

# CONVERGENCE ANALYSIS OF SPARSE BAYESIAN LEARNING UNDER APPROXIMATE INFERENCE TECHNIQUES

*Christo Kurisummoottil Thomas, Dirk Slock*

EURECOM, Sophia-Antipolis, France, Email: {kurisumm,slock}@eurecom.fr

## ABSTRACT

Sparse Bayesian Learning (SBL) is an efficient and well-studied framework for sparse signal recovery. SBL relies on a parameterized prior on the sparse signal to be estimated. The prior is chosen (with estimated hyperparameters) such that it encourages sparsity in the representation of the signal. However, SBL doesn't scale with problem dimensions due to the computational complexity associated with matrix inversion. To address this issue, there exists low complexity methods based on approximate Bayesian inference. Various state of the art approximate inference methods are based on variational Bayesian (VB) inference or message passing algorithms such as belief propagation (BP) or expectation propagation. Moreover, these approximate inference methods can be unified under the optimization of Bethe free energy with appropriate constraints. SBL allows to treat more general signal models by the use of hierarchical prior formulation which eventually becomes more sparsity inducing than e.g., Laplacian prior. In this paper, we study the convergence behaviour of the mean and variance of the unknown parameters in SBL under approximate Bayesian inference.

## 1. INTRODUCTION

The signal model for the recovery of a time varying sparse signal can be formulated as,  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{v}$ , where  $\mathbf{y}$  is the observations or data,  $\mathbf{A}$  is called the measurement or the sensing matrix which is known and is of dimension  $N \times M$  with  $N < M$ .  $\mathbf{x}$  contains only  $K$  non-zero entries, with  $K \ll M$ . In Bayesian inference, the sparse Bayesian learning (SBL) algorithm was first proposed by [1,2]. SBL is based on a two or three layer hierarchical prior on the sparse coefficients  $\mathbf{x}$ . The priors for the hyperparameters (precision parameters) are chosen such that it induces sparsity allowing majority of the coefficients to tend towards zero. Nevertheless, matrix inversion step involved in SBL at each iteration makes it a computationally complex algorithm even for moderately large datasets. This is the motivation behind looking for alternative solutions or approximate inference methods which has computational requirements proportional to the number of sparse coefficients.

Belief propagation (BP) based SBL algorithm [3] is more computationally efficient than the original algorithm. Due to space limitations we skip the detailed discussion and instead refer the readers to a more detailed discussion on the various approximate inference methods for SBL to our paper [4].

Various studies on convergence analysis of Gaussian BP (GaBP) can be found in [5–8]. Although BP achieves great empirical success [9], not enough rigorous work exist to characterize the convergence behaviour of BP in loopy networks. In [10], convergence condition for GaBP is provided which requires the underlying distribution to be walk-summable. Their convergence analysis is based

on the Gaussian Markov random field (GMRF) based decomposition, where the underlying distribution is expressed in terms of the pairwise connection between the variables.

### 1.1. Contributions of this paper

- We first review low complexity sparse Bayesian learning methods based on SAVE (space alternating variational estimation) [11, 12], BP, MF (mean field), EP (expectation propagation) or a combination of them, which can be unified under the optimization of Bethe free energy (BFE).
- We evaluate the convergence points of SBL solutions based on BP, MF or EP and derive the conditions under which they converge to the exact LMMSE (linear minimum mean squared error) estimates.
- Existing low complexity solutions derived from BP such as AMP or its generalizations converge only for a limited class of measurement matrices  $\mathbf{A}$ . So it becomes imperative again to analyze the convergence behaviour of approximate inference methods such as BP (from which AMP or related algorithms are derived) or variational Bayesian (VB) or EP under more general measurement matrices  $\mathbf{A}$  and we try to address this problem in this paper.
- Utilizing the large system analysis derived in [13], we show that the MSE (mean squared error) of BP converges to the exact LMMSE under i.i.d measurement matrix  $\mathbf{A}$  compared to the previous works (for e.g. AMP) which shows the exactness to LMMSE replica prediction method which is heuristic.

## 2. SBL PROBABILISTIC MODEL

<sup>1</sup>In Bayesian compressive sensing, a two-layer hierarchical prior is assumed for the  $\mathbf{x}$  as in [1]. The hierarchical prior is such that it encourages the sparsity property of  $\mathbf{x}$  or of innovation sequences  $\mathbf{v}$ .

$$p(\mathbf{x}|\mathbf{\Gamma}) = \prod_{i=1}^M \mathcal{CN}(\mathbf{0}, \mathbf{\Gamma}^{-1}), \quad \mathbf{\Gamma} = \text{diag}(\alpha_i). \quad (1)$$

We assume a Gamma prior for  $\mathbf{\Gamma}$ ,  $p(\mathbf{\Gamma}) = \prod_{i=1}^M p(\alpha_i/a, b) =$

$\prod_{i=1}^M \Gamma^{-1}(a) b^a \alpha_i^{a-1} e^{-b\alpha_i}$ . The inverse of noise variance  $\gamma$  is also

assumed to have a Gamma prior,  $p(\gamma/c, d) = \Gamma^{-1}(c) d^c \gamma^{c-1} e^{-d\gamma}$ , such that the marginal pdf of  $\mathbf{x}$  (student-t distribution) becomes more sparsity inducing than e.g. a Laplacian prior. The advantage is that the whole machinery of linear MMSE estimation can be exploited,

<sup>1</sup>Notations: The operator  $(\cdot)^H$  represents the conjugate transpose or conjugate for a matrix or a scalar respectively. In the following, the pdf of a complex Gaussian random variable  $x$  with mean  $\mu$  and variance  $\sigma^2$  is given by  $\mathcal{CN}(x; \mu, \nu)$ .  $KL(q||p)$  represents the Kullback-Leibler distance between the two distributions  $q, p$ .  $\mathbf{A}_{n,\cdot}$  represents the  $n^{th}$  row of  $\mathbf{A}$ .  $\text{blkdiag}(\cdot)$  represents blockdiagonal part of a matrix.  $\text{diag}(\mathbf{X})$  or  $\text{diag}(\mathbf{x})$  represents a vector obtained by the diagonal elements of the matrix  $\mathbf{X}$  or the diagonal matrix obtained with the elements of  $\mathbf{x}$  in the diagonal respectively.  $\mathbf{1}_M$  represents a vector of length  $M$  with all ones as the elements. For a matrix  $\mathbf{A}$ ,  $\mathbf{A} \geq 0$  implies it is non-negative (all the elements of  $\mathbf{A}$  are non-negative).

such as e.g., the Kalman filter. But this is embedded in other layers making things eventually non-Gaussian. Now the likelihood distribution can be written as,  $p(\mathbf{y}/\mathbf{x}, \gamma) = (2\pi)^{-N} \gamma^N e^{-\frac{\gamma}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2}$ . To make these priors non-informative, we choose them to be small values  $a = c = b = d = 10^{-5}$ . We define the unknown parameter vector  $\theta = \{\mathbf{x}, \Gamma, \gamma\}$  and  $\theta_i$  represents each scalar in  $\theta$ .

### 3. BETHE FREE ENERGY OPTIMIZATION

The fixed points of the standard BP algorithm are shown to be the stationary points of the BFE. However, for the MF approximation in variational Bayes [14], the approximate posteriors are shown to be converging to a local minimum of the MF free energy which is an approximation of the BFE. However, we observe in [11] that for estimation of the signals from interference corrupted observations, MF is a poor choice since it doesn't give the accurate posterior variance (posterior variance of  $x_i$  is observed to be independent of the error variances of other  $x_l, l \neq i$ ). Assume that the posterior be represented as,  $p(\theta) = \frac{1}{Z} \prod_{a \in \mathcal{A}_{BP}} f_a(\theta_a) \prod_{b \in \mathcal{A}_{MF}} f_b(\theta_b)$ ,

where  $\mathcal{A}_{BP}, \mathcal{A}_{MF}$  represent the set of nodes belonging to the BP part and MF part respectively with  $\mathcal{A}_{BP} \cap \mathcal{A}_{MF} = \emptyset$ .  $Z$  represents the normalization variable.  $\mathcal{N}(i), \mathcal{N}(a)$  represent the number of neighbouring nodes of any variable node  $i$  or factor node  $a$ .  $\mathcal{N}_{BP}(i)$  represents the number of neighbouring nodes of  $i$  which belong to the BP part, similarly  $\mathcal{N}_{MF}(i)$  is defined. Also, we define  $\mathcal{I}_{MF} = \bigcup_{a \in \mathcal{A}_{MF}} \mathcal{N}(a), \mathcal{I}_{BP} = \bigcup_{a \in \mathcal{A}_{BP}} \mathcal{N}(a)$ . The optimization of the resulting free energy obtained by the combination of BP and MF [4, eq.(2)] (Note that we use an abuse of notation and let  $q_i(\theta_i)$  represents the belief about  $\theta_i$  (the approximate posterior)) leads to the following message passing expressions. Let  $m_{a \rightarrow i}$  represents the message passed from any factor node  $a$  to variable node  $i$  and  $n_{i \rightarrow a}$  represents the message passed from any variable node  $i$  to factor node  $a$ . The fixed point equations are,

$$\begin{aligned} q_i(\theta_i) &= z_i \prod_{a \in \mathcal{N}_{BP}(i)} m_{a \rightarrow i}^{BP}(\theta_i) \prod_{a \in \mathcal{N}_{MF}(i)} m_{a \rightarrow i}^{MF}(\theta_i), \\ n_{i \rightarrow a}(\theta_i) &= \prod_{a \in \mathcal{N}_{BP}(i) \setminus a} m_{a \rightarrow i}(\theta_i) \prod_{a \in \mathcal{N}_{MF}(i)} m_{a \rightarrow i}(\theta_i), \\ m_{a \rightarrow i}^{MF}(\theta_i) &= \exp(\langle \ln f_a(\theta_a) \rangle_{q_i}) \prod_{j \in \mathcal{N}(a) \setminus i} n_{j \rightarrow a}(\theta_j), \\ m_{a \rightarrow i}^{BP}(\theta_i) &= \left( \int \prod_{j \in \mathcal{N}(a) \setminus i} n_{j \rightarrow a}(\theta_j) f_a(\theta_a) \prod_{j \neq i} d\theta_j \right), \end{aligned} \quad (2)$$

where  $\langle \cdot \rangle_q$  represents the expectation w.r.t distribution  $q$ .

The constraints in BFE can often be too complex to yield computationally tractable messages ( $m_{a \rightarrow i}, n_{a \rightarrow i}$ ), the following constraint relaxation leads to EP [15].

$$\begin{aligned} E_{q_a}(t(\theta_i)) &= E_{q_i}(t(\theta_i)), \text{ leads to,} \\ m_{a \rightarrow i}^{BP}(\theta_i) &= \frac{\text{Proj}_{\phi} \left( \int \prod_{j \in \mathcal{N}(a)} n_{j \rightarrow a}(\theta_j) f_a(\theta_a) \prod_{j \neq i} d\theta_j \right)}{n_{i \rightarrow a}(\theta_i)}, \end{aligned} \quad (3)$$

where  $\phi$  represents the family of distributions characterized by the sufficient statistics  $t(\theta_i)$ . In the following sections, we give a brief overview of the large system analysis techniques we propose to use to evaluate the convergence behaviour of SBL using BP/VB/EP based inference.

#### 3.1. SBL using Belief Propagation

We first review the BP messages being passed between the variable nodes and factor nodes corresponding to the factor graph in Figure 1. All the messages (beliefs or continuous pdfs) passed between them are all Gaussian [3]. So in message passing (MP), it suffices to represent them by two parameters, which are the mean and variance of the beliefs. Also, for the first instance, we assume that all the hyperparameters are known. We remark that the estimation of hyperparameters can be done using VB as in [11]. Below, indices  $m, n$  is

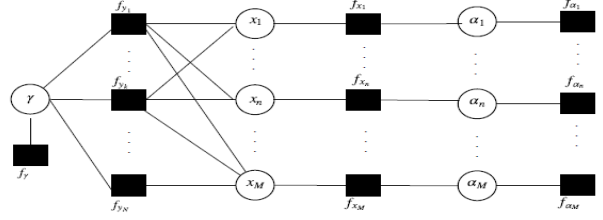


Fig. 1. Factor Graph for the static SBL.

used for representing variable nodes and  $i, k$  is used for representing factor nodes. We represent  $S_{n,k}$  as the inverse variance (precision) of the message passed from variable node  $n$  (corresponding to  $x_n$ ) to factor node  $k$  (corresponds to  $y_k$ ) and  $M_{n,k}$  be the mean of the message passed from  $n$  to  $k$ , total  $NM$  of them. Similarly  $S_{k,n}, M_{k,n}$  for messages from  $k$  to  $n$ . Let  $A_{k,n}$  represents the  $(k, n)^{th}$  element of  $\mathbf{A}$ . We start with the message passing expressions derived in [3].

$$\begin{aligned} S_{n,k} &= \alpha_n + \sum_{i \neq k} S_{i,n}, \quad M_{n,k} = S_{n,k}^{-1} \sum_{i \neq k} S_{i,n} M_{i,n}, \\ S_{k,n} &= A_{k,n}^2 \left( \frac{1}{\gamma} + \sum_{m \neq n} A_{k,m}^2 S_{m,k}^{-1} \right)^{-1}, \\ M_{k,n} &= A_{k,n}^{-1} \left( y_k - \sum_{m \neq n} A_{k,m} M_{m,k} \right), \end{aligned} \quad (4)$$

Note that instead of BP, if we use MF for the estimation of  $\mathbf{x}$ , the expressions above would remain the same except  $S_{k,n}$  which gets written as  $S_{k,n} = A_{k,n}^2 \gamma$ . This can be interpreted as, MF does not take into account the error variances in other  $x_m, m \neq n$  while passing the belief about  $x_n$  from any factor node  $y_k$  and hence it is suboptimal. Further, substituting  $S_{n,k}$  in  $S_{k,n}$ ,  $S_{k,n} = A_{k,n}^2 \left( \frac{1}{\gamma} + \sum_{m \neq n} A_{k,m}^2 (\alpha_m + \sum_{i \neq k} S_{i,m})^{-1} \right)^{-1}$ , so this is now only in terms of the message variances in the direction  $k$  to  $n$ . Finally, the belief (estimates) computed for each  $x_n$  is,

$$\sigma_n^2 = (\alpha_n + \sum_i S_{i,n})^{-1}, \quad \mu_n = \sigma_n^2 \left( \sum_i S_{i,n} M_{i,n} \right). \quad (5)$$

Further we simplify the messages and beliefs using the results from random matrix theory, for the simplest case of i.i.d  $\mathbf{A}$  in the large system regime where  $M, N \rightarrow \infty$  at a fixed ratio  $\frac{N}{M} > 0$  (represented in short as  $\xrightarrow[a.s]{M \rightarrow \infty}$ ). For the large system analysis, we use Theorem 1 and Lemma 4 from [13]. We briefly summarize the Lemma's here. Lemma 4 in Appendix VI of [13] states that  $\mathbf{x}_N^H \mathbf{A}_N \mathbf{x}_N \xrightarrow[N \rightarrow \infty]{a.s.} (1/N) \text{tr} \mathbf{A}_N$  when the elements of  $\mathbf{x}_N$  are iid with variance  $1/N$  and independent of  $\mathbf{A}_N$ , and similarly when  $\mathbf{y}_N$  is independent of  $\mathbf{x}_N$ , that  $\mathbf{x}_N^H \mathbf{A}_N \mathbf{y}_N \xrightarrow[N \rightarrow \infty]{a.s.} 0$ . Theorem 1 from [13] implies that any terms of the form  $\frac{1}{N} \text{tr} \{ (\mathbf{A}_N - z \mathbf{I}_N)^{-1} \}$ , where  $\mathbf{A}_N$  is the summation of independent rank one matrices with covariance matrix  $\Theta_i$  is equal to the unique positive solution of  $e_j = \frac{1}{N} \text{tr} \{ \left( \sum_{i=1}^K \frac{\Theta_i}{1+e_i} - z \mathbf{I}_N \right)^{-1} \}$ . In the large system limit, we can approximate (neglecting terms of  $\mathcal{O}(A_{i,j}^2)$ )  $S_{n,k} = \alpha_n + \sum_i S_{i,n} = S_n$ , independent of  $k$ . Further we define  $\mathbf{S} = \text{diag}(S_n)$ . Considering the term  $S_{k,n} = A_{k,n}^2 \left( \frac{1}{\gamma} + \sum_{m \neq n} A_{k,m}^2 S_{m,k}^{-1} \right)^{-1}$ , in the large system it can be approximated by  $S_{k,n} = A_{k,n}^2 \left( \frac{1}{\gamma} + \mathbf{A}_{k,:} \mathbf{S}^{-1} \mathbf{A}_{k,:}^T \right)^{-1}$ .

$\mathbf{A}_{k,:} \mathbf{S}^{-1} \mathbf{A}_{k,:}^T \xrightarrow[a.s.]{M \rightarrow \infty} \frac{1}{M} \text{tr} \{ \mathbf{S}^{-1} \} = \tau'_{BP}$ . From (5), it follows that  $MSE = \tau_{BP} = \text{tr} \{ \mathbf{S}^{-1} \}$ .  $\mathbf{A}_{k,:}$  represents the  $k^{th}$  row of  $\mathbf{A}$ . Further we obtain,  $S_n = \alpha_n + \left( \frac{1}{\gamma} + \tau'_{BP} \right)^{-1} \sum_i A_{i,n}^2$ ,  $\sum_i A_{i,n}^2 \xrightarrow[a.s.]{M \rightarrow \infty} 1$ , thus  $S_n = \alpha_n + \left( \frac{1}{\gamma} + \tau'_{BP} \right)^{-1}$ . Finally we can conclude that,  $\tau'_{BP}$  can be obtained as the unique positive solution of the following fixed point equation,

$$\tau'_{BP} = \sum_{n=1}^M \left( \alpha_n + \left( \frac{1}{\gamma} + \tau'_{BP} \right)^{-1} \right)^{-1}. \quad (6)$$

Next step is to simplify the expression for LMMSE posterior covariance in the large system limit using similar techniques as above. The posterior covariance can be written as,

$$\begin{aligned} \Sigma_L &= \Gamma^{-1} - \Gamma^{-1} \mathbf{A}^T (\mathbf{A} \Gamma^{-1} \mathbf{A}^T + \frac{1}{\gamma})^{-1} \mathbf{A} \Gamma^{-1}, \\ \mathbf{A}^T (\mathbf{A} \Gamma^{-1} \mathbf{A}^T + \frac{1}{\gamma})^{-1} \mathbf{A} &\stackrel{M \rightarrow \infty}{(a)} \mathbf{D}, \mathbf{D}_{i,i} = \frac{e}{1 + \frac{e}{\alpha_i}}, \end{aligned} \quad (7)$$

where (a) follows from Theorem 1 in [13] and  $e$  is defined as the unique positive solution of the following fixed point equation,

$$e = \left( \frac{1}{N} \sum_{i=1}^M \frac{\alpha_i^{-1}}{1 + \frac{e}{\alpha_i}} + \frac{1}{\gamma} \right)^{-1}, \text{tr}\{\Sigma_L\} = MSE = \sum_{i=1}^M \frac{\alpha_i^{-1} e}{1 + \frac{e}{\alpha_i}},$$

$$\text{From } e, \frac{1}{e} - \frac{1}{\gamma} = \frac{1}{N} \sum_{i=1}^M \frac{\alpha_i^{-1}}{1 + \frac{e}{\alpha_i}} = \frac{1}{N} MSE = \frac{\tau}{N} = \tau',$$

$$\frac{1}{e} = \frac{1}{\gamma} + \tau', \tau' = \frac{1}{N} \sum_{i=1}^M \frac{\alpha_i^{-1} (\frac{1}{\gamma} + \tau')}{\frac{1}{\gamma} + \tau' + \frac{1}{\alpha_i}} = \frac{1}{N} \sum_{i=1}^M \frac{1}{\alpha_i + (\frac{1}{\gamma} + \tau')^{-1}}. \quad (8)$$

Comparing (6) and (8), it can be observed that the MSE under BP,  $\tau_{BP}$  and the MMSE  $\tau$  can be obtained as a unique positive solution of the same fixed point equation. This implies that in the large system limit, under i.i.d  $\mathbf{A}$ , if BP converges, the MSE of SBL (assuming the hyperparameters are fixed or known) converges to the exact MMSE. Another remark is that the above large system analysis based on [13] can be applied to more general measurement matrices case, with rows of  $\mathbf{A}$  being restricted to have different covariance matrices, i.e.  $\mathbf{E}(\mathbf{A}_{i,:}^H \mathbf{A}_{i,:}) = \Theta_i$ .

Certain remarks comparing the existing convergence conditions for belief propagation is as follows. In [5], Jian Du et al. shows that depending on the underlying graphical structure (GMRF or factor graph based factorization) GaBP may exhibit different convergence properties. They prove that the convergence condition for the mean provided based on the factor graph representation encompasses much larger class of models than those given by the GMRF based walk-summable condition [10]. Further they show that GaBP always converges if the factor graph is a union of single loop and a forest. Moreover, they also analyze the convergence of the inverse of the message variances (message information matrix) and analytically show that with arbitrary positive semidefinite matrix initialization, the message information matrix converges to a unique positive definite matrix. So we can conclude that for BP there is a decoupling between the dynamics of the variance updates and that of the mean updates. And that we know that the mean converges to the LMMSE estimate under certain conditions. But it is to be mentioned that the convergence conditions and convergence values for the variance are more tricky, still requires rigorous analysis to characterize its behaviour, which is the main motivation behind this paper.

### 3.2. Iterations in Matrix Form

Let us denote  $d(\mathbf{A})$  as the vector with entries as the diagonal elements of  $\mathbf{A}$ .  $\mathbf{B}$  is defined as the matrix with entries as  $A_{i,j}^2$ . Let  $\mathbf{L}$  (of size  $M \times N$ ),  $\mathbf{S}$ ,  $\mathbf{M}$  (of size  $N \times M$ ) be the matrix with entries  $S_{n,k} M_{n,k}$ ,  $S_{k,n}$  and  $M_{k,n}$ , respectively. Defining  $\mathbf{T}$  to be a matrix of size  $M \times N$ , with entries as the inverse variance of the Gaussian messages transmitted from the variable nodes,  $S_{n,k}$ , we obtain,

$$\begin{aligned} \mathbf{T} &= (d(\Gamma) + \mathbf{S}^T \mathbf{1}_N) \otimes \mathbf{1}_M^T - \mathbf{S}^T, \\ \mathbf{L} &= d(\mathbf{S}^T \mathbf{M}) \otimes \mathbf{1}_N^T - (\mathbf{S} \circ \mathbf{M})^T, \mathbf{L}' = \mathbf{T}^{-1} \circ \mathbf{L}. \end{aligned} \quad (9)$$

We denote any matrix  $\mathbf{A}_{inv}$  as a matrix with entries as the element wise inverse of the matrix  $\mathbf{A}$ . Similarly, for the messages at the factor nodes, define  $\mathbf{C}$  to be the matrix with entries  $A_{k,n}^2 S_{k,n}^{-1}$ ,

$$\begin{aligned} \mathbf{C} &= \left( \frac{1}{\gamma} \mathbf{1}_N + d(\mathbf{B} \mathbf{T}_{inv}) \right) \otimes \mathbf{1}_M^T - \mathbf{B} \circ \mathbf{T}_{inv}^T, \mathbf{S} = \mathbf{C}_{inv} \circ \mathbf{B}, \\ \mathbf{V} &= (\mathbf{y} - d(\mathbf{A} \mathbf{L}')) \otimes \mathbf{1}_M^T + \mathbf{A} \circ \mathbf{L}'^T, \mathbf{M} = \mathbf{A}_{inv} \circ \mathbf{V}, \end{aligned} \quad (10)$$

where  $\mathbf{V}$  being the matrix with entries  $A_{k,n} M_{k,n}$ . The computational complexity of all the matrix operations above is  $\mathcal{O}(MN)$ , since the number of computations in the Hadamard product or Kronecker products in the above expressions is only  $MN$ . Assuming the number of iterations required to converge is  $N_{it}$ , the total complexity of the BP algorithm can be written as  $N_{it} \mathcal{O}(MN)$ .

### 3.3. Convergence Analysis of BP

In this subsection, we consider the convergence analysis of the mean and variance of the messages passed in BP. For the ease of analysis, we consider a simplified case, where we neglect terms of the order  $\mathcal{O}(A_{i,j}^2)$  under the large system limit  $M, N \rightarrow \infty$ . Hence the precisions of the posteriors passed  $A_{k,n}^{-1} S_{k,n}$ ,  $S_{n,k}$  in (4) can be approximated as  $S_n = \alpha_n + \sum_i S_{i,k}$  and  $A_{k,n}^{-2} S_{k,n} = (\frac{1}{\gamma} + \sum_m A_{k,m}^2 S_{m,k}^{-1})^{-1} \triangleq S_k$ . In fact,  $S_n, A_{k,n}^{-2} S_{k,n}$  represent the precision variables in the input and output stages of the GAMP algorithm derived in [16, Algorithm 1]. Using theorem 1 in [16], we can show that for any non-negative matrix  $B \succeq 0$ ,  $S_n, S_k$  converge to a positive value. However, we remark that it remains to be understood to which value these precision variables converge (and hence the posterior variance  $\sigma_n^2$ ) and it is left as a future work.

Further we look at the convergence behaviour of the mean value of the posteriors passed across the graph  $M_{k,n}$ . Substituting the value of  $M_{m,k}$  in the expression of  $M_{k,n}$  in (4), we obtain,

$$M_{k,n} = A_{k,n}^{-1} (y_k - \sum_{m \neq n} \sum_{i \neq k} A_{k,m} A_{i,m}^2 S_m^{*-1} S_i^* M_{i,m}), \quad (11)$$

where  $S_i^*, S_m^*$  are the converged values of the precision variables  $S_i, S_m$ , respectively. Defining  $\mathbf{m}^{(t)}$  as a vector of length  $MN$ , representing the values  $M_{k,n}$  at iteration  $t$ . So  $\mathbf{m}^{(t)} = [M_{1,1}, M_{1,2}, \dots, M_{1,M}, \dots, M_{N,1}, \dots, M_{N,M}]^T$ . Also, we define  $\mathbf{N}$  to be a diagonal matrix of length  $MN \times MN$  with entries  $A_{k,n}^{-1}$  and  $\mathbf{M}$  to be a  $MN \times MN$  matrix with  $((i-1)M + m)^{th}$  entry of the  $k^{th}$  row of  $\mathbf{M}$  being defined as  $A_{k,m} A_{i,m}^2 S_m^{-1} S_i$ , but equal to zero when either  $i = k$  or  $m = n$  or  $i = k$  and  $m = n$ .

$$\mathbf{m}^{(t+1)} = -\mathbf{M} \mathbf{m}^{(t)} + \mathbf{N} (\mathbf{y} \otimes \mathbf{1}_M). \quad (12)$$

The above iterations (12) converges if  $\rho(\mathbf{M}) < 1$ .

### 3.4. Scalar Iterations

Further defining the following terms,

$$Z_{k,n} = (y_k - \sum_{m \neq n} A_{k,m} M_{m,k}), \text{ So } M_{k,n} = A_{k,n}^{-1} Z_{k,n}. \quad (13)$$

Also, assume that in the large system limit,  $M_{n,k}$  can be written as,  $M_{n,k} = M_n + \delta_{n \rightarrow k}$ , where  $\delta_{n \rightarrow k}$  is of the  $\mathcal{O}(\frac{1}{\sqrt{N}})$ . This approximation follows from writing  $M_{n,k} = S_{n,k}^{-1} \sum_{i \neq k} S_{i,n} M_{i,n} = S_{n,k}^{-1} \sum_i S_{i,n} M_{i,n} - M_{k,n}$ . Substituting  $M_{n,k}$  in  $Z_{k,n}$ ,  $Z_{k,n} = (y_k - \sum_m A_{k,m} M_m - \sum_m A_{k,m} \delta_{m \rightarrow k} + A_{k,n} M_n + \mathcal{O}(\frac{1}{N})) = Z_k + \delta_{k \rightarrow n}$ , all the terms containing  $A_{i,j}^2$  or  $A_{i,j} \delta_{j \rightarrow i}$  becomes  $\mathcal{O}(\frac{1}{N})$  and  $\delta_{k \rightarrow n} = A_{k,n} M_n$ , also here  $Z_k = (y_k - \sum_m A_{k,m} M_m - \sum_m A_{k,m} \delta_{m \rightarrow k})$ .

$$\begin{aligned} M_{n,k} &= S_{n,k}^{-1} (\frac{1}{\gamma} + \tau'_{BP})^{-1} \sum_{i \neq k} A_{i,n} Z_{i,n} \\ &= S_n^{-1} (\frac{1}{\gamma} + \tau'_{BP})^{-1} \sum_{i \neq k} A_{i,n} Z_{i,n}. \end{aligned} \quad (14)$$

As in the papers by Montanari et. al. [17], for general priors, it is possible to write  $M_{n,k} = f_n(\sum_{i \neq k} A_{i,n} Z_{i,n})$ . Here  $f_n$  is a linear function for the Gaussian case (i.e.  $f_n(x) = S_n^{-1} (\frac{1}{\gamma} + \tau)^{-1} x_n$ ).

So if we consider the case of Gamma priors for  $\alpha$  etc, then this parameterization in terms of an  $f$  becomes easy to write the recursions. Now doing a first order Taylor series approximation of  $f$  around  $\sum_i A_{i,n} Z_{i,n}, M_{n,k} = f_n(\sum_i A_{i,n} Z_{i,n}) -$

$$A_{k,n} Z_{k,n} f'_n(\sum_i A_{i,n} Z_{i,n}), f'_n \text{ being derivative evaluated at } \sum_i A_{i,n} Z_{i,n}. \text{ Further substituting for } Z_{i,n} \text{ from (13),}$$

$$M_{n,k} = M_n + \delta_{n->k}, M_n = f_n(\sum_i A_{i,n} Z_i + \sum_i A_{i,n} \delta_{i->n})$$

$$\text{and } \delta_{n->k} = -A_{k,n} Z_k f'_n(\sum_i A_{i,n} Z_i).$$

Note that term  $A_{k,n} \delta_{k,n}$  becomes  $O(\frac{1}{N})$ . Substituting for  $\delta_{i->n}$  and with the large system approximation  $\sum_i A_{i,n}^2 - > 1, M_n = f_n(\sum_i A_{i,n} Z_i + \sum_i A_{i,n}^2 M_n) = f_n(\sum_i A_{i,n} Z_i + M_n)$ , Now further writing as a vector  $M$  (with each element  $M_n, \forall n$ ).  $M = f(A^T Z + M)$ , which is the AMP recursion for the mean and  $f_n(\cdot)$  represents each of the scalar components in  $f(\cdot)$ . Also in (13), substituting for  $\delta_{n->k}$  from (15),

$$Z_k = (y_k - \sum_m A_{k,m} M_m) + (\frac{1}{\delta}) Z_k (\frac{1}{n}) \sum_m f'_m(\sum_i A_{i,m} Z_i)$$

$$= (y_k - \sum_m A_{k,m} M_m) + \frac{1}{M} Z_k \sum_m f'_m(\sum_i A_{i,m} Z_i),$$

where  $(\frac{1}{M}) Z_k \sum_m f'_m(\sum_i A_{i,m} Z_i)$  is the Onsager term. (16)

#### 4. SBL USING MEAN FIELD APPROXIMATION

For MF or SAVE (space alternating variational estimation) [11], to obtain the free energy  $F(q) = U(q) - H(q)$ ,

$$U(q) = -E_q \ln p(\mathbf{x} | \mathbf{y}) = E_q (\frac{1}{2} \mathbf{x}^T \Sigma_L^{-1} \mathbf{x} - 2 \mathbf{y}^T \mathbf{A} \mathbf{x}) =$$

$$\frac{1}{2} \boldsymbol{\mu}^T \Sigma_L^{-1} \boldsymbol{\mu} - 2 \mathbf{y}^T \mathbf{A} \boldsymbol{\mu} + \sum_i \sigma_i^2 (\Sigma_L^{-1})_{i,i} + c_1,$$

$$H(q) = - \sum_i E_{q_i} \ln q_i = \frac{1}{2} \sum_i \ln \sigma_i^2 + c_2,$$

$c_i$  being constants, independent of  $\boldsymbol{\mu}$  and  $\sigma_i^2$ , also  $q_i(x_i) = \mathcal{N}(\mu_i, \sigma_i^2), \boldsymbol{\mu} = \hat{\mathbf{x}} = [\mu_1, \dots, \mu_M]^T$ . Now the MF free energy can be written as,

$$F(q) = \frac{1}{2} \boldsymbol{\mu}^T \Sigma_L^{-1} \boldsymbol{\mu} - 2 \mathbf{y}^T \mathbf{A} \boldsymbol{\mu} + \sum_i \sigma_i^2 (\Sigma_L^{-1})_{i,i} + \frac{1}{2} \sum_i \ln \sigma_i^2 + c.$$

It can be noticed that  $F(q)$  is a convex function w.r.t  $\boldsymbol{\mu}$  and  $\sigma_i^2$ , further optimizing this w.r.t  $\boldsymbol{\mu}$  leads to  $\boldsymbol{\mu} = \Sigma_L \mathbf{A}^T \mathbf{y}$  and  $\sigma_i^2 = \frac{1}{(\Sigma_L^{-1})_{i,i}}$ . So we can conclude that the mean converges to LMMSE in the case of SAVE while the variance is not exact. Further, we analyze the convergence conditions. The SAVE iterations for  $\boldsymbol{\mu}$  follow,

$$\text{Let } \mathbf{D} = \text{diag}(\gamma \mathbf{A}^T \mathbf{A} + \boldsymbol{\Gamma}), \mathbf{H} = \text{offdiag}(\gamma \mathbf{A}^T \mathbf{A}),$$

$$\mathbf{x}^{(t+1)} = -\mathbf{D}^{-1} \mathbf{H} \mathbf{x}^{(t)} + \mathbf{D}^{-1} \gamma \mathbf{A}^T \mathbf{y}, \quad (19)$$

In (2), we observed that MF can also be implemented as message passing in a factor graph. Hence, it is evident from the above expression that the factor graph representation for SAVE corresponds to the case when all the  $y_i$ 's are treated jointly and all the  $x_i$ 's at the scalar level. Noting that LMMSE estimate of  $\mathbf{x}$  can be written as the solution of  $\mathbf{J} \mathbf{x} = \mathbf{b}$ , with  $\mathbf{J} = \gamma \mathbf{A}^T \mathbf{A} + \boldsymbol{\Gamma}$  and  $\mathbf{b} = \gamma \mathbf{A}^T \mathbf{y}$ . In fact, SAVE corresponds to the Jacobi iterations [18] for solving this linear system with the splitting of  $\mathbf{J} = \mathbf{D} - \mathbf{H}$ , which converges to the true value only if  $\rho(\mathbf{D}^{-1} \mathbf{H}) < 1$ , where  $\rho$  represents the spectral radius. Further, we observe that if we rewrite the SAVE iterations as,  $x_i^{(t+1)} = \sigma_i^2 \mathbf{A}_i^T (\mathbf{y} - \mathbf{A}_{i-} \mathbf{x}_{i-}^{(t+1)} - \mathbf{A}_{i+} \mathbf{x}_{i+}^{(t)}) \gamma$ , where in the update of  $x_i$  at iteration  $(t+1)$  we include the updated values of  $x_k, k = 1, \dots, i-1$ . These updated recursions correspond to Gauss-Siedel method [18] for solving the linear system  $\mathbf{J} \mathbf{x} = \mathbf{b}$ . In Gauss-Siedel version,  $\mathbf{J}$  is split as  $\mathbf{J} = \mathbf{D} - \mathbf{L} - \mathbf{U}$ , where  $\mathbf{L}$  being a matrix which represents the lower triangular portion

of  $\mathbf{H}$  and  $\mathbf{U}$  representing the upper triangular portion. Hence for Gauss-Siedel, the SAVE iterations (19) can be rewritten as,  $\mathbf{x}^{(t+1)} = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{U} \mathbf{x}^{(t)} + (\mathbf{D} - \mathbf{L})^{-1} \gamma \mathbf{A}^T \mathbf{y}$ . Certain remarks on the convergence behaviour follows as below,

#### Remarks:

- From [18], if  $\mathbf{J}$  is an  $M$ -matrix, then Jacobi and Gauss-Siedel iterations for SAVE converge to the true values  $\mathbf{x}^* = \mathbf{J}^{-1} \mathbf{b}$ , for any arbitrary  $\mathbf{b}$ . For  $\mathbf{J}$  to be an  $M$ -matrix, it should be nonsingular and  $\mathbf{A}^{-1} \succeq 0$ . Moreover the off-diagonal elements,  $a_{ij} < 0, \forall i, j, j \neq i$ . Also, the diagonal elements of  $\mathbf{J}$  represented by  $\mathbf{D}$  is nonnegative and nonsingular.
- Another sufficient condition for convergence follows from the diagonal dominance theorem in [18], which says that if  $\mathbf{J}$  is strictly or irreducibly diagonally dominant then  $\hat{\mathbf{x}}$  converges to  $\mathbf{x}^*$ .
- To further accelerate the convergence, one possibility is to employ the successive over-relaxation method (SOR) [18], in which case, the SAVE iterations gets modified as follows.  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \omega(\bar{\mathbf{x}}^{(t+1)} - \mathbf{x}^{(t)})$ , where  $\bar{\mathbf{x}}^{(t+1)}$  corresponds to the Jacobi SAVE iterations (19) or the Gauss-Siedel iterations.
- To fix the convergence of SAVE (when  $\rho(\mathbf{D}^{-1} \mathbf{H}) > 1$ ), we can use the diagonal loading method similar to [19]. The modified iterations (with a diagonal loading factor matrix  $\boldsymbol{\Lambda}$ ) can be written as,

$$(\mathbf{D} + \boldsymbol{\Lambda}) \mathbf{x}^{(t+1)} = -(\mathbf{H} - \boldsymbol{\Lambda}) \mathbf{x}^{(t)} + \gamma \mathbf{A}^T \mathbf{y}, \implies$$

$$\mathbf{x}^{(t+1)} = -(\mathbf{D} + \boldsymbol{\Lambda})^{-1} (\mathbf{H} - \boldsymbol{\Lambda}) \mathbf{x}^{(t)} + (\mathbf{D} + \boldsymbol{\Lambda})^{-1} \gamma \mathbf{A}^T \mathbf{y},$$

The convergence condition gets modified as  $\rho((\mathbf{D} + \boldsymbol{\Lambda})^{-1} (\mathbf{H} - \boldsymbol{\Lambda})) < 1$ . Another point worth noting here is that, if the power delay profile  $\boldsymbol{\Gamma}$  is also estimated using VB as in [11], then we can write  $\mathbf{D} = \gamma \text{diag}(\mathbf{A}^T \mathbf{A}) + \hat{\boldsymbol{\Gamma}}$ , where  $\hat{\boldsymbol{\Gamma}} = \boldsymbol{\Gamma} + \tilde{\boldsymbol{\Gamma}}$ . In this case,  $\tilde{\boldsymbol{\Gamma}}$  may represent an automatic correction factor (diagonal loading) to force convergence of SAVE for cases where  $\rho(\mathbf{D}^{-1} \mathbf{H}) > 1$ .

#### 4.1. Sparsity Analysis with SAVE

In this subsection, we focus on the sparsity analysis of the SAVE iterations described above. We use the approach described in [20, 21], where they compute the stationary point of the precision components  $\alpha_i$ . The expression for mean value of  $\alpha_i$  (for the resulting Gamma posterior from [11]) is,  $\hat{\alpha}_i = \frac{a + \frac{1}{2}}{\left(\frac{\langle x_i^2 \rangle}{2} + b\right)}$ , where,  $\langle x_i^2 \rangle \succeq \hat{x}_i^2 +$

$\sigma_i^2$ . Further substituting for  $\hat{x}_i^2$  in  $\hat{\alpha}_i$ ,

$$\hat{\alpha}_i^{-1} \stackrel{(a)}{=} \frac{\gamma^2}{(\gamma \mathbf{A}_i^T \mathbf{A}_i + \hat{\alpha}_i)^2} [\text{tr}\{\mathbf{y} \mathbf{y}^T \mathbf{A}_i \mathbf{A}_i^T\} + \text{tr}\{\mathbf{A}_i^T \mathbf{A}_i \Sigma_i \mathbf{A}_i \mathbf{A}_i^T\}] +$$

$$\frac{1}{\gamma \mathbf{A}_i^T \mathbf{A}_i + \hat{\alpha}_i},$$

We define  $c_i = \text{tr}\{\mathbf{y} \mathbf{y}^T \mathbf{A}_i \mathbf{A}_i\}, d_i = \text{tr}\{\mathbf{A}_i^T \mathbf{A}_i \Sigma_i \mathbf{A}_i \mathbf{A}_i\}$ , where  $\Sigma_i$  is a diagonal matrix with entries  $\sigma_n^2, \forall n \neq i$ . Also, we made the large system approximation ( $M, N \rightarrow \infty$ ) that  $\mathbf{A}_i^T \mathbf{y} \hat{\mathbf{x}}_i^H \mathbf{A}_i^H \mathbf{A}_i \rightarrow \text{tr}\{E\{\hat{\mathbf{x}}_i^H \mathbf{A}_i^H \mathbf{A}_i \mathbf{A}_i^T \mathbf{y}\}\} = 0$ . After some algebraic manipulations, solving (21) which is of the form  $\alpha_i^{-1} = \mathcal{F}(\alpha_i)$  leads to the following stationary point for  $\alpha_i$ ,

$$\hat{\alpha}_i = \begin{cases} \frac{\gamma(\mathbf{A}_i^T \mathbf{A}_i)^2}{\gamma(c_i + d_i) - \mathbf{A}_i^T \mathbf{A}_i}, & \text{if } \gamma(c_i + d_i) > \mathbf{A}_i^T \mathbf{A}_i \\ \infty, & \text{if } \gamma(c_i + d_i) \leq \mathbf{A}_i^T \mathbf{A}_i \end{cases} \quad (22)$$

The above threshold condition can be intuitively interpreted as follows:  $c_i + d_i$  can be interpreted as the signal power in  $\mathbf{y}' = \mathbf{y} - \mathbf{A}_{i-} \mathbf{x}_{i-}$ . Hence the threshold above checks whether the signal-to-noise ratio of the residual signal (after the matched filtering by  $\mathbf{A}_i$ ) is greater than 1. As observed in [20], this should further accelerate the convergence of the SAVE iterations.

## 5. BAYESIAN SAGE (BSAGE)

In this section, we consider a Bayesian version of the space alternating generalized EM (SAGE) algorithm proposed in [22, 23]. In BSAGE, we consider the estimation of  $x_i$  by fixing the other variables and splitting  $x_k = \hat{x}_k + \tilde{x}_k, \forall k \neq i$ . We define  $\Sigma_{\tilde{i}}$  is the diagonal matrix with entries as the posterior variances  $\sigma_k^2, k \neq i$ . So we write the observation model as,

$$\mathbf{y} - \mathbf{A}_{\tilde{i}} \hat{\mathbf{x}}_{\tilde{i}} \triangleq \mathbf{y}_i = \mathbf{A}_i x_i + \mathbf{A}_{\tilde{i}} \tilde{\mathbf{x}}_{\tilde{i}} + \mathbf{v}, \quad (23)$$

Further we obtain the LMMSE estimate of  $x_i$  as,

$$\begin{aligned} \sigma_i^2 &= \alpha_i + \mathbf{A}_i^T (\mathbf{A}_{\tilde{i}} \Sigma_{\tilde{i}} \mathbf{A}_{\tilde{i}}^T + \frac{1}{\gamma} \mathbf{I}_N)^{-1} \mathbf{A}_i, \\ \hat{x}_i &= \sigma_i^2 \mathbf{A}_i^T (\mathbf{A}_{\tilde{i}} \Sigma_{\tilde{i}} \mathbf{A}_{\tilde{i}}^T + \frac{1}{\gamma} \mathbf{I}_N)^{-1} \mathbf{y}_i \end{aligned} \quad (24)$$

We further define,  $\mathbf{E}_i$  as the diagonal matrix with  $i^{\text{th}}$  entry  $\frac{1}{\alpha_i}$  and rest of the elements same as  $\Sigma$ . Also, define  $\mathbf{V}_i = \mathbf{A} (\mathbf{E}_i^{-1} \gamma^{-1} + \mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ . Further applying matrix inversion lemma [13] and substituting for  $\mathbf{y}_i$ , we obtain,

$$\hat{x}_i = \frac{\gamma}{\alpha_i} \mathbf{A}_i^T \mathbf{y} - \mathbf{A}_i^T \mathbf{V}_i \frac{\gamma}{\alpha_i} \mathbf{y} - \frac{\gamma}{\alpha_i} \mathbf{A}_i^T \mathbf{A}_i \hat{\mathbf{x}}_{\tilde{i}} + \mathbf{A}_i^T \mathbf{V}_i \frac{\gamma}{\alpha_i} \mathbf{A}_i \hat{\mathbf{x}}_{\tilde{i}}. \quad (25)$$

Further, in order to write it in the vector form, we define the matrix  $\mathbf{B}^T$  (of size  $M \times N$ ) with the rows as  $\mathbf{A}_i^T \mathbf{V}_i$ . We obtain the expressions in the vector form as,

$$\begin{aligned} \hat{\mathbf{x}}^{(k+1)} &= -\mathbf{M} \hat{\mathbf{x}}^{(k)} + \mathbf{N} \mathbf{y}, \text{ where, } \mathbf{M} = \gamma \mathbf{\Gamma}^{-1} (\mathbf{H} - \mathbf{L}), \\ \mathbf{L} &= (\mathbf{B}^T \mathbf{A} - \text{diag}(\mathbf{B}^T \mathbf{A})), \mathbf{H} = (\mathbf{A}^T \mathbf{A} - \text{diag}(\mathbf{A}^T \mathbf{A})), \\ \mathbf{N} &= \gamma \mathbf{\Gamma}^{-1} (\mathbf{A} - \mathbf{B})^T. \end{aligned}$$

The per-iteration complexity of BSAGE is also  $\mathcal{O}(M^2 N)$ , hence same as BP. The convergence condition can be written as  $\rho(\mathbf{M}) < 1$ . Further comparing the convergence conditions for SAVE and BSAGE,  $\rho_{\text{SAVE}} = \rho([\gamma \text{diag}(\mathbf{A}^T \mathbf{A}) + \mathbf{\Gamma}]^{-1} \text{offdiag}(\gamma \mathbf{A}^T \mathbf{A}))$  and  $\rho_{\text{BSAGE}} = \rho(\mathbf{\Gamma}^{-1} \text{offdiag}(\gamma (\mathbf{A} - \mathbf{B})^T \mathbf{A}))$ . It can be observed that if  $\mathbf{A}^T \mathbf{A}$  is diagonally dominant (which is also one of the conditions for the convergence of SAVE to the true means), then the effect of the offdiagonal terms of  $(\mathbf{A} - \mathbf{B})^T \mathbf{A}$  or  $\mathbf{A}^T \mathbf{A}$  is negligible and the dominating factor is the first term in the expression of  $\rho$ . Since  $[\gamma \text{diag}(\mathbf{A}^T \mathbf{A}) + \mathbf{\Gamma}]^{-1} < \mathbf{\Gamma}^{-1}$ , we can conclude that  $\rho_{\text{SAVE}} < \rho_{\text{BSAGE}}$  explaining the faster convergence of SAVE as noted in [11] and [4].

## 6. CONCLUSIONS

Motivated by the need for low complexity solutions for sparse signal recovery, we looked at various approximate inference techniques for SBL whose complexity is of the order of the length of the sparse signal. In this paper, we attempt to provide convergence analysis for SBL under approximate inference techniques such as VB, BP or EP. However, much remains to be done. The convergence values of the posterior variances for BP still needs to be understood. One possible future direction is to analyze the convergence behaviour with estimated hyperparameters. Another extension of the present work is when the dictionary matrix is unknown, for e.g. structured dictionary matrices as in [24, 25].

## 7. REFERENCES

- [1] M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine," *J. Mach. Learn. Res.*, vol. 1, 2001.
- [2] D. P. Wipf and B. D. Rao, "Sparse Bayesian Learning for Basis Selection," *IEEE Trans. on Sig. Proc.*, Aug. 2004.
- [3] X. Tan and J. Li, "Computationally efficient sparse Bayesian learning via belief propagation," *IEEE Trans. on Sig. Proc.*, vol. 58, no. 4, Apr. 2013.
- [4] C. K. Thomas and D. Slock, "Low complexity static and dynamic sparse bayesian learning combining BP, VB and EP message passing," in *Asilomar Conf. on Sig., Sys., and Comp.*, CA, USA, 2019.

- [5] J. Du et al., "Convergence Analysis of Distributed Inference with Vector-Valued Gaussian Belief Propagation," *Jrnl. of Mach. Learn. Res.*, April 2018.
- [6] J. Du et al., "Convergence analysis of the information matrix in Gaussian Belief Propagation," in *IEEE Intl. Conf. on Acoustics, Speech, and Sig. Process.*, New Orleans, LA, USA, 2017.
- [7] Q. Su and Y. Wu, "Convergence Analysis of the Variance in Gaussian Belief Propagation," *IEEE Trans. on Sig. Process.*, Oct. 2014.
- [8] B. Cseke and T. Heskes, "Properties of Bethe Free Energies and Message Passing in Gaussian Models," *Jrnl. of Art. Intell. Res.*, May 2011.
- [9] K. P. Murphy et al., "Loopy belief propagation for approximate inference: an empirical study," in *In 15th Conf. Uncert. in Art. Intell. (UAI)*, Stockholm, Sweden, 1999.
- [10] D. M. Malioutov, J. K. Johnson, and A. S. Willsky, "Walk-Sums and Belief Propagation in Gaussian Graphical Models," *Jrnl. of Mach. Learn. Res.*, Oct. 2006.
- [11] C. K. Thomas and D. Slock, "SAVE - space alternating variational estimation for sparse Bayesian learning," in *Data Science Workshop*, 2018.
- [12] C. K. Thomas and D. Slock, "Gaussian variational Bayes Kalman filtering for dynamic sparse Bayesian learning," in *ITISE*, 2018.
- [13] S. Wagner et al., "Large system analysis of linear precoding in MISO broadcast channels with limited feedback," *IEEE Trans. Inf. Theory*, vol. 58, no. 7, pp. 4509–4538, July 2012.
- [14] M. J. Beal, "Variational algorithms for approximate bayesian inference," in *Thesis, Univeristy of Cambridge, UK*, May 2003.
- [15] T. P. Minka, "Expectation propagation for approximate Bayesian inference," in *in Proc. of Conf. on Uncert. in Art. Intell. (UAI)*, San Francisco, CA, USA, 2001.
- [16] S. Rangan et al., "On the Convergence of Approximate Message Passing With Arbitrary Matrices," *IEEE Trans. on Info. Theo.*, Sept. 2019.
- [17] Mohsen Bayati and Andrea Montanari, "The Dynamics of Message Passing on Dense Graphs, with Applications to Compressed Sensing," *IEEE Trans. on Inf. Theory*, vol. 57, no. 2, February 2011.
- [18] J. M. Ortega, "Numerical Analysis: A Second Course," in *SIAM*, Philadelphia, 1990.
- [19] J. K. Johnson et al., "Fixing Convergence of Gaussian Belief Propagation," in *IEEE Intl. Symp. on Info. Theo.*, 2009.
- [20] D. Shutin et al., "Fast variational sparse bayesian learning with automatic relevance determination for superimposed signals," *IEEE Trans. on Sig. Process.*, vol. 59, no. 12, December 2011.
- [21] Michael E. Tipping and Anita C. Faul, "Fast marginal likelihood maximisation for sparse Bayesian models," in *AISTATS*, January 2003.
- [22] J. A. Fessler and A. Hero, "Space-alternating generalized expectation-maximization algorithm," *IEEE Trans. on Sig. Process.*, Oct. 1994.
- [23] Bernard H. Fleury et al., "Channel Parameter Estimation in Mobile Radio Environments Using the SAGE Algorithm," *IEEE J. on Sel. Areas in Commun.*, vol. 17, no. 3, March 1999.
- [24] C. K. Thomas and D. Slock, "Space Alternating Variational Estimation and Kronecker Structured Dictionary Learning," in *ICASSP*, 2019.
- [25] C. K. Thomas and D. Slock, "SAVED - Space alternating variational estimation for sparse Bayesian learning with parametric dictionaries," in *Asilomar Conf. on Sig., Sys., and Comp.*, 2018.