c|\boldsymbol{\theta}) &= f(\mathbf{r}, c|\boldsymbol{\theta}) - \mathbf{r} \\ &= \sum_{i=1}^k a_i f_i(\mathbf{r}|\boldsymbol{\vartheta}_i(c)) - \mathbf{r} \end{aligned} \quad (5.29)$$

Optimizing the weights a_i to achieve the minimum MSE requires differentiation of ε over a_i and solving for all $i \in (1, \dots, k)$.

$$\begin{aligned} \frac{d\varepsilon}{da_i} &= \langle 2g(\mathbf{r}, c|\boldsymbol{\theta}) \frac{d}{da_i} g(\mathbf{r}, c|\boldsymbol{\theta}) + 2c \frac{d}{da_i} g'(\mathbf{r}|\boldsymbol{\theta}) \rangle = 0 \\ \iff \sum_{j=1}^k \langle a_j f_j(\mathbf{r}|\boldsymbol{\vartheta}_j(c)) f_i(\mathbf{r}|\boldsymbol{\vartheta}_i(c)) \rangle &= -c \langle f'_i(\mathbf{r}|\boldsymbol{\vartheta}_i(c)) \rangle \end{aligned} \quad (5.30)$$

All equations can be summarized in the following matrix form

$$\underbrace{\begin{bmatrix} \langle f_1^2 \rangle & \cdots & \langle f_1 f_k \rangle \\ \vdots & \ddots & \vdots \\ \langle f_k f_1 \rangle & \cdots & \langle f_k^2 \rangle \end{bmatrix}}_{\mathcal{F}} \underbrace{\begin{bmatrix} a_1 \\ \vdots \\ a_k \end{bmatrix}}_{\mathcal{A}} = -c \underbrace{\begin{bmatrix} \langle f'_1 \rangle \\ \vdots \\ \langle f'_k \rangle \end{bmatrix}}_{\mathcal{D}} \quad (5.31)$$

The linear system can then be solve by

$$\mathcal{A} = -c\mathcal{F}^{-1}\mathcal{D} \quad (5.32)$$

With only two or three basis functions, \mathcal{F} is trivial to invert. In summary, the linear kernel weights can be easily optimized by solving a linear system of equations. We will demonstrate later that this linear parameterization is very advantageous in terms of the computational com-

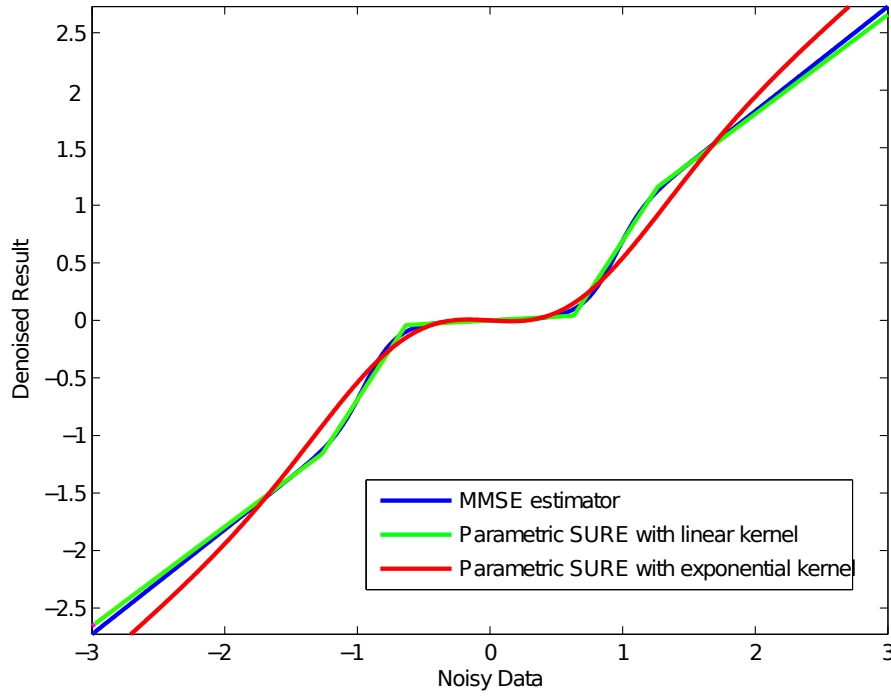


Figure 5.4: MMSE estimator and parametric SURE for the noisy Bernoulli-Gaussian data. The noise variance c is 0.1. The reconstruction error for the MMSE estimator, the SURE estimator with the first piecewise linear kernel and the SURE estimator with the exponential kernel are 0.020615, 0.020788 and 0.022047, respectively.

plexity. This aforementioned approach is in spirit similar to the SURE-LET algorithm in [124], only that the optimization is done recursively at each iteration for the parametric SURE-AMP algorithm.

5.3.4 Denoising performance

To validate our proposed kernel families and the parameter optimization scheme, we compare the optimized parametric denoisers alongside the MMSE estimator for BG and k -dense signals.

In Fig. 5.4 we can see that with just three kernel functions from the first piecewise linear kernel family and the suggested parameter optimization in (5.25), the constructed denoiser achieves an excellent agreement with the MMSE estimator for the noisy BG data. The MSE difference between the denoised signal using the SURE based parametric denoiser and the Bayesian optimal denoising is negligible. The exponential kernel family also does a good job at capturing the key structure of the MMSE estimator, especially in the vicinity of small values where most of the data concentrates.

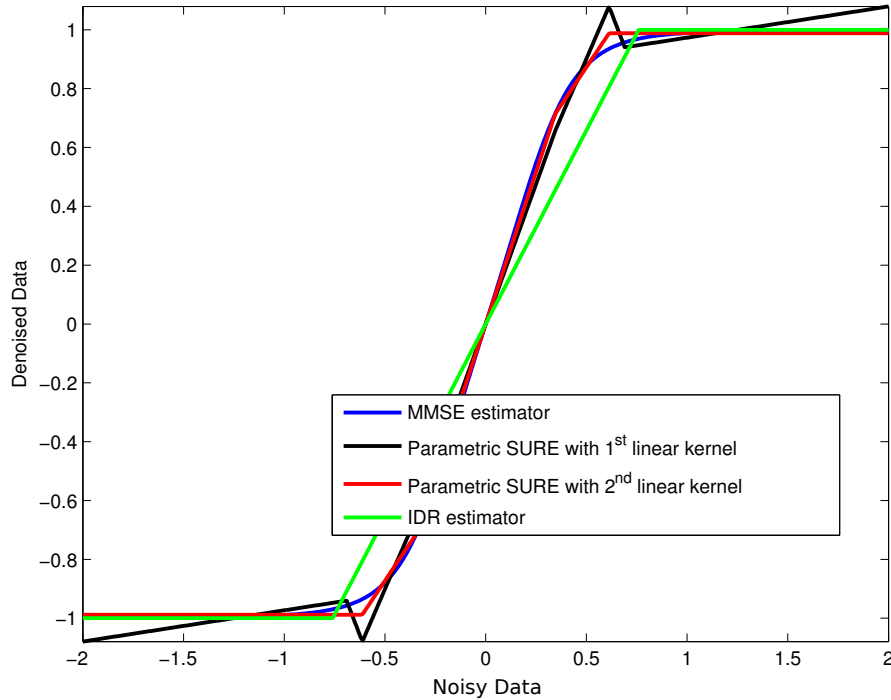


Figure 5.5: MMSE estimator and parametric SURE for the noisy k -dense data. The noise variance c is 0.1. The reconstruction error for the MMSE estimator, the SURE with the second piecewise linear kernel, the SURE estimator with the first piecewise linear kernel and the IDR denoiser are 0.0243 and 0.0248, 0.0251 and 0.0315 respectively.

In Fig. 5.5 we compare the MMSE estimator, the SURE based parametric denoisers with the proposed two piecewise linear kernel families, and the IDR estimator for the k -dense signal denoising. As demonstrated in the plot, the denoiser constructed with the second piecewise linear kernel fits the MMSE estimator better because the kernels are tailored to the k -dense structure. The first piecewise linear kernel based denoiser performs slightly worse because of the unbounded $f_3(\cdot)$ in eq. (5.20). The IDR denoiser is a piecewise linear function with just one hinge point. Thus it misses the subtle transition between the small and large elements and performs the worst among the three.

To check the denoising power of the proposed kernel families for heavy tailed signals, we present the averaged MSE for the Student's-t signal denoising in Table 5.1. Since there is not an explicit form for the MMSE estimator for the Student's-t prior, we compare the SURE based parametric denoiser with the GMD model based denoiser, which is the MMSE estimator for the 4-state GMD used to approximate the Student's-t distribution. It essentially is the key denoising approach implemented by the EM-GM-GAMP algorithm. Each figure reported in Table 5.1 is an average over 100 realizations with the signal length being 5000. The SURE based denoiser

Effective noise level c	0.01	0.1	1	5	10	50	100
MMSE estimator for 4-state GM	9.9655e-3	0.0958	0.7285	2.1788	3.2088	6.5801	8.6543
Exponential kernel denoiser	9.9948e-3	0.0967	0.7200	2.1504	3.1606	6.9979	9.6347
Piecewise linear kernel denoiser	9.9383e-3	0.0955	0.7191	2.1560	3.1764	6.6554	8.6245

Table 5.1: *Denoising comparison for noisy Student's-t signal with various denoisers*

with the exponential kernel and the first piecewise linear kernel both deliver similar denoising performance as the MMSE estimator for the 4-state GMD approximation, if not better. This implies that the corresponding parametric SURE-AMP algorithm should be competitive with the EM-GM-GAMP for the Student's-t signal reconstruction.

5.4 Numerical Results

In this section, the reconstruction performance and computational complexity of the parametric SURE-AMP algorithm, using the three types of kernel families introduced in Section 5.3, are compared with other CS reconstruction algorithms. In particular, we experiment with the BG, k -dense and Student's-t signals to demonstrate the reconstruction power and efficiency of the parametric SURE-AMP algorithm.

5.4.1 Noisy signal recovery

We first present the reconstruction quality for noisy signal recovery. For all simulations, we fixed the signal dimension to $n = 10000$. Each numerical point in the plots is an average of 100 Monte Carlo realizations. To have a fair comparison, the noise level is defined in the measurement domain and quantified as

$$SNR_y = 10 \log_{10} \frac{\|\Phi \mathbf{x}\|_2^2}{\|\boldsymbol{\xi}\|_2^2} \quad (5.33)$$

The reconstruction quality is evaluated in terms of the signal to noise ratio in the signal domain, defined as

$$SNR_x = 10 \log_{10} \frac{\|\mathbf{x}\|_2^2}{\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2} \quad (5.34)$$

The elements of the measurement matrix Φ are drawn i.i.d. from $\mathcal{N}(\Phi_{ij}; 0, m^{-1})$ and the matrix columns are normalized to one. For all reconstruction algorithms, the convergence tolerance is set as 10^{-6} . The maximum iteration number is set as 100.

5.4.1.1 Bernoulli-Gaussian prior

The BG data for the simulation are draw i.i.d. from

$$p(x) = 0.1\mathcal{N}(x; 0, 1) + 0.9\delta(x) \quad (5.35)$$

We choose the noise level to be $SNR_y = 25$ dB. For comparison, we show the performance of the parametric SURE-AMP algorithm with both first piecewise linear kernel and the exponential kernel family, the EM-BG-GAMP algorithm ¹, the ℓ_1 -AMP algorithm ² and the genie BAMP ³ algorithm. The reconstruction quality SNR_x for various sampling ratios are illustrated in Fig. 5.6.

It is obvious that the parametric SURE-AMP algorithm with the first piecewise linear kernel exhibits the near-optimal construction: for $\gamma \geq 0.24$, the difference between the parametric SURE-AMP algorithm which is blind to the signal prior and the genie BAMP algorithm is negligible. It also adequately demonstrates that SURE is a perfect surrogate for the MSE measure and the intrinsic signal property can be effectively exploited by the SURE-based denoiser. Moreover, it shows again that the proposed hinge point selection strategy in (5.25) works very well regardless of the effective noise level. Compared with the EM-BG-GAMP algorithm, it delivers roughly 2 dB better recovery for $0.24 \leq \gamma \leq 0.3$. This is probably because the EM-BG-GAMP gets stuck at local minima for smaller sampling ratio. For $\gamma > 0.36$, EM-BG-GAMP also delivers reconstruction performance that is very close to the genie BAMP result. It is because the kernels used in EM-BG-GAMP to fit the data are essentially the prior for generating the data. For the parametric SURE-AMP algorithm with the exponential kernels, it is roughly 1 dB worse than its counterpart with the first piecewise linear kernel and the Bayesian optimal reconstruction for $\gamma \geq 0.26$. This comes as no surprise as we have already seen in Fig. 5.4 that the denoiser based on exponential kernels doesn't capture the MMSE estimator structure for data with large magnitude. Nevertheless, it still demonstrates significant improvement over the ℓ_1 -AMP reconstruction for which no statistical property of the original signal is exploited.

¹A special case of the EM-GM-GAMP algorithm which approximates the signal prior with a mixture of Bernoulli and Gaussian distributions. We use the implementation from <http://www2.ece.ohio-state.edu/~vilaj/EMGMAMP/EMGMAMP.html>.

²We use the implementation from <http://people.epfl.ch/ulugbek.kamilov>.

³The true signal prior $p(\mathbf{x}_o)$ is assumed known for the BAMP reconstruction. It is served as the upper bound for SNR_x .

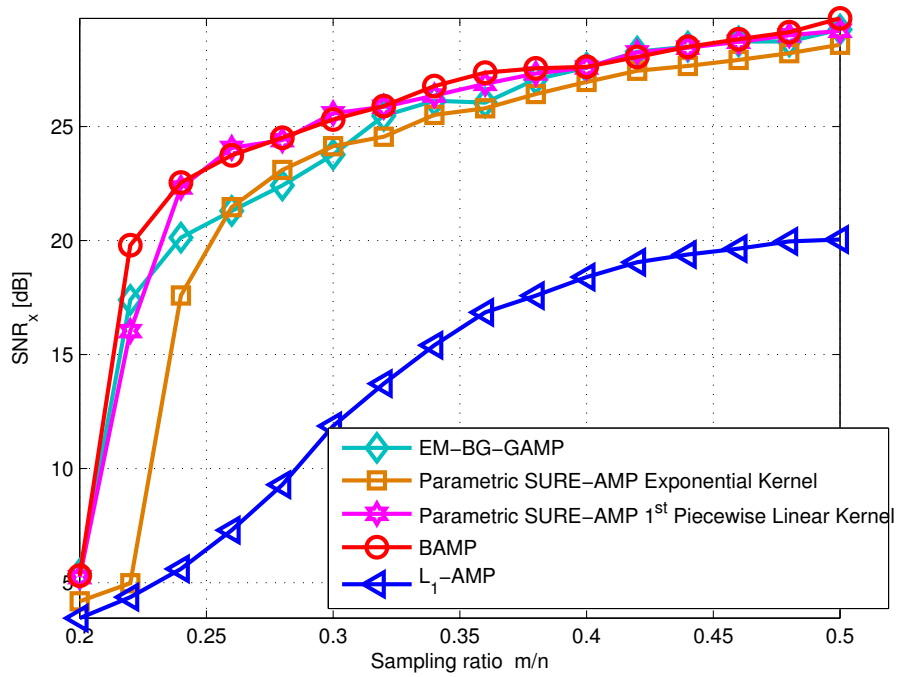


Figure 5.6: SNR_x versus sampling ratio for CS recovery of noisy Bernoulli-Gaussian data.

5.4.1.2 k -Dense signal

In [129], extensive simulations have been conducted to compare the IDR algorithm performance with the state-of-art algorithms for the noisy k -dense signal reconstruction. Thus in this chapter, we use the same setting and mainly compare the parametric SURE-AMP using two piecewise linear kernel families with the IDR, EM-GM-GAMP and the genie BAMP algorithm. The k -dense signal is generated i.i.d. from

$$p_{\text{KD}}(x) = 0.45\delta(x + 1) + 0.45\delta(x - 1) + 0.1\mathcal{U}(-1, 1) \quad (5.36)$$

The results are demonstrated in Fig. 5.7. The noise level is $SNR_y = 28$ as in [129]. For the EM-GM-GAMP algorithm we found that as the number of Gaussian components increase, the reconstruction quality gets better. Thus we used 20-state GMD to fit the k -dense prior, which is the largest number allowed for the EM-GM-GAMP MATLAB package. The parametric SURE-AMP with the second piecewise linear kernel is only slightly worse than the genie BAMP reconstruction. There is roughly 0.5 dB difference between the two for $\gamma > 0.5$. Comparing to the IDR reconstruction, there is a consistent 2 dB improvement for $\gamma \geq 0.55$. This reconstruction quality gain is predictable as we have already demonstrated in the denoising section in Fig. 5.5. It is also reasonable that the first piecewise linear kernel based parametric SURE-AMP does not perform as well as IDR and the second piecewise linear kernel. It is in general 2 dB

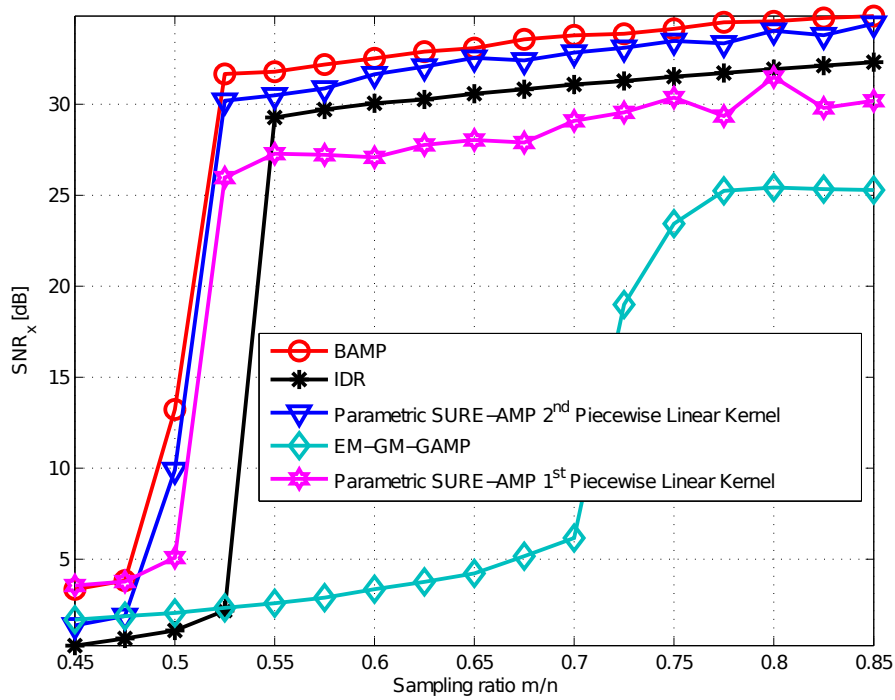


Figure 5.7: SNR_x versus sampling ratio for CS recovery of noisy k -dense data.

worse than the IDR and 5 dB worse than the genie BAMP bench mark. This is mainly because the first piecewise linear kernel fails to correct the large coefficients to be ± 1 . However, it still greatly outperforms the EM-GM-GAMP algorithm. The failure of the EM-GM-GAMP in this case is probably because the algorithm gets stuck at some local minima when fitting the prior. This example confirms the advantageous motivation for the parametric SURE-AMP algorithm: minimizing the MSE is the direct approach to obtain the best reconstruction.

5.4.1.3 Student-t prior

To investigate the parametric SURE-AMP performance for signals that are not strictly sparse, we consider the Student's-t prior as a heavily-tailed distribution example. The signal is draw i.i.d. according to the following distribution.

$$p_{\mathbf{T}}(x_o) = \frac{\Gamma((q+1)/2)}{\sqrt{q\pi}\Gamma(q/2)}(1+x_o^2)^{-(q+1)/2} \quad (5.37)$$

where q controls the distribution shape. It has been demonstrated in [8] that the Student's-t distribution is an excellent model to capture the statistical behaviour of the DCT coefficients for natural images. In the simulation, we set $q = 1.67$, $SNR_y = 25$ dB as in [54]. The parametric SURE-AMP using both exponential and the first piecewise linear kernel family are

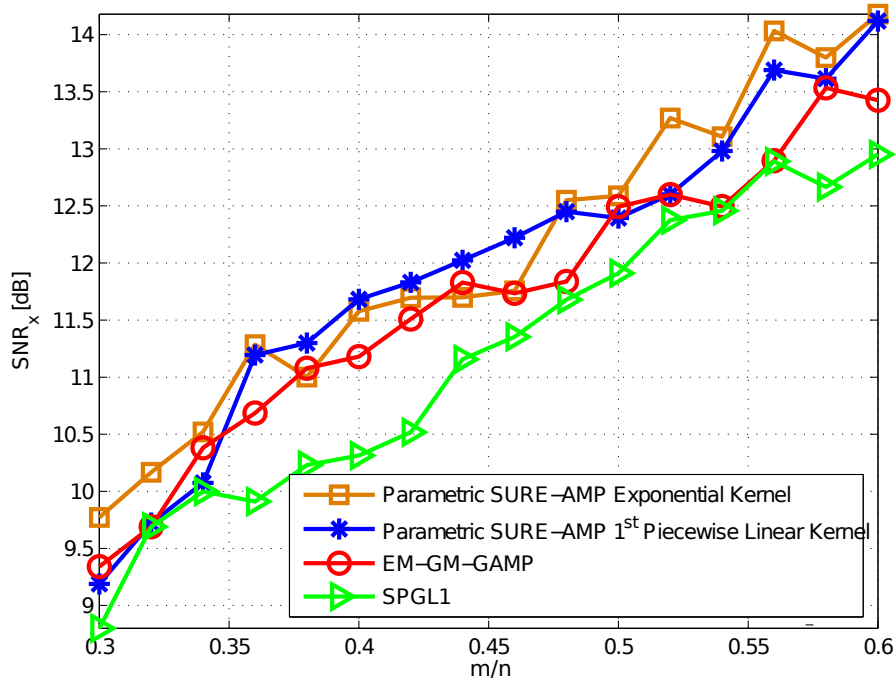


Figure 5.8: SNR_x versus sampling ratio for CS recovery of noisy student- t data.

compared with the EM-GM-GAMP algorithm and LASSO via SPGL1⁴ [130]. As we can see from Fig. 5.8, the parametric SURE-AMP and EM-GM-GAMP have the similar reconstruction performance. This can be expected from the denoising comparison in the previous section. None of them achieves significant improvement over the ℓ_1 -minimization approach though. This is probably because the signal is not very compressible [8]. The parametric SURE-AMP and EM-GM-GAMP algorithm are likely to be already near the Bayesian optimal performance.

5.4.2 Runtime comparison

The parametric SURE-AMP algorithm with the simple kernel functions and linear parameterization does not only achieve the near optimal reconstruction, more importantly, it significantly reduces the computational complexity. The authors in [54] have compared the EM-GM-GAMP algorithm with most of the existing CS algorithms that are blind of the prior and proved EM-GM-GAMP is the most efficient among them all. Thus in this section, we will use the EM-GM-GAMP runtime performance as the bench mark to evaluate the efficiency of the parametric SURE-AMP algorithm. For this, we fixed $\gamma = 0.5$, $SNR_y = 25$ dB and varied the signal length n from 10000 to 100000. For the EM-GM-GAMP algorithm, we set the EM tolerance

⁴We run the SPGL1 in the “BPDN” mode. The MATLAB package can be found in <http://www.cs.ubc.ca/labs/scl/spgl1>.

to 10^{-5} and the maximum EM iterations to 20. The runtime for noisy recovery of the BG, k -dense and Student's-t data are plotted in Fig. 5.9, Fig. 5.10 and Fig. 5.11 respectively. Every point in the plots is an average over 100 realizations. The algorithms tested here are the same as described before.

For both the EM-GM-GAMP and parametric SURE-AMP algorithm, there is the computational cost for the matrix multiplication of the vector with the measurement matrix Φ and Φ^T at each iteration. However, we observed a dramatic runtime improvement across all tested signal lengths for three signal priors. The parametric SURE-AMP is more than 20 times faster than the EM-GM-GAMP scheme. The algorithm efficiency can be attributed to the simple form of the kernel functions, the linear parameterization of the SURE-based denoiser and the reduced number of iterations. Consider the runtime comparison of the BG data reconstruction. The total number of the EM-BG-GAMP iterations is roughly twice as many as that of the parametric SURE-AMP algorithm. Moreover, the per-iteration computational cost is much more expensive for EM-BG-GAMP since fitting the signal prior requires many EM iterations. While for each parametric SURE-AMP iteration, only a 3 dimensional linear system needs to be solved to optimize the adaptive estimator. When compared with the ℓ_1 -AMP, the runtime for each parametric SURE-AMP iteration is approximately the same. The improved runtime performance here comes from the effective denoising so that fewer iterations are required for the parametric SURE-AMP to converge. The best runtime performance of the IDR algorithm for the k -dense data in Fig. 5.10 is understandable since it only applies an adaptive thresholding function at each iteration and has no parameter optimization procedure.

5.5 Conclusion

In this chapter, the parametric SURE-AMP is presented as a novel compressed sensing algorithm, which directly minimizes the MSE of the recovered signal at each iteration. Motivated from the fact that the AMP can be cast as an iterative Gaussian denoising algorithm, we propose to utilize the adaptive SURE based parametric denoiser within the AMP iteration. The optimization of the parameters is achieved by minimizing the SURE, which is an unbiased estimate of the MSE. More importantly, the minimization of SURE depends purely on the noisy observation, which in the large system limit fundamentally eliminates the need of the signal prior. This is also the first time that it has been employed for the CS reconstruction. The parametric SURE-AMP with the proposed three kernel families have demonstrated almost the same

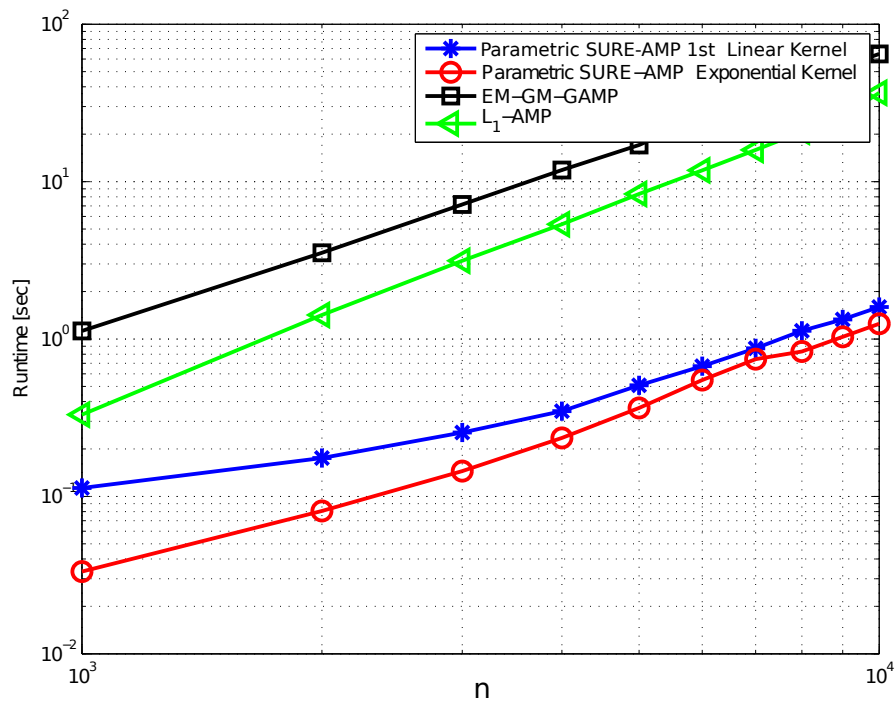


Figure 5.9: Runtime versus signal dimension for CS recovery of noisy Bernoulli-Gaussian data.

reconstruction quality as the BAMP algorithm, where the true signal prior is provided. It also outperforms the EM-GM-GAMP algorithm in terms of the computational cost. Directions for further research would involve considering other type of kernel families and the rigorous proof for the state evolution dynamics.

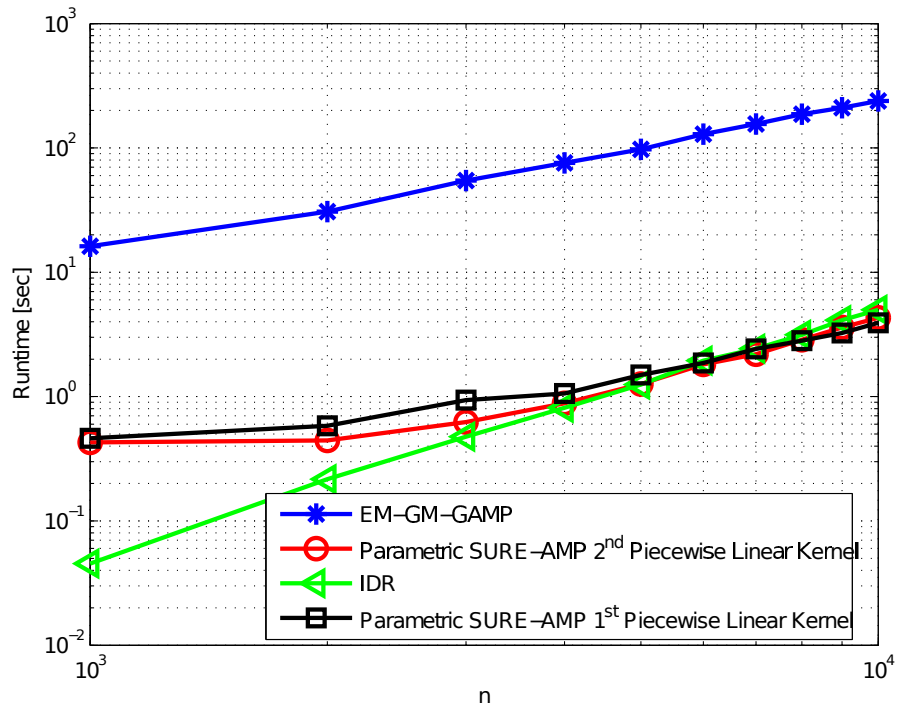


Figure 5.10: Runtime versus signal dimension for CS recovery of noisy k -dense data.

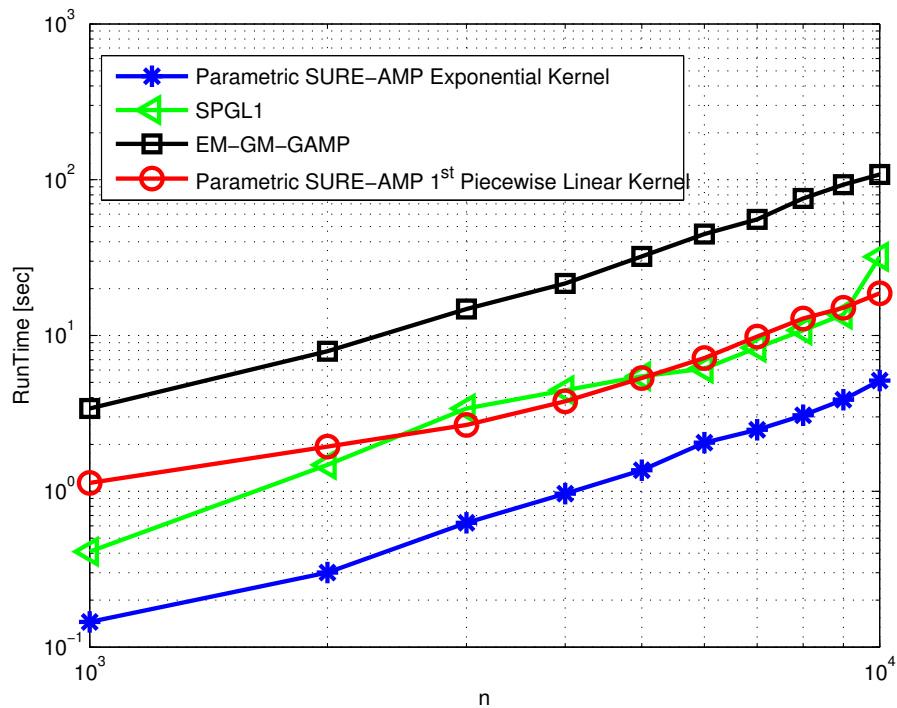


Figure 5.11: Runtime versus signal dimension for CS recovery of noisy student- t data.

Chapter 6

Conclusions and Future work

6.1 Conclusion

In this thesis we have investigated two aspects of the compressed sensing reconstruction: the measurement matrix design and enhancing one of the CS reconstruction algorithms. In Chapter 1, three main contributions made in this thesis were listed. Here we revisit these points and summarise the main advances and findings.

While the CS reconstruction power for natural images in the transformed domain is well recognized, there is very little literature discussing how to optimize the measurement matrix in a tractable manner and the performance bound for CS imaging. The main problem we addressed in Chapter 3 is that for the multi-resolution image model and the independent bandwise sampling strategy, what is the optimal sample allocation for the measurement matrix? To be specific, with a fixed number of samples, we would like to know how many samples should be allocated for each band to achieve the reconstruction with the least MSE. To this end, we proposed the sample distortion framework to quantify and bound the reconstruction MSE for any combination of measurement matrix and recovery algorithm at different sampling ratios. We subsequently derived the hybrid zeroing matrix to convexify the SD function and the greedy sample allocation based on the convexified SD function for the multi-resolution image model. The value of the study in Chapter 3 is twofolded. For one thing it confirmed the advantages of the structured measurement matrix over the homogeneous one for multi-resolution CS imaging. For another it theoretically quantifies and bounds the recovery performance of various reconstruction algorithms for natural images. The key insight is that with the optimized sample allocation, the reconstruction error decays at the same rate for both CS and the best linear reconstruction for the multi-resolution image model. Thus the reconstruction gain of CS algorithms over the linear techniques is fundamentally limited.

The structure of the measurement matrix can have a considerable impact on the CS reconstruction quality. In the CS literature, the spatially coupled measurement matrix was first proposed to push the perfect reconstruction sampling ratio as low as the sparsity level for sparse signals. Empirical simulation and theoretical analysis both confirm that with the specially designed

measurement matrix, it is possible to reach the theoretical limit for the exact recovery of sparse signals. While most of the work in the literature concentrates on achieving the perfect recovery, this thesis focused on designing the measurement matrix that improves the overall sample distortion performance for both sparse and compressible signals. The proposed modulated matrix, especially the two block design, can be seen as a special case of the spatially coupled seeded matrix. The key feature of the modulated matrix is that it has a much simpler 1-D state evolution dynamics to predict its asymptotic behaviour in the large system limit. This has a big impact on optimizing the matrix configuration. Compared to the seeded matrix, it has much fewer free parameters to tune. Moreover, this simple matrix construction does not necessarily degrade its reconstruction power. Extensive simulation for both sparse and compressible signals in Chapter 4 validates the proposal.

While utilizing the ℓ_1 -AMP and BAMP for signal reconstruction in this thesis, we noticed the performance gap between the two as well as the possibility to achieve the BAMP recovery without requiring the explicit signal prior in advance. Since the AMP based algorithms have the recursive Gaussian denoising nature, we proposed to utilize the SURE based parametric least square estimator family to deal with the signal denoising. The corresponding parametric SURE-AMP algorithm leverages the Gaussian behaviour of the AMP residual at each iteration and adaptively selects the denoiser in accordance with the effective Gaussian noise variance. The parametric SURE-AMP algorithm uses neither the ℓ_1 -minimization nor the Bayesian method, but employs the parametric approach with the Stein's unbiased risk estimate for reconstruction. With proper parameterization of the denoiser, we are able to fundamentally eliminate the need of the signal prior while still achieving the Bayesian optimal recovery.

6.2 Open Problem and Further Work

In this section we discuss some of the unanswered questions related to the work in this thesis and suggest several possible directions for further pursuit.

- *Extend Sample Allocation:*

It has been demonstrated in Chapter 3 that the combination of the sample allocation and the exploitation of the dependency across wavelet scales delivers the improved reconstruction for natural images. However, there is still a gap to be filled between the best achieved recovery and the theoretical MSE bound. This might be overcome by considering the truly optimized sample allocation for the wavelet tree structure or other more

sophisticated image models. Or it might be possible to approach the theoretical bound by allocating samples in an adaptive manner based on the reconstruction feedback from previous steps.

Another interesting research area is to extend the sample allocation scheme to practical imaging problems, e.g. magnetic resonance imaging (MRI) and computed tomography (CT). The difficulty in designing a practical measurement matrix is that we will not have all the freedom to allocate samples as in the theoretical analysis. For both MRI and CT, the data acquisition is performed in the Fourier domain rather than the wavelet domain. Since the wavelet bands typically do not have a finite k -space support, the analytical results in Chapter 3 cannot be applied directly for MRI and CT imaging. The physical constraints of medical devices also need to be taken into consideration for sampling pattern design. Despite all the difficulties, lots of heuristic CS sampling strategies for MRI indicate that concentrating most of the samples in the low frequency components while randomly sampling in the high frequencies benefits the overall construction. Thus systematically optimizing the sample allocation with constraints would have great practical value. It would also be interesting to extend the sample allocation scheme to 3-D to achieve acceleration for the MRI scanning.

- *Modulated Matrix Design:*

In this thesis we have demonstrated the advantageous reconstruction results using the modulated matrix for both sparse and compressible signals. However, only heuristic parameter settings for the two block matrix are presented without giving a parameter optimization scheme. Thus, one possible research direction would be to derive a systematic approach for optimizing the modulated matrix parameters. One could adopt the empirical strategy in [4] to obtain the optimal matrix setting, utilizing the state evolution dynamics to empirically test the reconstruction performance for various parameter combinations and select the best among them. To analytically optimize the modulated matrix configuration, new mathematical ideas need to be leveraged to exploit the state evolution equations.

Another open question for the two block matrix design is the lack of theoretical proof for the optimality of the hybrid zeroing matrix for compressible signals without first order phase transition. We believe that analysing the state evolution dynamics might lead to a further breakthrough. Moreover, in this thesis we only considered the combination of two direct deltas as the rescaling distribution for the modulated matrix. More sophisticated

distribution and parameter settings also need to be considered to obtain the improved reconstruction performance.

- *Parametric SURE-AMP Algorithm:*

As we have pointed out in Chapter 5, the empirical Monte Carlo simulation exhibits excellent agreement with the proposed state evolution equations for the parametric SURE-AMP algorithm. The next problem is to conduct the theoretical analysis to validate the SE prediction for the asymptotic behaviour. To obtain a rigorous proof for the SE dynamics, one might be able to leverage the technique used in [73] since they also considered an adaptation module within the AMP framework. Alternatively, one might find inspiration from the original state evolution analysis for AMP in [55]. The other unsolved problem for the parametric SURE-AMP algorithm is the theoretical validation of the greedy approach for jointly optimizing the denoisers across all iterations. We believe that the proof would rely on analysing the SE dynamics.

In Chapter 5, we only considered denoisers constructed by the exemplary piecewise linear and exponential kernel families with no more than three kernel functions. More sophisticated and general kernel format can be exploited to increase the flexibility and accuracy of the denoising functions. Different forms of parameterization of the denoisers could also be deployed. Additionally, other off-the-shelf denoising algorithms would be utilized to enhance the AMP reconstruction as suggested in [51].

The comparison of the parametric SURE-AMP and the EM-GM-GAMP algorithm is conducted in terms of the reconstruction quality and the computational complexity. The fundamental difference between the two is that one resorts to minimize the reconstruction MSE while the other concentrates on minimizing the Kullback-Leibler divergence to fit the signal prior. Although we noticed the differences, we did not address the intrinsic relationship between the two approaches. One possible research direction would be to understand under what conditions or with what form of kernel functions fitting the prior would be equivalent to optimizing the reconstruction MSE.

- *Enhancing AMP Algorithm:*

While in this thesis the AMP based algorithms are employed as the major CS reconstruction tool, there is some fundamental work to be done to extend the AMP framework to a more general problem setting. Currently the theoretical analysis for AMP relies heavily on the Gaussian assumption of the measurement matrix, which can be restrictive when applied to practical problems. Although in [14] different matrix ensembles including the

Rademacher, partial Fourier and USE (columns i.i.d uniformly distributed on the unit sphere) have been shown to follow the SE analysis of AMP, ideally we would like AMP to work for arbitrary matrices, i.e. the complex measurement matrix or the Fourier matrix with non-uniform sample allocation, to suit practical applications. This involves a better understanding of the generic message passing approximation for the graphic representation of the CS system, especially the derivation of the “Onsager” reaction term.

In Chapter 3 the signal sparsity in the wavelet domain has been utilized to enhance the AMP reconstruction. The graphic model for the CS observation with the tree structure prior is presented. The corresponding turbo reconstruction scheme is successful in exploiting both the signal prior and the sparse property. It might be beneficial if more sparse models on various orthogonal bases can be employed to aid the AMP reconstruction, i.e. using two types of wavelet or exploring in both wavelet and Fourier domain. Thus one possible avenue of research is to generate the graphic model for the CS measurement with more than one sparse representation and derive the appropriate message passing for the new system. The potential problem involved would be how to incorporate the estimate from different sparse bases to obtain a better reconstruction.

Appendix A

Approximation of the factor-to-variable Message

In this appendix, the detailed steps of approximating the factor-to-variable message $m_{j \rightarrow i}(x_i)$ in (2.33) with (2.37) is provided. First, from the Hubbard-Stratonovich transform we have

$$e^{-\frac{x^2}{2a}} = \sqrt{\frac{1}{2\pi a}} \int e^{-\frac{t^2}{2a} + \frac{-\sqrt{-1}xt}{a}} dt \quad (\text{A.1})$$

Let $x = \sum_{q \neq i} \Phi_{jq} x_q$, applying (A.1) we have

$$e^{-\frac{(\sum_{q \neq i} \Phi_{jq} x_q)^2}{2\sigma_\xi^2}} = \frac{1}{\sqrt{2\pi\sigma_\xi^2}} \int dt e^{-\frac{t^2}{2\sigma_\xi^2} + \frac{\sqrt{-1}t \sum_{q \neq i} \Phi_{jq} x_q}{\sigma_\xi^2}} \quad (\text{A.2})$$

From (2.33) we have

$$m_{j \rightarrow i}(x_i) = \frac{1}{Z_{j \rightarrow i}} e^{-\frac{(y_j - \Phi_{ji} x_i)^2}{2\sigma_\xi^2}} \int \prod_{q \neq i} dx_q e^{-\frac{(\sum_{q \neq i} \Phi_{jq} x_q)^2}{2\sigma_\xi^2}} e^{\frac{(y_j - \Phi_{ji} x_i) \sum_{q \neq i} \Phi_{jq} x_q}{\sigma_\xi^2}} m_{q \rightarrow j}(x_q) \quad (\text{A.3})$$

Combining (A.2) with (A.3) we arrive at

$$m_{j \rightarrow i}(x_i) = \frac{1}{Z_{j \rightarrow i} \sqrt{2\pi\sigma_\xi^2}} e^{-\frac{(y_j - \Phi_{ji} x_i)^2}{2\sigma_\xi^2}} \int dt e^{-\frac{t^2}{2\sigma_\xi^2}} \prod_{q \neq i} \left[\int dx_q m_{q \rightarrow j}(x_q) e^{\frac{\Phi_{jq} x_q}{\sigma_\xi^2} (y_j - \Phi_{ji} x_i + it)} \right] \quad (\text{A.4})$$

We then expand the last exponential in (A.4) around zero to further simplify the message $m_{j \rightarrow i}(x_i)$.

$$\begin{aligned}
 & \int dx_q \quad m_{q \rightarrow j}(x_q) e^{\frac{\Phi_{jq} x_q}{\sigma_\xi^2} (y_j - \Phi_{ji} x_i + it)} \\
 \stackrel{(a)}{=} & \int dx_q \quad m_{q \rightarrow j}(x_q) \left[1 + \frac{\Phi_{jq} x_q}{\sigma_\xi^2} (y_j - \Phi_{ji} x_i + it) + \frac{\Phi_{jq}^2 x_q^2}{2\sigma_\xi^4} (y_j - \Phi_{ji} x_i + it)^2 \right] \\
 = & 1 + \frac{\Phi_{jq}}{\sigma_\xi^2} (y_j - \Phi_{ji} x_i + it) \int dx_q \quad x_q m_{q \rightarrow j}(x_q) + \frac{\Phi_{jq}^2}{2\sigma_\xi^4} (y_j - \Phi_{ji} x_i + it)^2 \int dx_q \quad x_q^2 m_{q \rightarrow j}(x_q) \\
 \stackrel{(b)}{=} & 1 + \frac{\Phi_{jq} \alpha_{q \rightarrow j}}{\sigma_\xi^2} (y_j - \Phi_{ji} x_i + it) + \frac{\Phi_{jq}^2}{2\sigma_\xi^4} (y_j - \Phi_{ji} x_i + it)^2 (\nu_{q \rightarrow j} + \alpha_{q \rightarrow j}^2) \\
 \stackrel{(c)}{=} & 1 + \frac{\Phi_{jq} \alpha_{q \rightarrow j}}{\sigma_\xi^2} (y_j - \Phi_{ji} x_i + it) + \frac{\Phi_{jq}^2}{2\sigma_\xi^4} (y_j - \Phi_{ji} x_i + it)^2 \nu_{q \rightarrow j} \\
 & + \frac{1}{2} \left[\frac{\Phi_{jq} \alpha_{q \rightarrow j}}{\sigma_\xi^2} (y_j - \Phi_{ji} x_i + it) + \frac{\Phi_{jq}^2}{2\sigma_\xi^4} (y_j - \Phi_{ji} x_i + it)^2 \nu_{q \rightarrow j} \right]^2 \\
 \stackrel{(d)}{=} & e^{\frac{\Phi_{jq} \alpha_{q \rightarrow j}}{\sigma_\xi^2} (y_j - \Phi_{ji} x_i + it) + \frac{\Phi_{jq}^2 \nu_{q \rightarrow j}}{2\sigma_\xi^4} (y_j - \Phi_{ji} x_i + it)^2}
 \end{aligned} \tag{A.5}$$

In the above derivation, we use the following assumption and observations

- (a) The Taylor expansion of the exponential around zero $e^x \approx 1 + x + \frac{x^2}{2}$
- (b) The mean and variance of the message $m_{q \rightarrow j}(x_q)$ are defined as $\alpha_{q \rightarrow j}, \nu_{q \rightarrow j}$ in (2.35) and (2.36) respectively as the new messages.
- (c) We assume all terms that are above the order $\mathcal{O}(1/m)$ are neglectable, thus $\Phi_{ji}^3 \approx \Phi_{ji}^4 \approx 0$.
- (d) The inverse of the Taylor expansion of the exponential around zero $1 + x + \frac{x^2}{2} \approx e^x$

Apply (A.5) to (A.4) we have

$$m_{j \rightarrow i}(x_i) = \frac{1}{Z^{j \rightarrow i} \sqrt{2\pi\sigma_\xi^2}} e^{-\frac{(y_j - \Phi_{ji} x_i)^2}{2\sigma_\xi^2}} \int dt \quad e^{-\frac{t^2}{2\sigma_\xi^2}} \prod_{q \neq i} \left[e^{\frac{\Phi_{jq} \alpha_{q \rightarrow j}}{\sigma_\xi^2} (y_j - \Phi_{ji} x_i + it) + \frac{\Phi_{jq}^2 \nu_{q \rightarrow j}}{2\sigma_\xi^4} (y_j - \Phi_{ji} x_i + it)^2} \right] \tag{A.6}$$

We then perform the Gaussian integral over t in (A.6). Define the following notation

$$A = \sum_{q \neq i} \Phi_{jq}^2 \nu_{q \rightarrow j}, \quad B = \sum_{q \neq i} \Phi_{jq} \alpha_{q \rightarrow j}, \quad N = y_j - \Phi_{ji} x_i \tag{A.7}$$

Then we have

$$\begin{aligned}
 & \frac{1}{\sqrt{2\pi\sigma_\xi^2}} e^{-\frac{(y_j - \Phi_{ji}x_i)^2}{2\sigma_\xi^2}} \int dt e^{-\frac{t^2}{2\sigma_\xi^2}} \prod_{q \neq i} \left[e^{\frac{\Phi_{jq}\alpha_{q \rightarrow j}}{\sigma_\xi^2} (y_j - \Phi_{ji}x_i + it) + \frac{\Phi_{jq}^2 \nu_{q \rightarrow j}}{2\sigma_\xi^4} (y_j - \Phi_{ji}x_i + it)^2} \right] \\
 &= \frac{1}{\sqrt{2\pi\sigma_\xi^2}} e^{-\frac{N^2}{2\sigma_\xi^2} e^{\frac{NB}{\sigma_\xi^2} + \frac{N^2A}{2\sigma_\xi^4}}} \int dt e^{-\left(\frac{1}{2\sigma_\xi^2} + \frac{A}{2\sigma_\xi^4}\right)t^2 + \left(\frac{B}{\sigma_\xi^2} + \frac{AN}{\sigma_\xi^4}\right)it} \\
 &\stackrel{(a)}{=} \sqrt{\frac{\sigma_\xi^2}{A + \sigma_\xi^2}} e^{-\frac{N^2}{2\sigma_\xi^2} e^{\frac{NB}{\sigma_\xi^2} + \frac{N^2A}{2\sigma_\xi^4}}} e^{\frac{(B\sigma_\xi^2 + AN)^2}{2\sigma_\xi^4(A + \sigma_\xi^2)}} \\
 &= \sqrt{\frac{\sigma_\xi^2}{A + \sigma_\xi^2}} e^{\frac{B^2}{2(A + \sigma_\xi^2)}} e^{\frac{2A^2 - \sigma_\xi^4}{2\sigma_\xi^4(A + \sigma_\xi^2)} N^2 + \frac{B(2A + \sigma_\xi^2)}{\sigma_\xi^2(A + \sigma_\xi^2)} N} \\
 &\stackrel{(b)}{=} \sqrt{\frac{\sigma_\xi^2}{A + \sigma_\xi^2}} e^{\frac{B^2}{2(A + \sigma_\xi^2)}} e^{-\frac{1}{2(A + \sigma_\xi^2)} N^2 + \frac{B}{A + \sigma_\xi^2} N} \\
 &\stackrel{(c)}{=} \sqrt{\frac{\sigma_\xi^2}{A + \sigma_\xi^2}} e^{\frac{B^2 - y_j^2 + 2By_j}{2(A + \sigma_\xi^2)}} e^{-\frac{\Phi_{ji}^2}{2(A + \sigma_\xi^2)} x_i^2} e^{\frac{\Phi_{ji}(y_j - B)}{A + \sigma_\xi^2} x_i}
 \end{aligned} \tag{A.8}$$

where we have used the following observations

- (a) Gaussian integral $\int dx e^{-ax^2 - 2bx} = \sqrt{\frac{\pi}{a}} e^{\frac{b^2}{a}}$.
- (b) We omit the terms that of the order $\mathcal{O}(1/m)$ and above, thus $AB \approx A^2 \approx 0$.
- (c) The pre-defined notation $N = y_j - \Phi_{ji}x_i$.

Finally, we define the following notation

$$A_{j \rightarrow i} = \frac{\Phi_{ji}^2}{\sigma_\xi^2 + \sum_{q \rightarrow i} \Phi_{jq}^2 \nu_{q \rightarrow j}} \tag{A.9}$$

$$B_{j \rightarrow i} = \frac{\Phi_{ji}(y_j - \sum_{q \neq i} \Phi_{jq} \alpha_{q \rightarrow j})}{\sigma_\xi^2 + \sum_{q \neq i} \Phi_{jq}^2 \nu_{q \rightarrow j}} \tag{A.10}$$

together with (A.8) we have the simplified version of the message $m_{j \rightarrow i}(x_i)$.

$$m_{j \rightarrow i}(x_i) = \frac{1}{\tilde{Z}_{j \rightarrow i}} e^{-\frac{A_{j \rightarrow i}}{2} x_i^2 + B_{j \rightarrow i} x_i} \quad \tilde{Z}_{j \rightarrow i} = \sqrt{\frac{2\pi}{A_{j \rightarrow i}}} e^{\frac{B_{j \rightarrow i}^2}{2A_{j \rightarrow i}}} \tag{A.11}$$

where $\tilde{Z}_{j \rightarrow i}$ is the normalization factor containing all x_i -independent terms.

Appendix B

Derivation of Equation (2.53)

With (2.44) we have

$$\begin{aligned}
\alpha_{i \rightarrow j} &= f_a\left(\frac{\sum_{p \neq j} B_{p \rightarrow i}}{\sum_{p \neq j} A_{p \rightarrow i}}, \frac{1}{\sum_{p \neq j} A_{p \rightarrow i}}\right) \\
&\stackrel{(a)}{=} f_a(R_i, \Sigma_i^2) + \left(\frac{\sum_{p \neq j} B_{p \rightarrow i}}{\sum_{p \neq j} A_{p \rightarrow i}} - \frac{\sum_p B_{p \rightarrow i}}{\sum_p A_{p \rightarrow i}}\right) \frac{\partial f_a}{\partial \varepsilon}(\varepsilon_i, c_i) \\
&\quad + \left(\frac{1}{\sum_{p \neq j} A_{p \rightarrow i}} - \frac{1}{\sum_p A_{p \rightarrow i}}\right) \frac{\partial f_a}{\partial \varepsilon}(\varepsilon_i, c_i) \\
&\stackrel{(b)}{=} \alpha_i + \frac{\sum_p A_{p \rightarrow i} (\sum_p B_{p \rightarrow i} - B_{j \rightarrow i}) - \sum_p B_{p \rightarrow i} (\sum_p A_{p \rightarrow i} - A_{j \rightarrow i})}{\sum_p A_{p \rightarrow i} \sum_{p \neq j} A_{p \rightarrow i}} \frac{\partial f_a}{\partial \varepsilon}(\varepsilon_i, c_i) \\
&= \alpha_i + \frac{A_{j \rightarrow i} (\sum_{p \neq j} B_{p \rightarrow i} + B_{j \rightarrow i}) - B_{j \rightarrow i} (\sum_{p \neq j} A_{p \rightarrow i} + A_{j \rightarrow i})}{\sum_p A_{p \rightarrow i} \sum_{p \neq j} A_{p \rightarrow i}} \frac{\partial f_a}{\partial \varepsilon}(\varepsilon_i, c_i) \\
&\stackrel{(c)}{=} \alpha_i - B_{j \rightarrow i} \Sigma_i^2 \frac{\partial f_a}{\partial \varepsilon}(\varepsilon_i, c_i) \\
&\stackrel{(d)}{=} \alpha_i - B_{j \rightarrow i} f_c(\varepsilon_i, c_i) = \alpha_i - B_{j \rightarrow i} \nu_i
\end{aligned} \tag{B.1}$$

where we have used the following observations

- (a) We use the Taylor expansion for the two-variable function

$$f(x, y) = f(a, b) + (x - a) \frac{\partial f}{\partial x}(a, b) + (y - b) \frac{\partial f}{\partial y}(a, b) \tag{B.2}$$

- (b) Omit the correction term $\frac{A_{j \rightarrow i}}{\sum_p A_{p \rightarrow i} \sum_{p \neq j} A_{p \rightarrow i}} \approx 0$.
- (c) Omit the correction term $\frac{A_{j \rightarrow i} \sum_{p \neq j} B_{p \rightarrow i}}{\sum_p A_{p \rightarrow i} \sum_{p \neq j} A_{p \rightarrow i}} \approx 0$.
- (d) With the definition in (2.42) and (2.43), we have the following relationship

$$f_c(\varepsilon, c) = c \frac{d}{d\varepsilon} f_a(\varepsilon, c) \tag{B.3}$$

The proof is provided in Appendix C.

Appendix C

Relationship of Conditional Mean and Variance for Gaussian Corrupted Data

Theorem 6. *Let x be the signal that we are interested, $r = x + w$ is the Gaussian noise corrupted data with $w \sim \mathcal{N}(w; 0, \Sigma^2)$. Then the MMSE estimator $F(r, \Sigma^2) = \mathbb{E}(x|r)$ and the conditional variance $G(r, \Sigma^2) = \text{Var}(x|r)$ has the following relationship*

$$G(r, \Sigma^2) = \Sigma^2 F'(r, \Sigma^2) \quad (\text{C.1})$$

where F' is the derivative of F with respect to r .

Proof. Given the definition of the conditional mean and variance

$$F(r, \Sigma^2) = \frac{1}{p(r)} \int xp(r|x)p(x)dx \quad (\text{C.2})$$

$$G(r, \Sigma^2) = \frac{1}{p(r)} \int x^2p(r|x)p(x)dx - F^2(r, \Sigma^2) \quad (\text{C.3})$$

we have

$$\begin{aligned} \Sigma^2 F'(r, \Sigma^2) &= \Sigma^2 \frac{dF}{dr} \\ &= \frac{\Sigma^2}{p(r)} \int xp'(r|x)p(x)dx - \frac{\Sigma^2 p'(r)}{p(r)} \int \frac{xp(r|x)p(x)}{p(r)} dx \end{aligned} \quad (\text{C.4})$$

Since $p(r|x) = \mathcal{N}(r; x, \Sigma^2)$, we have

$$p'(r|x) = p(r|x) \frac{x-r}{\Sigma^2} \quad (\text{C.5})$$

From Miysawa [127] we know that the least squares estimator can be written entirely in terms of the measurement density.

$$F(r, \Sigma^2) = r + \Sigma^2 \frac{p'(r)}{p(r)} \quad (\text{C.6})$$

Combining (C.4) with (C.5) and (C.6) we have

$$\begin{aligned}\Sigma^2 F'(r, \Sigma^2) &= \frac{1}{p(r)} \int xp(r|x)(x-r)p(x)dx - (F(r, \Sigma^2) - r) \int \frac{xp(r|x)p(x)}{p(r)} dx \\ &= \frac{1}{p(r)} \int x^2p(r|x)p(x)dx - F^2(r, \Sigma^2) \\ &= G(r, \Sigma^2)\end{aligned}\tag{C.7}$$

which completes the proof. □

Appendix D

Proof of the Entropy Based Bound

Proof:

Without loss of generality we will assume that Φ is an orthogonal projection operator and we denote by Φ^\perp the orthogonal projection onto the null space of Φ . We can then split the signal \mathbf{x} into its observed and unobserved components: $\mathbf{y} = \Phi\mathbf{x}$ and $\mathbf{z} = \Phi^\perp\mathbf{x}$. Since we directly observe \mathbf{y} we need only consider the component of the decoder that estimates \mathbf{z} , $\Delta^{(\mathbf{z})} : \mathbb{R}^m \rightarrow \mathbb{R}^{n-m}$. We can then estimate \mathbf{x} as:

$$\hat{\mathbf{x}} = \Delta(\mathbf{y}) = \Phi^T\mathbf{y} + [\Phi^\perp]^T \Delta^{(\mathbf{z})}(\mathbf{y}) \quad (\text{D.1})$$

We can further write the squared error distortion in terms of $\Delta^{(\mathbf{z})}(\mathbf{y})$ as

$$D = \frac{1}{n} \int p(\mathbf{y}) \int p(\mathbf{z}|\mathbf{y}) \|\mathbf{z} - \Delta^{(\mathbf{z})}(\mathbf{y})\|_2^2 d\mathbf{z}d\mathbf{y} \quad (\text{D.2})$$

Now consider the following decomposition of the differential entropy $h(\mathbf{x})$ of the vector \mathbf{x} :

$$\begin{aligned} h(\mathbf{x}) &= h(\mathbf{y}) + h(\mathbf{z}|\mathbf{y}) \\ &= h(\mathbf{y}) + h(\mathbf{z} - \Delta^{(\mathbf{z})}(\mathbf{y})|\mathbf{y}) \\ &\leq h(\mathbf{y}) + h(\mathbf{z} - \Delta^{(\mathbf{z})}(\mathbf{y})) \\ &\leq \frac{m}{2} \log_2 2\pi e + \frac{n-m}{2} \log_2 2\pi e n D / (n-m) \end{aligned} \quad (\text{D.3})$$

where we have used the following observations

- (line 2) The decoder is a deterministic function of \mathbf{y} and therefore the differential entropy of $h(\mathbf{z} - \Delta^{(\mathbf{z})}(\mathbf{y})|\mathbf{y}) = h(\mathbf{z}|\mathbf{y})$.
- (line 3) The conditional entropy is bounded by the marginal entropy: $h(\mathbf{x}|\mathbf{y}) \leq h(\mathbf{x})$.

- (line 4) The entropy of a random variable with a fixed covariance is bounded by the entropy of a Gaussian with the same covariance. Similarly the entropy of a random vector $\mathbf{v} \in \mathbb{R}^{n-m}$ under the constraint that $\mathbb{E}\{\mathbf{v}^T \mathbf{v}\} = nD$ is bounded by the entropy of a Gaussian random vector with covariance $\frac{nD}{(n-m)I}$.

The principle here is that the optimal projection should maximize the entropy of the observed component $h(\mathbf{y})$ while the decoder, $\Delta(\mathbf{y})$, should minimize the distortion possible. This is similar to the concept of information sensing proposed in [86].

Substituting $\gamma = m/n$ into (D.3) gives:

$$h(x) \leq \frac{1-\gamma}{2} \log_2 2\pi e \frac{D}{1-\gamma} + \frac{\gamma}{2} \log_2 2\pi e \quad (\text{D.4})$$

where we have used the i.i.d assumption to write $h(\mathbf{x}) = nh(x)$. This can then be rearranged to give the EBB.

Appendix E

Derivation of the Hierarchical Bayesian Model for GGD

Proof:

Here, we derive the hierarchical Bayesian model to describe the GGD, which is then used to bound the MSE performance described in the main text in Sec. 3.2.2.2. We introduce two latent variables c_1 and c_2 to simplify the expression of GGD:

$$c_1 = \frac{\alpha}{2\sqrt{\beta}\sigma\Gamma(\frac{1}{\alpha})} \quad c_2 = (\sqrt{\beta}\sigma)^\alpha \quad (\text{E.1})$$

Then the pdf of GGD can be written as

$$p_{\text{GGD}}(x) = c_1 \exp\left(-\frac{|x|^\alpha}{c_2}\right) \quad (\text{E.2})$$

Let $p(x|\tau) = \mathcal{N}(x; 0, \tau)$. To establish the hierarchical model, we need to find the prior $p(\tau)$ which satisfies:

$$\int_0^\infty \mathcal{N}(x; 0, \tau) p(\tau) d\tau = c_1 \exp\left(-\frac{|x|^\alpha}{c_2}\right) \quad (\text{E.3})$$

Using the substitution $g(\tau) = \frac{1}{\sqrt{2\pi\tau}} p(\tau)$, $m = \frac{x^2}{2}$ and $t = \frac{\sqrt{2}^\alpha}{c_2}$, the question becomes solving $g(\tau)$ subject to

$$\int_0^\infty \exp\left(-\frac{m}{\tau}\right) g(\tau) d\tau = c_1 \exp\left(-tm^{\frac{\alpha}{2}}\right) \quad (\text{E.4})$$

let $z = \frac{1}{\tau}$ and $G(z) = g(\tau)|_{\tau=\frac{1}{z}}$, we further transform the problem to find $G(z)$ subject to

$$\int_0^\infty \exp(-zm) \frac{G(z)}{z^2} dz = c_1 \exp\left(-tm^{\frac{\alpha}{2}}\right) \quad (\text{E.5})$$

Applying the integral formula [131]: if $\int_0^\infty e^{-zt} y(t) dt = f(z)$, then $y(t) = \mathcal{L}^{-1}(f(z))$, we obtain

$$\frac{G(z)}{z^2} = c_1 \mathcal{L}^{-1} \exp\left(-\frac{x^\alpha}{c_2}\right) \quad (\text{E.6})$$

where $\mathcal{L}^{-1}(\cdot)$ is the inverse Laplace transform. The inversion of Laplace transform in (E.6) can be solved numerically [132]. From here we obtain the MBB for the GGD data in Fig. 3.2.

Bibliography

- [1] D. Donoho, “High-dimensional centrally symmetric polytopes with neighborliness proportional to dimension,” *Discrete & Computational Geometry*, vol. 35, no. 4, pp. 617–652, 2006.
- [2] D. Donoho and J. Tanner, “Neighborliness of randomly-projected simplices in high dimensions,” in *Proceedings of the National Academy of Sciences*, vol. 102, no. 27, 2005, pp. 9452–9457.
- [3] D. Martin, C. Fowlkes, D. Tal, and J. Malik, “A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics,” in *Proc. 8th Int. Conf. on Computer Vision (ICCV)*, vol. 2, 2001, pp. 416–423.
- [4] F. Krzakala, M. Mézard, F. Sausset, Y. Sun, and L. Zdeborová, “Statistical-physics-based reconstruction in compressed sensing,” *Phys. Rev. X*, vol. 2, pp. 021 005(1–18), May 2012. [Online]. Available: <http://link.aps.org/doi/10.1103/PhysRevX.2.021005>
- [5] E. Candes, “compressed sensing makes every pixel count,” *what’s happening in the mathematical sciences*, vol. 7, pp. 114–127, 2009.
- [6] A. Gilbert, M. Iwen, and M. Strauss, “Group testing and sparse signal recovery,” in *42nd Asilomar Conference on Signals, Systems and Computers*, Oct 2008, pp. 1059–1063.
- [7] E. Candes and M. Wakin, “An introduction to compressive sampling,” *IEEE Signal Processing Magazine*, vol. 25, pp. 21–30, 2008.
- [8] R. Gribonval, V. Cehver, and M. Davies, “Compressible distributions for high dimensional statistics,” *IEEE Transactions on Information Theory*, vol. 58, no. 8, pp. 5016–5034, Aug 2012.
- [9] H. Choi and R. Baraniuk, “Wavelet statistical models and Besov spaces,” in *Proc. of SPIE Tech. Conf. on Wavelet Applicat. in Signal Process. VII.* Denver, Colo., Jul. 1999, pp. 489–501.
- [10] S. Mallat, “A theory for multiresolution signal decomposition: the wavelet representation,” *IEEE Trans. on Pattern Anal. Mach. Intell.*, vol. 11, no. 7, pp. 674–693, Jul 1989.

-
- [11] V. Cevher, "Learning with compressible priors," in *Proc. Neural Information Processing Systems*, 2008.
- [12] G. Kramer, I. Marić, and R. Yates, "Cooperative communications," *Foundations and Trends in Networking*, vol. 1, no. 3-4, pp. 271–425, June 2007.
- [13] D. Donoho and J. Tanner, "Counting faces of randomly projected polytopes when the projection radically lowers dimension," *Journal of the American Mathematical Society*, vol. 22, no. 1, pp. 1–53, 2009.
- [14] D. Donoho, A. Maleki, and A. Montanari, "Message-passing algorithms for compressed sensing," *Proc. of the Nat. Academy of Sciences*, vol. 106, no. 45, pp. 18 914–18 919, 2009.
- [15] I. Maravić, J. Kusuma, and M. Vetterli, "Low-sampling rate uwb channel characterization and synchronization," *J. Comm. Netw.*, vol. 5, pp. 319–327, 2003.
- [16] R. Baranuik, M. Daverport, R. DeVore, and M. Wakin, "a simple proof of the restricted isometry property for random matrices," *Constr. Approx.*, vol. 28, pp. 253–263, 2008.
- [17] S. Mendelson, A. Pajor, and N. Tomczak-Jaegermann, "Uniform uncertainty principle for bernoulli and subgaussian ensembles," *Constr. Approx.*, vol. 28, pp. 277–289, 2008.
- [18] J. A. Tropp and A. C. Gilbert, "signal recovery from random measurements via orthogonal matching pursuit," *IEEE Trans. on Inform. Theory*, vol. 53, pp. 4655–4666, 2007.
- [19] L. Welch, "Lower bounds on the maximum cross correlation of signals," *IEEE Trans. Inform. Theory*, vol. 20, pp. 397–399, 1974.
- [20] Y. C. Eldar and G. Kutyniok, Eds., *Compressed Sensing: Theory and Applications*. Cambridge Univ. Press, 2012.
- [21] S. Muthukrishnan, *Data streams: Algorithms and applications*. Now Publishers Inc, 2005.
- [22] S. Chen, D. Donoho, and M. Saunders, "Atomic decomposition by basis pursuit," *SIAM J. Sci Comp.*, vol. 20, pp. 33–61, 1999.
- [23] D. Donoho and X. Huo, "Uncertainty principle and ideal atomic decomposition," *IEEE Trans on Information Theory*, vol. 47, pp. 2845–2862, 2001.

-
- [24] D. L. Donoho and M. Elad, "Optimally sparse representation in general (nonorthogonal) dictionaries via ℓ_1 minimization," *Proceedings of the National Academy of Sciences*, vol. 100, no. 5, pp. 2197–2202, 2003.
- [25] D. Donoho, M. Elad, and N. Temlyakov, "Stable recovery of sparse overcomplete representations in the presence of noise," *IEEE Transactions on Information Theory*, vol. 52, no. 1, pp. 6–18, 2006.
- [26] M. Elad and M. Bruckstein, "A generalized uncertainty principle and sparse representation in pairs of bases," *IEEE Transactions on Information Theory*, vol. 48, no. 9, pp. 2558–2567, 2002.
- [27] J.-J. Fuchs, "On sparse representations in arbitrary redundant bases," *Information Theory, IEEE Transactions on*, vol. 50, no. 6, pp. 1341–1344, 2004.
- [28] R. Gribonval and M. Nielsen, "Sparse representations in unions of bases," *Information Theory, IEEE Transactions on*, vol. 49, no. 12, pp. 3320–3325, 2003.
- [29] A. Tropp, "Just relax: Convex programming methods for identifying sparse signals in noise," *IEEE Transactions on Information Theory*, vol. 52, no. 3, pp. 1030–1051, 2006.
- [30] Y. Tsaig and D. Donoho, "Breakdown of equivalence between the minimal ℓ_1 -norm solution and the sparsest solution," *Signal Processing*, vol. 86, no. 3, pp. 533–548, 2006.
- [31] S. Becker, J. Bobin, and E. Candès, "Nesta: a fast and accurate first-order method for sparse recovery," *SIAM Journal on Imaging Sciences*, vol. 4, no. 1, pp. 1–39, 2011.
- [32] S. Becker, E. Candès, and M. Grant, "Templates for convex cone problems with applications to sparse signal recovery," *Mathematical Programming Computation*, vol. 3, no. 3, pp. 165–218, 2011.
- [33] A. Galatsanos, N. and Katsaggelos, "Methods for choosing the regularization parameter and estimating the noise variance in image restoration and their relation," *Image Processing, IEEE Transactions on*, vol. 1, no. 3, pp. 322–336, 1992.
- [34] Y. C. Eldar, "Generalized sure for exponential families: Applications to regularization," *Signal Processing, IEEE Transactions on*, vol. 57, no. 2, pp. 471–481, 2009.
- [35] S. Mallat and Z. Zhang, "Matching pursuits with time-frequency dictionaries," *IEEE Transactions on Signal Processing*, vol. 41, no. 12, pp. 3397–3415, Dec 1993.

-
- [36] D. Donoho, Y. Tsaig, I. Drori, and J.-L. Starck, “Sparse solution of underdetermined systems of linear equations by stagewise orthogonal matching pursuit,” *IEEE Transactions on Information Theory*, vol. 58, no. 2, pp. 1094–1121, Feb 2012.
- [37] D. Needell and A. Tropp, “Cosamp: Iterative signal recovery from incomplete and inaccurate samples,” *Applied and Computational Harmonic Analysis*, vol. 26, no. 3, pp. 301–321, 2009.
- [38] T. Blumensath and M. Davies, “Iterative hard thresholding for compressed sensing,” *Applied and Computational Harmonic Analysis*, vol. 27, no. 3, pp. 265–274, 2009.
- [39] —, “Gradient pursuits,” *IEEE Transactions on Signal Processing*, vol. 56, no. 6, pp. 2370–2382, June 2008.
- [40] —, “Iterative thresholding for sparse approximations,” *Journal of Fourier Analysis and Applications*, vol. 14, no. 5-6, pp. 629–654, 2008.
- [41] W. Dai and O. Milenkovic, “Subspace pursuit for compressive sensing signal reconstruction,” *Information Theory, IEEE Transactions on*, vol. 55, no. 5, pp. 2230–2249, May 2009.
- [42] I. Daubechies, M. Defrise, and C. De Mol, “An iterative thresholding algorithm for linear inverse problems with a sparsity constraint,” *Communications on pure and applied mathematics*, vol. 57, no. 11, pp. 1413–1457, 2004.
- [43] M. Davenport and M. Wakin, “Analysis of orthogonal matching pursuit using the restricted isometry property,” *Information Theory, IEEE Transactions on*, vol. 56, no. 9, pp. 4395–4401, Sept 2010.
- [44] D. Donoho and Y. Tsaig, “Fast solution of l_1 -norm minimization problems when the solution may be sparse,” *IEEE Trans on Information Theory*, vol. 54, pp. 4789–4812, 2008.
- [45] M. Elad, B. Matalon, J. Shtok, and M. Zibulevsky, “A wide-angle view at iterated shrinkage algorithms,” in *Optical Engineering+ Applications*. International Society for Optics and Photonics, 2007, pp. 670 102–670 102.
- [46] P. Indyk and M. Ruzic, “Near-optimal sparse recovery in the l_1 norm,” in *Foundations of Computer Science, 2008. FOCS 08. IEEE 49th Annual IEEE Symposium on*, Oct 2008, pp. 199–207.

- [47] D. G. and S. Verdu, “Randomly spread cdma: asymptotics via statistical physics,” *IEEE Transactions on Information Theory*, vol. 51, no. 6, pp. 1983–2010, June 2005.
- [48] J. Boutros and G. Caire, “Iterative multiuser joint decoding: Unified framework and asymptotic analysis,” *IEEE Trans. on Inf. Theory*, vol. 48, pp. 1772–1793, 2002.
- [49] D. Guo and C.-C. Wang, “Asymptotic mean-square optimality of belief propagation for sparse linear systems,” in *Proc. IEEE Inf. Theory Workshop*, 2006.
- [50] A. Maleki, “Approximate message passing algorithms for compressed sensing,” Ph.D. dissertation, Stanford University, 2011.
- [51] C. A. Metzler, A. Maleki, and R. G. Baraniuk, “From denoising to compressed sensing,” arXiv:1406.4175v3 [cs.IT], July 2014.
- [52] D. Donoho, A. Maleki, and A. Montanari, “How to design message passing algorithms for compressed sensing,” 2011. [Online]. Available: <http://www.ece.rice.edu/mam15/bpist.pdf>
- [53] F. Krzakala, M. Mézard, F. Sausset, Y. Sun, and L. Zdeborová, “Probabilistic reconstruction in compressed sensing: Algorithms, phase diagrams, and threshold achieving matrices,” *J. Stat. Mech.*, vol. P08009, Aug. 2012.
- [54] J. Vila and P. Schniter, “Expectation-maximization gaussian-mixture approximate message passing,” *Signal Processing, IEEE Transactions on*, vol. 61, no. 19, pp. 4658–4672, Oct 2013.
- [55] M. Bayati and A. Montanari, “The dynamics of message passing on dense graphs, with applications to compressed sensing,” *IEEE Trans. on Inf. Theory*, vol. 57, no. 2, pp. 764–785, Feb. 2011.
- [56] D. Donoho, A. Maleki, and A. Montanari, “Message passing algorithms for compressed sensing: I. motivation and construction,” in *IEEE Inf. Theory Workshop (ITW)*. Dublin, Ireland, 2010, pp. 1–5.
- [57] J. D. Blanchard, C. Cartis, J. Tanner, and A. Thompson, “Phase transitions for greedy sparse approximation algorithms,” *Applied and Computational Harmonic Analysis*, vol. 30, no. 2, pp. 188–203, 2011.
- [58] A. Montanari, “Graphical models concepts in compressed sensing,” *Compressed Sensing: Theory and Applications*, pp. 394–438, 2012.

- [59] D. Donoho and J. Tanner, “Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing,” *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 367, no. 1906, pp. 4273–4293, 2009.
- [60] F. Kschischang, B. Frey, and H.-A. Loeliger, “Factor graphs and the sum-product algorithm,” *IEEE Trans. on Inf. Theory*, vol. 47, no. 2, pp. 498–519, Feb. 2001.
- [61] D. MacKay, *Information Theory, Inference and Learning Algorithms*. Cambridge Univ. Press, 2002.
- [62] J. Pearl, *Probabilistic reasoning in intelligent systems*. Morgan Kaufmann Publishers Inc., 1988.
- [63] F. Jensen, *An Introduction to Bayesian Networks*. New York: Springer-Verlag, 1996.
- [64] B. Frey, *Graphical Models for Machine Learning and Digital Communication*. MA:MIT Press, 1998.
- [65] J. Yedidia, W. Freeman, and Y. Weiss, “Understanding belief propagation and its generalizations,” Mitsubishi Electric Res. Labs, Tech, Rep. TR2001-022, Tech. Rep., 2002.
- [66] R. Cowell, A. Dawid, and S. Lauritzen, *Probabilistic Networks and Expert Systems*. New York: Springer-Verlag, 2003.
- [67] F. Jensen, *Bayesian Networks and Decision Graphs*. New York: Springer-Verlag, 2002.
- [68] R. Koetter, A. Singer, and M. Tüchler, “Turbo equalization,” *IEEE Signal Processing Magazine*, vol. 21, no. 1, pp. 67–80, Jan 2004.
- [69] S. Som and P. Schniter, “Compressive imaging using approximate message passing and a markov-tree prior,” *IEEE Trans. on Signal Process.*, vol. 60, no. 7, pp. 3439–3448, July 2012.
- [70] S. Rangan, “Generalized approximate message passing for estimation with randomly linear mixing,” in *2011 IEEE International Symposium on Information Theory (ISIT)*, July 2011.
- [71] D. Thouless, P. Anderson, and R. Palmer, “Solution of solvable model of a spin glass,” *Philosophical Magazine*, vol. 35, no. 3, pp. 593–601, 1977.

-
- [72] A. Maleki and A. Montanari, "Analysis of approximate message passing algorithm," in *Information Sciences and Systems (CISS), 2010 44th Annual Conference on*, March 2010, pp. 1–7.
- [73] U. Kamilov, S. Rangan, A. K. Fletcher, and M. Unser, "Approximate message passing with consistent parameter estimation and applications to sparse learning," *IEEE Transactions on Information Theory*, pp. 2969–2985, 2014.
- [74] T. Richardson and R. Urbanke, *Modern Coding Theory*. Cambridge Univ. Press, 2008.
- [75] M. Mézard and A. Montanari, *Information, Physics, and Computation*. Oxford Press, 2009.
- [76] M. Mézard, G. Parisi, and A. Virasoro, *Spin-Glass Theory and Beyond*. World Scientific, Singapore, 1987.
- [77] T. Tanaka, "A statistical-mechanics approach to large-system analysis of cdma multiuser detectors," *IEEE Transactions on Information Theory*, vol. 48, no. 11, pp. 2888–2910, Nov 2002.
- [78] Y. Kabashima, T. Wadayama, and T. Tanaka, "A typical reconstruction limit for compressed sensing based on lp-norm minimization," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2009, no. 09, p. L09003, 2009.
- [79] S. Rangan, A. K. Fletcher, and V. K. Goyal, "Asymptotic analysis of map estimation via the replica method and compressed sensing," in *Neural Information Processing Systems (NIPS)*, 2009, pp. 1545–1553.
- [80] S. Ganguli and H. Sompolinsky, "Statistical mechanics of compressed sensing," *Physical review letters*, vol. 104, no. 18, p. 188701, 2010.
- [81] G. Dongning, D. Baron, and S. Shamai, "A single-letter characterization of optimal noisy compressed sensing," in *2009. 47th Annual Allerton Conference on Communication, Control and Computing (Allerton)*, Sept 2009, pp. 52–59.
- [82] A. Maleki and D. Donoho, "Optimally tuned iterative reconstruction algorithms for compressed sensing," *Selected Topics in Signal Processing, IEEE Journal of*, vol. 4, no. 2, pp. 330–341, April 2010.
- [83] D. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, Apr. 2006.

- [84] Y. Tsaig, “Sparse solution of underdetermined linear systems: algorithms and applications,” Ph.D. dissertation, Dept. of Computational Mathematics & Engineering, Stanford University, Stanford, CA, 2007.
- [85] J. Fowler, S. Mun, and E. Tramel, “Multiscale block compressed sensing with smoother projected landweber reconstruction,” in *19th European signal process. conf. (EUSIPCO)*. Barcelona, Spain, 2011, pp. 564–568.
- [86] H. Chang, Y. Weiss, and W. Freeman, “Informative sensing of natural images,” in *16th IEEE Int. Conf. on Image Process. (ICIP)*. Cairo, Egypt, 2009, pp. 3025–3028.
- [87] H. Chang, “Informative sensing: theory and applications,” Ph.D. dissertation, Dept. of Elect. Eng. and Comput. Sci., MIT, Cambridge, Mass., 2012. [Online]. Available: <http://hdl.handle.net/1721.1/74890>
- [88] R. Linsker, *An application of the principle of maximum information preservation to linear systems*, 1st ed. Burlington, Mass.: Morgan Kaufmann Publishers Inc., 1989.
- [89] B. Kashin, “Diameters of some finite-dimensional sets and classes of smooth functions,” *Math. USSR Izvestiya*, vol. 11, pp. 317–333, 1977.
- [90] M. Maiorov, “Discretization of the diameter problem,” *Uspekhi Mat. Nauk*, vol. 30, pp. 179–180, 1975.
- [91] M. Bayati and A. Montanari, “The lasso risk for gaussian matrices,” *IEEE Trans. on Inf. Theory*, vol. 58, no. 4, pp. 1997–2017, Apr. 2012.
- [92] M. Bayati, M. Lelarge, and A. Montanari, “Universality in polytope phase transitions and iterative algorithms,” in *IEEE Int. Symp. on Inf. Theory Process. (ISIT)*, Cambridge, MA, 2012, pp. 1643–1647.
- [93] M. Crouse, R. Nowak, and R. Baraniuk, “Wavelet-based statistical signal processing using hidden markov models,” *IEEE Trans on Signal Process.*, vol. 46, no. 4, pp. 886–902, Apr. 1998.
- [94] M. Duarte, M. Wakin, and G. Baraniuk, “Wavelet-domain compressive signal reconstruction using a hidden markov tree model,” in *IEEE Int. Conf. on Acoust., Speech and Signal Process. (ICASSP)*, Las Vegas, NV, 2008, pp. 5137–5140.

- [95] P. Moulin and J. Liu, "Analysis of multiresolution image denoising schemes using generalized gaussian and complexity priors," *IEEE Trans. on Inf. Theory*, vol. 45, no. 3, pp. 909–919, Apr. 1999.
- [96] C. Bouman and K. Sauer, "A generalized gaussian image model for edge-preserving map estimation," *IEEE Trans. on Image Process.*, vol. 2, no. 3, pp. 296–310, July 1993.
- [97] P. Schniter, "Exploiting structured sparsity in bayesian experimental design," in *4th IEEE Int. Workshop on Computational Advances in Multi-Sensor Adaptive Process. (CAMSAP)*, San Juan, Puerto Rico, 2011, pp. 357–360.
- [98] L. He and L. Carin, "Exploiting structure in wavelet-based bayesian compressive sensing," *IEEE Trans. on Signal Process.*, vol. 57, no. 9, pp. 3488–3497, Sept. 2009.
- [99] Y. Kim, M. Nadar, and A. Bilgin, "Wavelet-based compressed sensing using a gaussian scale mixture model," *IEEE Trans. on Image Process.*, vol. 21, no. 6, pp. 3102–3108, June 2012.
- [100] A. Averbuch, S. Dekel, and S. Deutsch, "Adaptive compressed image sensing using dictionaries," *SIAM J. on Imaging Sciences*, vol. 5, no. 1, pp. 57–89, 2012.
- [101] C. Shannon and W. Weaver, *The mathematical theory of communication*, 1st ed. Urbana, Illinois: University of Illinois Press, 1949.
- [102] S. Rangan, P. Schniter, E. Riegler, A. Fletcher, and V. Cevher, "Fixed points of generalized approximate message passing with arbitrary matrices," arXiv:1301.6295v2 [cs.IT], 2013.
- [103] A. Hyvärinen, J. Karhunen, and E. Oja, *Independent Component Analysis*. John Wiley & Sons Inc, 2001.
- [104] Y. Wu and S. Verdú, "Optimal phase transitions in compressed sensing," *IEEE Trans. on Inf. Theory*, vol. 58, no. 10, pp. 6241–6263, Oct. 2012.
- [105] S. Mallat, *A Wavelet Tour of Signal Processing*, 2nd ed. San Diego, CA: Academic Press, 1999.
- [106] T. Cover and J. Thomas, *Elements of Information Theory*, 2nd ed. Hoboken, N.J.: Wiley-Interscience, 2006.

- [107] G. Lorentz, M. Golitschek, and Y. Makovoz, *Constructive approximation: advanced problems*, 1st ed. Berlin, Germany: Springer-Verlag Berlin and Heidelberg GmbH & Co. K, 1996. [Online]. Available: <http://opac.inria.fr/record=b1091520>
- [108] L. He, H. Chen, and L. Carin, “Tree-structured compressive sensing with variational bayesian analysis,” *IEEE Signal Processing Letters*, vol. 17, no. 3, pp. 233–236, March 2010.
- [109] P. Schniter, “Turbo reconstruction of structured sparse signals,” in *2010 44th Annual Conference on Information Sciences and Systems (CISS)*, March 2010, pp. 1–6.
- [110] A. Dempster, N. Laird, and D. Rubin, “Maximum likelihood from incomplete data via the EM algorithm,” *J. of the Royal Statistical Soc., Series B*, vol. 39, no. 1, pp. 1–38, 1977.
- [111] S. Kudekar and H. Pfister, “The effect of spatial coupling on compressive sensing,” in *2010 48th Annual Allerton Conference on Communication, Control and Computing (Allerton)*, Sept 2010, pp. 347–353.
- [112] S. Kudekar, T. Richardson, and R. Urbanke, “Threshold saturation via spatial coupling: Why convolutional ldpc ensembles perform so well over the bec,” *IEEE Transactions on Information Theory*, vol. 57, no. 2, pp. 803–834, Feb 2011.
- [113] —, “Spatially coupled ensembles universally achieve capacity under belief propagation,” in *2012 IEEE International Symposium on Information Theory Proceedings (ISIT)*, July 2012, pp. 453–457.
- [114] D. Donoho, A. Javanmard, and A. Montanari, “Information-theoretically optimal compressed sensing via spatial coupling and approximate message passing,” in *2012 IEEE International Symposium on Information Theory Proceedings (ISIT)*, July 2012, pp. 1231–1235.
- [115] J. Barbier, F. Krzakala, M. Mezard, and L. Zdeborova, “Compressed sensing of approximately-sparse signals: Phase transitions and optimal reconstruction,” in *2012 50th Annual Allerton Conference on Communication, Control and Computing (Allerton)*, Oct 2012, pp. 800–807.
- [116] H. Jegou, T. Furon, and J. Fuchs, “Anti-sparse coding for approximate nearest neighbor search,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Mar. 2012, pp. 2029–2032.

-
- [117] V. Chandrasekaran, B. Recht, A. Parrilo, and A. Willsky, “The convex geometry of linear inverse problems,” *Foundations of Computational Mathematics*, vol. 12, no. 6, pp. 805–849, 2012.
- [118] D. Donoho, A. Maleki, and A. Montanari, “Message passing algorithms for compressed sensing: Ii. analysis and validation,” in *IEEE Inform. Theory Workshop*, 2010.
- [119] D. Donoho and I. Johnstone, “Adapting to unknown smoothness via wavelet shrinkage,” *J. American Stat. Assoc.*, vol. 90, pp. 1200–1224, 1995.
- [120] J. Pesquet and D. Leporini, “A new wavelet estimator for image denoising,” in *6th IEEE Int. Conf. on Image Process. and its Applicat.*, 1997, pp. 249–253.
- [121] A. Benazza-Benyahia and J. C. Pesquet, “Building robust wavelet estimators for multi-component images using steins principle,” *IEEE Trans. Image Proc.*, vol. 14, pp. 1814–1830, 2005.
- [122] F. Luisier, T. Blu, and M. Unser, “Sure-based wavelet thresholding integrating inter-scale dependencies,” pp. 1457–1460, Oct 2006.
- [123] ———, “A new sure approach to image denoising: Interscale orthonormal wavelet thresholding,” *Image Processing, IEEE Transactions on*, vol. 16, no. 3, pp. 593–606, March 2007.
- [124] T. Blu and F. Luisier, “The sure-let approach to image denoising,” *IEEE Trans. on Image Process.*, vol. 16, no. 11, pp. 2778–2786, Nov 2007.
- [125] M. Raphan and P. Simoncelli, “Least squares estimation without priors or supervision,” *Neural computation*, vol. 23, no. 2, pp. 374–420, 2011.
- [126] C. Stein, “Estimation of the mean of a multivariate normal distribution,” *The annals of Statistics*, pp. 1135–1151, 1981.
- [127] K. Miyasawa, “An empirical bayes estimator of the mean of a normal population,” *Bull. Inst. Internat. Statist.*, vol. 38, no. 181-188, pp. 1–2, 1961.
- [128] R. Gribonval, “Should penalized least squares regression be interpreted as maximum a posteriori estimation?” *Signal Processing, IEEE Transactions on*, vol. 59, no. 5, pp. 2405–2410, May 2011.

- [129] N. Goertz, C. Guo, A. Jung, M. Davies, and G. Doblinger, “Iterative recovery of dense signals from incomplete measurements,” submitted to IEEE signal processing letters.
- [130] E. van den Berg and M. Friedlander, “Probing the pareto frontier for basis pursuit solutions,” *SIAM Journal on Scientific Computing*, vol. 31, no. 2, pp. 890–912, 2008. [Online]. Available: <http://link.aip.org/link/?SCE/31/890>
- [131] A. Polyanin and A. V. Manzhirov, *Handbook of integral equations*, 2nd ed. Boca Raton, Florida: Chapman and Hall/CRC Press, 2008.
- [132] J. Valsa and L. Brancik, “Approximate formulae for numerical inversion of laplace transforms.” *Int. J. Numer. Model.*, vol. 11, pp. 153–166, 1998.