

Reconciling AMP Algorithms derived from Belief Propagation or the Large System Limit Bethe Free Energy

Zilu Zhao¹, Fangqing Xiao¹, Christo K. Thomas² and Dirk Slock¹

¹Communication Systems Department, EURECOM, France

²NEWS Lab, Virginia Tech, VA, USA

{Zilu.Zhao, Fangqing.Xiao, Dirk.Slock}@eurecom.fr, ChristoKT@vt.edu

Abstract—When derived from the Bethe Free Energy (BFE) of the Generalized Linear Model (GLM), Approximate Message Passing (AMP) algorithms combine two asymptotic Large System Limit (LSL) simplifications which are asymptotic Gaussianity of extrinsics and large random matrix theory based asymptotic variance computations. In the provably convergent AMBGAMP algorithm, a LSL version of the BFE is derived. In Expectation Propagation (EP) style minimization, the LSL BFE cost function is augmented with Lagrangian terms for mean and variance consistency constraints, augmented with a quadratic version of the mean constraints as in the Method of Multipliers (MM). The mean Lagrange multipliers then get updated ADMM-style (Alternating Direction of MM). In this approach, the weights of the MM terms need to be carefully chosen, which is not part of the MM philosophy, and the Lagrange multipliers have no particular meaning. On the other hand, AMP can be derived by directly introducing LSL simplifications in the Belief Propagation (BP) algorithm that minimizes the original GLM BFE. This allows to relate extrinsic messages to posterior pdfs by first-order Taylor series expansion based perturbations. We also apply LSL approximations to the variances of the various Gaussians involved, which in fact leads to a rederivation of a fundamental LSL theorem describing the deterministic limit of posterior variances. We show that this LSL version of BP leads to BFE modifications that correspond to the augmented Lagrangian of the LSL BFE, explaining its weights and Lagrange Multipliers. These insights should facilitate the extension of AMP to more complex settings such as bilinear models.

I. INTRODUCTION

Sparse signal recovery is a fundamental problem in signal processing with a wide range of applications. Many of these problems can be framed as the task of estimating a latent vector \mathbf{x} based on a correlated observation vector \mathbf{y} [1]. In the Bayesian framework, the complexity of Canonical Methods such as MMSE and MAP experiences exponential growth as the dimension of the problem grows.

By exploiting the structure of the models, graphical model based methods prove to be effective. Belief Propagation (BP) transforms the global inference problem into a local inference problem as outlined by [2]. Expectation Propagation (EP) was introduced in [3] and has been shown to share a similar updating scheme as BP, but for computational efficiency, the messages in BP are projected into a suitable member of the family of exponential distributions [3].

In both [1] and [4], the authors unify EP and BP within the framework of minimizing variational free energy. They demonstrate the close relationship between the fixed points of various message-passing algorithms and the stationary points of Bethe Free Energy (BFE).

EP can serve as an inference method in the generalized linear model (GLM). However, the computational cost corresponds to propagating $2MN$ messages as in Fig. 1 when the data matrix \mathbf{A} is of size $M \times N$. Generalized Approximate Mes-

sage Passing (GAMP) [5] builds upon EP, but through the application of large system approximations (LSA), it effectively reduces the number of messages to $M + N$ extrinsics and (marginal) posteriors, providing a more computationally efficient approach.

In [6], the authors investigated the fixed points of the Generalized AMP (GAMP) algorithm for GLMs. They discovered that GAMP shares the same fixed point as the stationary points of the Large System Limit Bethe Free Energy (LSL BFE). In [7] we then proposed AMBGAMP which is guaranteed to converge. Building upon the works of [1], [8], [7], [9], [10], and [11], we present the contributions described in the abstract.

II. BETHE FREE ENERGY OF THE GENERALIZED LINEAR MODEL

A. Bethe Free Energy (BFE)

Consider a pdf factorization

$$p(\mathbf{x}, \mathbf{y}) \propto \prod_{\alpha} f_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha}), \quad (1)$$

where \mathbf{x}_{α} is a subvector of \mathbf{x} . In case of a tree-structured factor graph, an alternative equivalent form is [2]

$$p(\mathbf{x}|\mathbf{y}) = \frac{\prod_{\alpha} p(\mathbf{x}_{\alpha})}{\prod_i p(x_i)^{M_i-1}}, \quad (2)$$

where M_i is the number of subvectors \mathbf{x}_{α} that contain x_i . In (2), the $p(\mathbf{x}_{\alpha})$ and $p(x_i)$ are the exact factor (subvector) resp. variable marginals.

The concept of variational free energy suggests that to infer the marginals from a tree structured $p(\mathbf{x}, \mathbf{y})$ given in (1), we can use as trial distribution

$$q_{\mathbf{x}}(\mathbf{x}) = \frac{\prod_{\alpha} q_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha})}{\prod_i q_{x_i}(x_i)^{M_i-1}}. \quad (3)$$

The true marginals can be obtained by [1]

$$\begin{aligned} \min_{q_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha}), q_{x_i}(x_i)} F &= D[q(\mathbf{x}) \| \prod_{\alpha} f_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha})]; \\ \text{s.t. } \forall \alpha, \forall i \in \mathcal{I}_{\alpha}, q_{x_i}(x_i) &= \int q_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha}) d\mathbf{x}_{\bar{i}}, \end{aligned} \quad (4)$$

where we define the shorthand notation (for arbitrary nonnegative functions q, p) $D(q\|p) = \int q(x) \ln \frac{q(x)}{p(x)} dx$ (which is the Kullback-Leibler Divergence (KLD) in case of normalized q, p) and $\mathbf{x}_{\bar{i}}$ denotes all \mathbf{x} except x_i . The free energy can be expanded as

$$F = \sum_{\alpha} D[q_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha}) \| f_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha})] + \sum_i (M_i - 1) H[q_{x_i}(x_i)], \quad (5)$$

where $H(\cdot)$ denotes entropy in nats. Note that this representation only holds for a tree structured distribution. For general graphs that contain loops, (2) no longer holds. Thus, in cases with loops, (5) is only an approximation of the variational free energy. The expression (5) is instead called Bethe free energy.

B. BFE of the GLM for BP

We consider a GLM with

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i), \mathbf{z} = \mathbf{A}\mathbf{x}, p(\mathbf{y}|\mathbf{z}) = \prod_{j=1}^M p(y_j|z_j), \quad (6)$$

where the ratio N/M is a constant for large system considerations. We interpret the linear mixing as a conditional probability $p(\mathbf{z}|\mathbf{x}) = \delta(\mathbf{z} - \mathbf{A}\mathbf{x})$.

From this general linear model, a joint (loopy) factorization scheme comes up naturally:

$$p(\mathbf{x}, \mathbf{z}|\mathbf{y}) \propto p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{y}|\mathbf{z})\delta(\mathbf{z} - \mathbf{A}\mathbf{x})p(\mathbf{x}). \quad (8)$$

According to the definition of BFE (5), the associated BFE based on the joint factorization scheme (8) is calculated [1] as

$$F = D[q_{\mathbf{x}}(\mathbf{x})||p(\mathbf{x})] + D[q_{\mathbf{z}}(\mathbf{z})||p(\mathbf{y}|\mathbf{z})] + \sum_i H[q_{x_i}(x_i)] + D[b_{\mathbf{x},\mathbf{z}}(\mathbf{x}, \mathbf{z})||\delta(\mathbf{z} - \mathbf{A}\mathbf{x})] + \sum_j H[q_{z_j}(z_j)], \quad (9)$$

where $q_{\mathbf{x}}$, $q_{\mathbf{z}}$, $b_{\mathbf{x},\mathbf{z}}$, q_{x_i} and q_{z_j} are only approximate posteriors because of the loops in the factor graph. Since we need to minimize the BFE given by (9), the distribution function $b_{\mathbf{x},\mathbf{z}}(\mathbf{x}, \mathbf{z})$ must be of the form

$$b_{\mathbf{x},\mathbf{z}}(\mathbf{x}, \mathbf{z}) = b_{\mathbf{x}}(\mathbf{x})\delta(\mathbf{z} - \mathbf{A}\mathbf{x}), \quad (10)$$

to avoid an infinite value of the KLD, leading to $D[b_{\mathbf{x},\mathbf{z}}(\mathbf{x}, \mathbf{z})||\delta(\mathbf{z} - \mathbf{A}\mathbf{x})] = -H[b_{\mathbf{x}}]$. For BP, the BFE (9) needs to be minimized w.r.t. marginal consistency constraints $q_{\mathbf{x}}(x_i) = b_{\mathbf{x}}(x_i) = q_{x_i}(x_i)$, $q_{\mathbf{z}}(z_j) = q_{z_j}(z_j)$. Given the independent priors for \mathbf{x} , \mathbf{z} , minimization of the BFE leads to $q_{\mathbf{x}}(\mathbf{x}) = \prod_i q_{x_i}(x_i)$, $q_{\mathbf{z}}(\mathbf{z}) = \prod_j q_{z_j}(z_j)$. Furthermore, the maximization of $H[b_{\mathbf{x}}]$ under marginal constraints leads to $b_{\mathbf{x}}(\mathbf{x}) = \prod_i b_{x_i}(x_i)$. Together with the marginal constraints, this leads to the cancellation of the entropy terms in \mathbf{x} in the BFE, which becomes $F =$

$$\sum_i D[q_{x_i}(x_i)||p(x_i)] + \sum_j D[q_{z_j}(z_j)||p(y_j|z_j)] + \sum_j H[q_{z_j}(z_j)] \quad (11)$$

which needs to be minimized under the constraint $\mathbf{z} = \mathbf{A}\mathbf{x}$.

III. GAMP FROM LSL BELIEF PROPAGATION

In reGVAMP [12], [10], extrinsics in the GLM are built from the *equivalent Gaussian linear model*, which introduces *equivalent Gaussian priors* from Gaussian posterior approximations and Gaussian extrinsics.

GAMP exploits LSL simplifications of reGVAMP for a random \mathbf{A} with i.i.d. signs which leads to

- (i) Gaussianity of extrinsics (also in reGVAMP), and
 - (ii) independence of marginals (extra w.r.t. reGVAMP).
- (ii) leads to the large system simplifications of the variances, avoiding covariance matrix inverses. But also posterior and extrinsic estimates $\hat{\mathbf{x}}$, $\hat{\mathbf{z}}$ and \mathbf{r} , \mathbf{p} that are constructed by combining decoupled pieces of information. These estimates are non-linear MMSE and CWCU MMSE estimates in

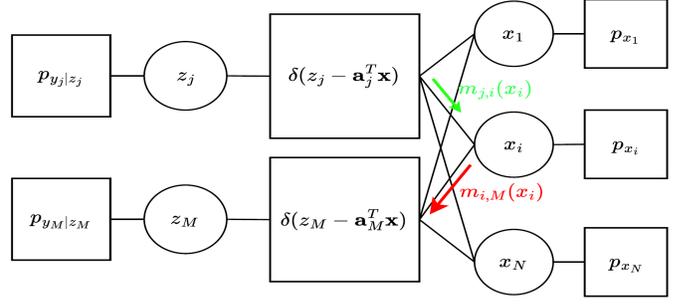


Fig. 1. Factor Graph for the GLM used by GAMP.

general. Extrinsic estimates are not obtained as linear perturbations of corresponding MMSE estimates because those are not necessarily close to each other. Rather the interplay between \mathbf{x} and \mathbf{z} is exploited with perturbations due to the small effect of a single term in \mathbf{A} in the LSL. In both reGVAMP and GAMP, we have:

Gaussian extrinsics: $e_{\mathbf{x}}(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{r}, \boldsymbol{\tau}_r)$, $e_{\mathbf{z}}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{p}, \boldsymbol{\tau}_p)$ and

Posterior marginals proportional to: $q_{\mathbf{x}}(\mathbf{x}) \sim p_{\mathbf{x}}(\mathbf{x})e_{\mathbf{x}}(\mathbf{x})$, $q_{\mathbf{z}}(\mathbf{z}) \sim p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z})e_{\mathbf{z}}(\mathbf{z})$ with Gaussian approximations $\mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \boldsymbol{\tau}_x)$, $\mathcal{N}(\mathbf{z}; \hat{\mathbf{z}}, \boldsymbol{\tau}_z)$ resp. (where \mathbf{r} , $\boldsymbol{\tau}_r$, \mathbf{p} , $\boldsymbol{\tau}_p$ etc. are vectors, e.g. $\mathcal{N}(\mathbf{x}; \mathbf{r}, \boldsymbol{\tau}_r)$ is short for $\mathcal{N}(\mathbf{x}; \mathbf{r}, \mathbf{D}_{\boldsymbol{\tau}_r})$, and $\mathbf{D}_{\boldsymbol{\tau}_r} = \text{diag}(\boldsymbol{\tau}_r)$). In [11], it is shown that LSL simplifications of BP lead to the GAMP algorithm (see e.g. [7]) which computes the indicated marginal posteriors that minimize the cost function [8] (to be minimized w.r.t. $q_{\mathbf{x}}$, $q_{\mathbf{z}}$ and \mathbf{u} later)

$$D(q_{\mathbf{x}}||p_{\mathbf{x}}e_{\mathbf{x}}/Z_{\mathbf{x}}) + D(q_{\mathbf{z}}||p_{\mathbf{z}}e_{\mathbf{z}}/Z_{\mathbf{z}}). \quad (12)$$

We get per component

$$\min_{q_{x_i}} D(q_{x_i}||g_{x_i}/Z_{x_i}) \Rightarrow q_{x_i} = g_{x_i}/Z_{x_i}, Z_{x_i} = \int g_{x_i}(x_i) dx_i, -\ln g_{x_i}(x_i) = f_{x_i}(x_i) + \frac{1}{2\tau_{r_i}}[(x_k - r_i)^2 - r_i^2]. \quad (13)$$

The partition function Z_{x_i} acts as cumulant generating function:

$$\tau_{r_i} \frac{\partial \ln Z_{x_i}}{\partial r_i} = \mathbb{E}(x_i|q_{x_i}) = \mathbb{E}(x_i|r_i, \tau_{r_i}) = \hat{x}_i, \tau_{r_i}^2 \frac{\partial^2 \ln Z_{x_i}}{\partial r_i^2} = \text{var}(x_i|r_i, \tau_{r_i}) = \tau_{x_i}. \quad (14)$$

We also get per component

$$\min_{q_{z_k}} D(q_{z_k}||g_{z_k}/Z_{z_k}) \Rightarrow q_{z_k} = g_{z_k}/Z_{z_k}, Z_{z_k} = \int g_{z_k}(z_k) dz_k, -\ln g_{z_k}(z_k) = f_{z_k}(z_k) + \frac{1}{2\tau_{p_k}^2}[(z_k - p_k)^2 - p_k^2]. \quad (15)$$

The partition function Z_{z_k} acts again as cumulant generating function:

$$-\frac{\partial \ln Z_{z_k}}{\partial p_k} = \mathbb{E}(z_k|q_{z_k}) = \mathbb{E}(z_k|p_k, \tau_{p_k}) = \hat{z}_k, \frac{\partial^2 \ln Z_{z_k}}{\partial p_k^2} = \text{var}(z_k|p_k, \tau_{p_k}) = \tau_{z_k}. \quad (16)$$

The LSL BP derivation also leads to the following identities

$$Z_z(p, y, \tau_p) = \int p_{y|z}(y|z) e^{-\frac{1}{2\tau_p}(z-p)^2} dz, \frac{\partial \ln Z_z}{\partial p} = \frac{Z'_z}{Z_z} = s = \frac{\hat{z}-p}{\tau_p}, \hat{z} = \frac{1}{Z_z} \int z p_{y|z}(y|z) e^{-\frac{1}{2\tau_p}(z-p)^2} dz, \frac{\partial^2 \ln Z_z}{\partial p^2} = -\tau_s = \frac{Z''_z}{Z_z} - \left(\frac{Z'_z}{Z_z}\right)^2 = -(1 - \tau_z/\tau_p)/\tau_p$$

and updates of the following quantities

$$\begin{aligned} \mathbf{p} &= \mathbf{A} \hat{\mathbf{x}} - \tau_k^p \cdot \mathbf{s}, \quad \tau_p - \mathbf{S} \tau_x \\ \mathbf{r} &= \hat{\mathbf{x}} + \tau_r \cdot \mathbf{A}^T \mathbf{s}, \quad \tau_r = \mathbf{1} / (\mathbf{S}^T \tau_s) \end{aligned} \quad (17)$$

where $\mathbf{S} = \mathbf{A} \cdot \mathbf{A}$ and we use the notations: $\|\mathbf{u}\|_{\tau}^2 = \sum_i u_i^2 / \tau_i$, element-wise multiplication as in $\mathbf{s} \cdot \tau$ and element-wise division as in $\mathbf{1} / \tau$, and $\mathbf{1}$ is a vector of ones.

IV. LSL BFE AND EP

After the LSL simplifications [11], the BFE from (11) with marginal pdf consistency constraints can be seen to become equivalent to the following LSL-BFE [7], [9] :

$$\begin{aligned} \min_{q_x, q_z, \tau_p, \mathbf{u}} D[q_x \| p_x] + D[q_z \| p_{y|z}] + \frac{1}{2} \sum_k \left[\frac{\text{var}(z_k | q_z)}{\tau_{pk}} + \ln(\tau_{pk}) \right] \\ \text{s.t. } E[\mathbf{z} | q_z] = \mathbf{A} \mathbf{u} \\ E[\mathbf{x} | q_x] = \mathbf{u} \\ \tau_p = \mathbf{S} \text{var}(\mathbf{x} | q_x). \end{aligned} \quad (18)$$

We will exploit some useful relations

$$\begin{aligned} \forall \tau, \mathbf{c}^T \text{var}(\mathbf{x} | q_x) &= \int \|\mathbf{x} - \mathbf{u}\|_{\mathbf{1} / \tau}^2 q_x(\mathbf{x}) d\mathbf{x} \\ \sum_k \frac{\text{var}(z_k | q_z)}{\tau_{pk}} &= \int \|\mathbf{z} - \mathbf{A} \mathbf{u}\|_{\tau_p}^2 q_z(\mathbf{z}) d\mathbf{z}. \end{aligned} \quad (19)$$

The Lagrangian of (18) becomes

$$\begin{aligned} L &= D[q_x \| p_x] + D[q_z \| p_{y|z}] + \frac{1}{2} \sum_k \left[\frac{\text{var}(z_k | q_z)}{\tau_{pk}} + \ln(\tau_{pk}) \right] \\ &+ \lambda_{\mu_x}^T \left(\mathbf{A} \mathbf{u} - \int \mathbf{z} q_z(\mathbf{z}) d\mathbf{z} \right) + \lambda_{\mu_z}^T \left(\mathbf{u} - \int \mathbf{x} q_x(\mathbf{x}) d\mathbf{x} \right) \\ &- \frac{1}{2} \lambda_{\tau}^T (\tau_p - \mathbf{S} \text{var}(\mathbf{x} | q_x)) \end{aligned} \quad (20)$$

The derivatives w.r.t. $q_x, q_z, \tau_p, \mathbf{u}$ become

$$\begin{aligned} \frac{\partial L}{\partial q_x} &= \ln(q_x) - \ln(p_x) - \lambda_{\mu_x}^T \mathbf{x} + \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_{\mathbf{1} / (\mathbf{S}^T \lambda_{\tau})}^2 \\ \frac{\partial L}{\partial q_z} &= \ln(q_z) - \ln(p_{y|z}) - \lambda_{\mu_z}^T \mathbf{z} + \frac{1}{2} \|\mathbf{z} - \mathbf{A} \mathbf{u}\|_{\tau_p}^2 \\ \frac{\partial L}{\partial \tau_{pk}} &\propto -\frac{\text{var}(z_k | q_z)}{\tau_{pk}^2} + \frac{1}{\tau_{pk}} - \lambda_{\tau_k} \\ \frac{\partial L}{\partial \mathbf{u}} &= -\mathbf{A}^T \mathbf{D}_{\tau_p}^{-1} (\hat{\mathbf{z}} - \mathbf{A} \mathbf{u}) + \mathbf{A}^T \lambda_{\mu_z} + \lambda_{\mu_x} \\ &- \mathbf{D}_{\mathbf{S}^T \lambda_{\tau}} (\hat{\mathbf{x}} - \mathbf{u}), \end{aligned} \quad (21)$$

where $\hat{\mathbf{z}} = \mathbb{E}[\mathbf{z} | q_z]$ and $\hat{\mathbf{x}} = \mathbb{E}[\mathbf{x} | q_x]$. Zeroing derivatives:

$$q_x(\mathbf{x}) \propto p_x(\mathbf{x}) e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_{\mathbf{1} / (\mathbf{S}^T \lambda_{\tau})}^2} e^{\lambda_{\mu_x}^T \mathbf{x}} \quad (22)$$

$$q_z(\mathbf{z}) \propto p_{y|z}(\mathbf{y} | \mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{A} \mathbf{u}\|_{\tau_p}^2} e^{\lambda_{\mu_z}^T \mathbf{z}} \quad (23)$$

$$\lambda_{\tau_k} = \frac{1}{\tau_{pk}} - \frac{\tau_{z_k}}{\tau_{pk}^2} \quad (24)$$

$$\begin{aligned} \left[\mathbf{A}^T \mathbf{D}_{\tau_p}^{-1} \mathbf{A} + \mathbf{D}_{\mathbf{S}^T \lambda_{\tau}} \right] \mathbf{u} &= \mathbf{A}^T \mathbf{D}_{\tau_p}^{-1} \hat{\mathbf{z}} \\ &+ \mathbf{D}_{\mathbf{S}^T \lambda_{\tau}} \hat{\mathbf{x}} - \mathbf{A}^T \lambda_{\mu_z} - \lambda_{\mu_x} \end{aligned} \quad (25)$$

where $\tau_{z_k} = \mathbb{E}[(z_k - \hat{z}_k)^2 | q_z]$. By satisfying the two mean constraints in (18), the equation (25) becomes

$$\mathbf{A}^T \lambda_{\mu_z} = -\lambda_{\mu_x} \quad (26)$$

A solution can be obtained by solving the system of 7 equations containing (22), (23), (24), (26) and the three constraint equations in (18).

V. ITERATIVE SOLUTION LEADING TO GAMP

We ignore pdf normalization for simplicity. Furthermore, we use red symbols to indicate parameters to be updated.

A. Update of λ_{μ_z}

Consider (23) and the two mean constraints in (18)

$$\mathbb{E} \left[\mathbf{z} | p_{y|z}(\mathbf{y} | \mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \hat{\mathbf{z}}\|_{\tau_p}^2} e^{\lambda_{\mu_z}^{(\text{new})T} \mathbf{z}} \right] \quad (27)$$

$$= \mathbb{E} \left[\mathbf{z} | p_{y|z}(\mathbf{y} | \mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{A} \hat{\mathbf{x}}\|_{\tau_p}^2} e^{\lambda_{\mu_z}^T \mathbf{z}} \right] = \hat{\mathbf{z}} \quad (28)$$

We first use (28) to obtain $\hat{\mathbf{z}}$. Then we use this newly obtained $\hat{\mathbf{z}}$ to update $\lambda_{\mu_z}^{(\text{new})}$, since we need to keep the exponential factor identical in order not to change the mean, i.e.,

$$e^{-\frac{1}{2} \|\mathbf{z} - \hat{\mathbf{z}}\|_{\tau_p}^2} e^{\lambda_{\mu_z}^{(\text{new})T} \mathbf{z}} = e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{A} \hat{\mathbf{x}}\|_{\tau_p}^2} e^{\lambda_{\mu_z}^T \mathbf{z}}. \quad (29)$$

If we want to bridge the GAMP from [7] and BFE, we can denote

$$\mathbf{p} = \mathbf{A} \hat{\mathbf{x}} + \mathbf{D}(\tau_p) \lambda_{\mu_z}. \quad (30)$$

With definition (30), the expression (28) can be written as

$$\hat{\mathbf{z}} = \mathbb{E} \left[\mathbf{z} | p_{y|z}(\mathbf{y} | \mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{p}\|_{\tau_p}^2} \right]. \quad (31)$$

Therefore, the updating for $\lambda_{\mu_z}^{(\text{new})}$ according to (29) becomes

$$\lambda_{\mu_z}^{(\text{new})} = \mathbf{D}(\tau_p^{-1}) (\mathbf{p} - \hat{\mathbf{z}}). \quad (32)$$

It is now clear that we can relate to the LSL BP GAMP above (or [7]) if we define

$$\mathbf{s} = -\lambda_{\mu_z}. \quad (33)$$

For further use, we also state the computation of $\tau_{\hat{\mathbf{z}}}$ explicitly:

$$\tau_{\hat{\mathbf{z}}} = \mathbb{E} \left[(\mathbf{z} - \hat{\mathbf{z}})^2 | p_{y|z}(\mathbf{y} | \mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{p}\|_{\tau_p}^2} \right] \quad (34)$$

$$= \mathbb{E} \left[(\mathbf{z} - \hat{\mathbf{z}})^2 | p_{y|z}(\mathbf{y} | \mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \hat{\mathbf{z}}\|_{\tau_p}^2} e^{\lambda_{\mu_z}^{(\text{new})T} \mathbf{z}} \right] \quad (35)$$

where \mathbf{z}^2 denotes element-wise square of vector \mathbf{z} . (34) and (35) result in the same solution.

B. Update of λ_{μ_x}

According to (26), we can update λ_{μ_x} by

$$\lambda_{\mu_x} = -\mathbf{A}^T \lambda_{\mu_z}. \quad (36)$$

To show the relation between this paper and [7], we define

$$\begin{aligned} \tau_{\mathbf{r}} &= \mathbf{1} / (\mathbf{S}^T \lambda_{\tau}) \\ \mathbf{r} &= \hat{\mathbf{x}}^{\text{old}} + \mathbf{D}_{\tau_{\mathbf{r}}} \mathbf{A}^T \mathbf{s} \end{aligned} \quad (37)$$

Then the updated posterior mean and variance becomes

$$\begin{aligned} \hat{\mathbf{x}} &= \mathbb{E} \left[\mathbf{x} | p_x(\mathbf{x}) e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_{\tau_{\mathbf{r}}}} \right] \\ \tau_{\hat{\mathbf{x}}} &= \mathbb{E} \left[(\mathbf{x} - \hat{\mathbf{x}})^2 | p_x(\mathbf{x}) e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_{\tau_{\mathbf{r}}}} \right], \end{aligned} \quad (38)$$

where we also used the mean constraint for \mathbf{x} in (18).

C. The update of λ_τ and τ_p

The updates of these two variables are quite straightforward. They are already explicitly given by (24) and the variance constraint in (18). To show the relation with GAMP in [7] explicitly, we can define τ_s and then get

$$\tau_s = \lambda_\tau \text{ from which } \tau_p = \mathbf{S}\tau_x, \tau_{s_k} = \frac{1}{\tau_{p_k}} - \frac{\tau_{z_k}}{\tau_{p_k}^2}. \quad (39)$$

VI. ITERATIVE SOLUTION LEADING TO AMBGAMP

GAMP does not use the extra variable \mathbf{u} in (18) (hence uses $\mathbf{u} = \hat{\mathbf{x}}$) and as result is an algorithm that does not correspond to alternating optimization of a BFE, with the resulting convergence issues. For AMBGAMP, we keep variable \mathbf{u} , and use (22)-(26) along with the three constraints in (18) as a system of 8 equations to be solved.

A. Update of λ_{μ_z}

Consider (23) and the mean constraint in (18)

$$\mathbb{E} \left[\mathbf{z} | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \hat{\mathbf{z}}\|_{\tau_p}^2} e^{\lambda_{\mu_z}^{(\text{new})T} \mathbf{z}} \right] \quad (40)$$

$$= \mathbb{E} \left[\mathbf{z} | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{A}\mathbf{u}\|_{\tau_p}^2} e^{\lambda_{\mu_z}^T \mathbf{z}} \right] = \hat{\mathbf{z}} \quad (41)$$

To make the connection with AMBGAMP in [7], we define

$$\mathbf{p} = \mathbf{A}\mathbf{u} + \tau_p \cdot \lambda_{\mu_z} \quad (42)$$

Similar to the previous section, we have the update

$$\lambda_{\mu_z}^{(\text{new})} = (\mathbf{p} - \hat{\mathbf{z}}) / \tau_p. \quad (43)$$

Substitute (42) into (43), and we have

$$\lambda_{\mu_z}^{(\text{new})} = \lambda_{\mu_z} + (\mathbf{A}\mathbf{u} - \hat{\mathbf{z}}) / \tau_p. \quad (44)$$

If we define

$$\mathbf{s} = -\lambda_{\mu_z}, \quad (45)$$

it then follows

$$\mathbf{s}^{(\text{new})} = \mathbf{s} + (\hat{\mathbf{z}} - \mathbf{A}\mathbf{u}) / \tau_p. \quad (46)$$

For the convenience of the further discussion, we write the update for the posterior mean and variance of \mathbf{z}

$$\hat{\mathbf{z}} = \mathbb{E} \left[\mathbf{z} | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{p}\|_{\tau_p}^2} \right] \quad (47)$$

$$\tau_{\hat{\mathbf{z}}} = \mathbb{E} \left[(\mathbf{z} - \hat{\mathbf{z}})^2 | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{p}\|_{\tau_p}^2} \right]. \quad (48)$$

B. Update of λ_{μ_w}

We can use (26) and (45) to obtain the update

$$\lambda_{\mu_w} = -\mathbf{A}^T \lambda_{\mu_z} = \mathbf{A}^T \mathbf{s} \quad (49)$$

If we define (and note that $\lambda_\tau = \tau_s$)

$$\tau_r = \mathbf{1} / (\mathbf{S}^T \lambda_\tau), \mathbf{r} = \mathbf{u} + \tau_r \cdot \lambda_{\mu_w} = \mathbf{u} + \tau_r \cdot (\mathbf{A}^T \mathbf{s}) \lambda_{\mu_w}, \quad (50)$$

we have the explicit update for $\hat{\mathbf{x}}$ and $\tau_{\hat{\mathbf{x}}}$:

$$\begin{aligned} \hat{\mathbf{x}} &= \mathbb{E} \left[\mathbf{x} | p_{\mathbf{x}}(\mathbf{x}) e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_{\tau_r}^2} \right] \\ \tau_{\hat{\mathbf{x}}} &= \mathbb{E} \left[(\mathbf{x} - \hat{\mathbf{x}})^2 | p_{\mathbf{x}}(\mathbf{x}) e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_{\tau_r}^2} \right]. \end{aligned} \quad (51)$$

C. Update of \mathbf{u}

By combining (25) and (26), we get the solution

$$\mathbf{u} = \left[\mathbf{A}^T \mathbf{D}_{\tau_p}^{-1} \mathbf{A} + \mathbf{D}_{\tau_r}^{-1} \right]^{-1} (\mathbf{A}^T \mathbf{D}_{\tau_p}^{-1} \hat{\mathbf{z}} + \mathbf{D}_{\tau_r}^{-1} \hat{\mathbf{x}}). \quad (52)$$

AMBGAMP actually updates \mathbf{u} by applying SGD with step-size by linesearch to the quadratic cost function that (52) is the solution of. The update of λ_τ and τ_p are identical to the updates in GAMP in (39).

VII. CONCLUDING REMARKS

In this paper, we have shown that it is possible to derive the convergent AMBGAMP algorithm by analyzing the KKT conditions for optimizing the LSL BFE. And this while avoiding the quadratic augmentation terms of the Method of Moments, which require a very particular choice in their weights, and circumventing the ADMM-style update of a Lagrange multiplier. This is thanks to the introduction of the auxiliary variable \mathbf{u} in the mean consistency constraints. This \mathbf{u} is optimized to minimize the BFE equivalent in (12) and can be interpreted to be a MMSE estimate of an equivalent underlying Gaussian linear model. On the other hand, another solution to the LSL BFE, which eliminates \mathbf{u} via $\mathbf{u} = \hat{\mathbf{x}}$, leads to GAMP and corresponds to the original LSL BP based derivation, optimizing BFE with LSL approximations. Hence we have reconciled these seemingly different approaches.

The variance predictions in (AMB)GAMP are based on a sign i.i.d. model for \mathbf{A} , which leads to decorrelation and Gaussianity after multiplication of a vector with \mathbf{A} or \mathbf{A}^T , similar to spreading and despreading in CDMA. Another somewhat popular model for \mathbf{A} is the Right Rotationally Invariant class, in which (only) the right singular vectors of \mathbf{A} are modeled as random, and in particular as Haar distributed. This is the motivation for Vector AMP (VAMP) [13]. To keep complexity low however, VAMP has to restrict diagonal covariances to multiples of identity, which e.g. is not useful for Sparse Bayesian Learning [14]. GAMP-style low complexity algorithms can be derived also, but they require some correction terms in the variance predictions, stemming from the Haar distribution [15], [16].

VIII. ACKNOWLEDGEMENTS

EURECOM's research is partially supported by its industrial members: ORANGE, BMW, SAP, iABG, Norton LifeLock, by the Franco-German projects CellFree6G and 5G-OPERA, by the EU H2030 project CONVERGE, and by a Huawei France funded Chair towards Future Wireless Networks.

REFERENCES

- [1] D. Zhang, X. Song, W. Wang, G. Fettweis, and X. Gao, "Unifying Message Passing Algorithms Under the Framework of Constrained Bethe Free Energy Minimization," *IEEE Transactions on Wireless Communications*, vol. 20, no. 7, 2021.
- [2] M. J. Wainwright, M. I. Jordan *et al.*, "Graphical Models, Exponential Families, and Variational Inference," *Foundations and Trends® in Machine Learning*, vol. 1, no. 1–2, 2008.
- [3] T. Minka *et al.*, "Divergence Measures and Message Passing," Citeseer, Tech. Rep., 2005.
- [4] T. Heskes, M. Opper, W. Wiegand, O. Winther, and O. Zoeter, "Approximate Inference Techniques with Expectation Constraints," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2005, no. 11, 2005.

- [5] Q. Zou and H. Yang, "A Concise Tutorial on Approximate Message Passing," *arXiv preprint arXiv:2201.07487*, 2022.
- [6] S. Rangan, P. Schniter, E. Riegler, A. K. Fletcher, and V. Cevher, "Fixed Points of Generalized Approximate Message Passing with Arbitrary Matrices," *IEEE Transactions on Information Theory*, vol. 62, no. 12, 2016.
- [7] C. Kurisummootil Thomas, Z. Zhao, and D. Slock, "Towards Convergent Approximate Message Passing by Alternating Constrained Minimization of Bethe Free Energy," in *IEEE Information Theory Workshop (ITW)*, Saint Malo, France, 2023.
- [8] S. Rangan, A. Fletcher, P. Schniter, and U. Kamilov, "Inference for Generalized Linear Models via Alternating Directions and Bethe Free Energy Minimization," *IEEE Trans. Info. Theory*, Jan. 2017.
- [9] Z. Zhao and D. Slock, "Bethe Free Energy and Extrinsic in Approximate Message Passing," in *IEEE Asilomar Conf. Signals, Systems and Computers*, 2023.
- [10] Z. Zhao, F. Xiao, and D. Slock, "Vector approximate message passing for not so large N.I.I.D. generalized I/O linear models," in *IEEE Int'l Conf. Acoustics, Speech and Sig. Proc. (ICASSP)*, Seoul, 2024.
- [11] —, "Extrinsic and Linearized Component-Wise Conditionally Unbiased MMSE Estimation as in GAMP," in *IEEE Asilomar Conf. Signals, Systems and Computers*, 2024.
- [12] —, "Approximate Message Passing for Not So Large iid Generalized Linear Models," in *Proc. Int'l Workshop Signal Processing Advances in Wireless Comm's (SPAWC)*, Sept. 2023.
- [13] S. Rangan, P. Schniter, and A. K. Fletcher, "Vector Approximate Message Passing," *IEEE Trans. On Info. Theo.*, Oct. 2019.
- [14] C. K. Thomas and D. Slock, "SAVE - Space alternating variational estimation for sparse Bayesian learning," in *IEEE Data Science Workshop*, June 2018.
- [15] Z. Zhao and D. Slock, "Variance Predictions in VAMP/UAMP with Right Rotationally Invariant Measurement Matrices for iid Generalized Linear Models," in *European Sig. Proc. Conf. (EUSIPCO)*, Helsinki, Finland, 2023.
- [16] —, "Improved Variance Predictions in Approximate Message Passing," in *IEEE Int'l Workshop Machine Learning and Sig. Proc. (MLSP)*, Rome, Italy, 2023.