

VECTOR APPROXIMATE MESSAGE PASSING FOR NOT SO LARGE N.I.I.D. GENERALIZED I/O LINEAR MODELS

Zilu Zhao, Fangqing Xiao, Dirk Slock
 Communication Systems Department, EURECOM, France
 zilu.zhao@eurecom.fr, fangqing.xiao@eurecom.fr, dirk.slock@eurecom.fr

Abstract—Many signal processing problems involve a Generalized Linear Model (GLM), which is a type of linear model where the unknowns may be non-identically independently distributed (n.i.i.d.). Vector Approximate Message Passing for Generalized Linear Models (GVAMP) is a computationally efficient belief propagation technique used for Bayesian inference. However, the posterior variances obtained from GVAMP with limited complexity are only exact under the assumption of an independent and identically distributed (i.i.d.) prior, owing to the averaging operations involved. In numerous problems, it is beneficial not just to estimate the unknowns but also to obtain accurate posterior distributions. While VAMP, and especially AMP, are applicable to high-dimensional problems, many applications involve dimensions that are not excessively high, allowing for more complex operations. Furthermore, in finite dimensions, the asymptotic regime that leads to correct variances under certain measurement matrix model assumptions is not applicable. To overcome these challenges, we propose a revised version of GVAMP, named reGVAMP. This method provides a multivariate Gaussian posterior approximation, which includes inter-parameter correlations, and yields accurate posterior marginals requiring only the extrinsic distributions to become Gaussian.

I. INTRODUCTION

Recovering the input signal in the context of a *generalized linear model* (GLM) [1] is a fundamental challenge in the field of signal processing. Applications involving this problem include statistics regression [2], wireless communication [3], and machine learning [4]. For instance, within the realm of wireless communication, $p_{y|x}$ often characterizes receiver hardware and inference mechanisms, in the symbol detection where \mathbf{x} represents a vector of discrete symbols to be recovered, or in the channel estimation and localization, where \mathbf{x} comprises propagation channel parameters, \mathbf{A} could embody the propagation features, modulation/demodulation schemes, or pilot symbols, depending on the scenario. Viewing this from a Bayesian perspective, the estimation of the random vector $\mathbf{x} \in \mathbb{R}^N$ from observed measurements $\mathbf{y} \in \mathbb{R}^M$ is often tackled through methods such as maximum a posteriori (MAP) estimation and minimum mean square error (MMSE) estimation. However, in certain scenarios, these methods become unviable due to the optimization complexities inherent in MAP or the intractable integrals associated with MMSE. To tackle this challenge, one established method is approximate inference. Generalized Approximate Message Passing (GAMP) [5] stands out as a prominent and computationally efficient approach within the realm of GLM. Expanding upon the foundation of Approximate Message Passing (AMP) [6] used in standard linear models (SLM), GAMP showcases its prowess in recovering high-dimensional signals. The dynamics of GAMP align with those of AMP, both governed by a state

evolution. Nevertheless, GAMP encounters convergence issues when faced with an ill-conditioned measurement matrix.

A. Prior Work

In response, a solution comes as an extension to the Vector Approximate Message Passing (VAMP) [7] algorithm, tailored for GLM scenarios. This approach, named GVAMP [8], accommodates the challenge of an ill-conditioned matrix. It applies a vector-level Expectation Propagation [9] algorithm on a specially designed factorization scheme.

Compared to AMP which requires i.i.d. measurement matrix to converge, VAMP shines particularly for its converging behavior in the scenario of a right rotationally invariant measurement matrix, with a scalar state evolution having been rigorously established. It is worth noting that GLM-VAMP and its variants often assume a high system dimension, which drives efforts to avoid complex matrix inversions and necessitates additional approximations for efficiency.

However, the original (G)VAMP algorithm solely provides averaged variances. This limitation serves as the impetus for introducing the reVAMP (revisited- Vector Approximate Message Passing) [10] for a standard linear model with arbitrary additive Gaussian noise. It gives a Gaussian approximated posterior distribution including inter-parameter correlations.

B. Main Contribution

In this paper, reGVAMP (revisited- Generalized Vector Approximate Message Passing) algorithm is extended from reVAMP to cover GLM. We derive reGVAMP by utilizing the multivariate Gaussian marginalization under EP framework. In reGVAMP, each likelihood and prior is approximated by a Gaussian distribution. Like reVAMP, it results in an approximated posterior within the Gaussian family and gives a covariance matrix with inter-parameter correlations.

II. DERIVATION UNDER EP FRAMEWORK

We consider here a generalized linear model

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

where the ratio N/M is a constant. We interpret the linear mixing as a conditional probability

$$p_{\mathbf{z}|\mathbf{x}}(\mathbf{z}, \mathbf{x}) = \delta(\mathbf{z} - \mathbf{A}\mathbf{x}). \quad (2)$$

We obtain the posterior probability by

$$\begin{aligned} p_{\mathbf{x}, \mathbf{z} | \mathbf{y}}(\mathbf{x}, \mathbf{z}) &\propto p_{\mathbf{x}, \mathbf{y}, \mathbf{z}}(\mathbf{x}, \mathbf{y}, \mathbf{z}) \\ &= \prod_j p_{y_j | z_j}(z_j) p_{\mathbf{z} | \mathbf{x}}(\mathbf{x}, \mathbf{z}) \prod_i p_{x_i}(x_i). \end{aligned} \quad (3)$$

We want to approximate this true posterior by

$$\begin{aligned} q_{\mathbf{x}, \mathbf{z} | \mathbf{y}}(\mathbf{x}, \mathbf{z}) &\propto \prod_j q_{z_j | \mathbf{y}}(z_j) \prod_i q_{x_i | \mathbf{y}}(x_i) \\ &\propto \prod_j q_{y_j | z_j}(z_j) q_{\mathbf{y}_{\bar{j}}, z_j}(z_j) \prod_i q_{x_i}(x_i) q_{\mathbf{y} | x_i}(x_i), \end{aligned} \quad (4)$$

where for all the events E , q_E is the Gaussian approximation for p_E . We shall also see later that $q_{\mathbf{y}_{\bar{j}}, z_j}$ and $q_{\mathbf{y} | x_i}$ are fully determined by $q_{y_j | z_j}$ and q_{x_i} . Ideally, we wish to find an approximation by minimizing

$$\begin{aligned} KLD[p_{\mathbf{x}, \mathbf{z} | \mathbf{y}} \| q_{\mathbf{x}, \mathbf{z} | \mathbf{y}}] &= \sum_j KLD[p_{z_j | \mathbf{y}} \| q_{z_j | \mathbf{y}}] + \sum_i KLD[p_{x_i | \mathbf{y}} \| q_{x_i | \mathbf{y}}] + c \\ &= \sum_j KLD[p_{z_j | \mathbf{y}} \| q_{z_j | \mathbf{y}}] + \sum_i KLD[p_{x_i | \mathbf{y}} \| q_{x_i | \mathbf{y}}] + c, \end{aligned} \quad (5)$$

where we define $KLD(p \| q) = \int p(x) \log \frac{p(x)}{q(x)} dx$ to be the Kullback-Leibler divergence. We can minimize alternately w.r.t. the factors $q_{z_j | \mathbf{y}}$ and $q_{x_i | \mathbf{y}}$. The true posteriors for z_j can be written as

$$\begin{aligned} p_{z_j | \mathbf{y}}(z_j) &\propto \int p_{\mathbf{x}, \mathbf{y}, \mathbf{z}}(\mathbf{x}, \mathbf{y}, \mathbf{z}) d\mathbf{x} d\mathbf{z}_{\bar{j}} \\ &= p_{y_j | z_j}(z_j) \int p_{\mathbf{y}_{\bar{j}} | \mathbf{z}_{\bar{j}}}(\mathbf{z}_{\bar{j}}) p_{\mathbf{z} | \mathbf{x}}(\mathbf{x}, \mathbf{z}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} d\mathbf{z}_{\bar{j}} \\ &= p_{y_j | z_j}(z_j) p_{\mathbf{y}_{\bar{j}} | \mathbf{z}_{\bar{j}}}(z_j) p_{z_j}(z_j) = p_{y_j | z_j}(z_j) p_{\mathbf{y}_{\bar{j}}, z_j}(z_j), \end{aligned} \quad (6)$$

where we define the notation $\forall \mathbf{z} \forall j, \mathbf{z}_{\bar{j}}$ denotes all the elements in \mathbf{z} except the j -th element.

A similar procedure can be done to derive the true posterior for x_i ,

$$\begin{aligned} p_{x_i | \mathbf{y}}(x_i) &\propto p_{x_i}(x_i) \int p_{\mathbf{y} | \mathbf{z}}(\mathbf{z}) p_{\mathbf{z} | \mathbf{x}}(\mathbf{x}, \mathbf{z}) p_{\mathbf{x}_{\bar{i}}}(x_{\bar{i}}) d\mathbf{x}_{\bar{i}} d\mathbf{z} \\ &= p_{x_i}(x_i) p_{\mathbf{y} | x_i}(x_i). \end{aligned} \quad (7)$$

To make the computation tractable, we approximate the extrinsic $p_{\mathbf{y}_{\bar{j}}, z_j}$ and $p_{\mathbf{y} | x_i}(x_i)$ as Gaussian. If Gaussian distributions are close approximations for $p_{\mathbf{y}_{\bar{j}}, z_j}$ and $p_{\mathbf{y} | x_i}$, then due to the Central Limit Theory, we have at convergence

$$\begin{aligned} KLD[p_{z_j | \mathbf{y}} \| q_{z_j | \mathbf{y}}] &\simeq KLD[p_{y_j | z_j} \cdot q_{\mathbf{y}_{\bar{j}}, z_j} \| q_{y_j | z_j} \cdot q_{\mathbf{y}_{\bar{j}}, z_j}] \\ &= \mathbb{E} \left[\ln \frac{p_{y_j | z_j}(z_j)}{q_{y_j | z_j}(z_j)} \right] \\ &\sum_i KLD[p_{x_i | \mathbf{y}} \| q_{x_i | \mathbf{y}}] \simeq KLD[p_{x_i} \cdot q_{\mathbf{y} | x_i} \| q_{x_i} \cdot q_{\mathbf{y} | x_i}] \\ &= \mathbb{E} \left[\ln \frac{p_{x_i}(x_i)}{q_{x_i}(x_i)} \right], \end{aligned} \quad (8)$$

where

$$\begin{aligned} q_{\mathbf{y}_{\bar{j}}, z_j}(z_j) &\propto \frac{\iint q_{\mathbf{y} | \mathbf{z}}(\mathbf{z}) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) q_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} d\mathbf{z}_{\bar{j}}}{q_{y_j | z_j}(z_j)}; \\ q_{\mathbf{y} | x_i}(x_i) &\propto \frac{\iint q_{\mathbf{x}}(\mathbf{x}) q_{\mathbf{y} | \mathbf{z}}(\mathbf{z}) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) d\mathbf{z} d\mathbf{x}_{\bar{i}}}{q_{x_i}(x_i)}, \end{aligned} \quad (9)$$

and the two expectations are calculated with respect to the tilted posteriors

$$q_{z_j | \mathbf{y}}(z_j) \propto p_{y_j | z_j}(z_j) q_{\mathbf{y}_{\bar{j}}, z_j}(z_j) \quad (10)$$

$$q_{x_i | \mathbf{y}}(x_i) \propto p_{x_i}(x_i) q_{\mathbf{y} | x_i}(x_i) \quad (11)$$

respectively. It is worth noticing that these tilted posteriors can be identified as marginal posterior with Gaussian approximated extrinsic.

A. Output Node

At the output node, we optimize the first expression in (8). We denote

$$\begin{aligned} q_{\mathbf{y} | \mathbf{z}}(\mathbf{z}) &= \prod_j q_{y_j | z_j}(z_j) := \mathcal{N}(\mathbf{z} | \mathbf{m}_{\mathbf{z}}, \mathbf{D}_{\tau_{\mathbf{z}}}) \\ q_{\mathbf{x}}(\mathbf{x}) &= \prod_i q_{x_i}(x_i) := \mathcal{N}(\mathbf{x} | \mathbf{m}_{\mathbf{x}}, \mathbf{D}_{\tau_{\mathbf{x}}}), \end{aligned} \quad (12)$$

where

$$\mathbf{D}_{\tau_{\mathbf{z}}} = \text{diag}[\tau_{z_j}]; \quad \mathbf{D}_{\tau_{\mathbf{x}}} = \text{diag}[\tau_{x_i}]. \quad (13)$$

The diag operation either transforms a vector into a diagonal matrix or extracts the diagonal elements from a matrix into a vector identical to the MATLAB function.

The optimization of the first KLD in (8) is equivalent to matching the first and second order moment with respect to z_j of the tilted posterior given by (10). We first compute the marginalization step of the tilted posterior

$$\begin{aligned} q_{z_j | \mathbf{y}}(z_j) &\propto \int \int \frac{p_{y_j | z_j}(z_j)}{q_{y_j | z_j}(z_j)} q_{\mathbf{y} | \mathbf{z}}(\mathbf{z}) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) q_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} d\mathbf{z}_{\bar{j}} \\ &= \frac{p_{y_j | z_j}(z_j)}{q_{y_j | z_j}(z_j)} \int q_{\mathbf{y} | \mathbf{z}}(\mathbf{z}) \int \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) \mathcal{N}(\mathbf{x} | \mathbf{m}_{\mathbf{x}}, \mathbf{D}_{\tau_{\mathbf{x}}}) d\mathbf{x} d\mathbf{z}_{\bar{j}} \\ &= \frac{p_{y_j | z_j}(z_j)}{q_{y_j | z_j}(z_j)} \int \mathcal{N}(\mathbf{z} | \mathbf{m}_{\mathbf{z}}, \mathbf{D}_{\tau_{\mathbf{z}}}) \mathcal{N}(\mathbf{z} | \mathbf{A}\mathbf{m}_{\mathbf{x}}, \mathbf{A}\mathbf{D}_{\tau_{\mathbf{x}}}\mathbf{A}^T) d\mathbf{z}_{\bar{j}} \end{aligned} \quad (14)$$

From the second line to the third line, we consider an auxiliary joint probability

$$\tilde{p}_{\mathbf{x}, \mathbf{z}}(\mathbf{x}, \mathbf{z}) = p_{\mathbf{z} | \mathbf{x}}(\mathbf{x}, \mathbf{z}) \tilde{p}_{\mathbf{x}}(\mathbf{x}) := \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) \mathcal{N}(\mathbf{x} | \mathbf{m}_{\mathbf{x}}, \mathbf{D}_{\tau_{\mathbf{x}}}). \quad (15)$$

Since $\mathbf{z} = \mathbf{A}\mathbf{x}$ and $\mathbf{x} \sim \mathcal{N}(\mathbf{m}_{\mathbf{x}}, \mathbf{D}_{\tau_{\mathbf{x}}})$, we have $\mathbf{z} \sim \mathcal{N}(\mathbf{A}\mathbf{m}_{\mathbf{x}}, \mathbf{A}\mathbf{D}_{\tau_{\mathbf{x}}}\mathbf{A}^T)$. The integral in the last line can be viewed as a multivariate Gaussian marginalization of Gaussian joint posterior (with \mathbf{x} marginalized out)

$$\mathcal{N}(\mathbf{z} | \mathbf{m}_{\mathbf{z}}, \mathbf{D}_{\tau_{\mathbf{z}}}) \mathcal{N}(\mathbf{z} | \mathbf{A}\mathbf{m}_{\mathbf{x}}, \mathbf{A}\mathbf{D}_{\tau_{\mathbf{x}}}\mathbf{A}^T) \propto \mathcal{N}(\mathbf{z} | \mathbf{m}_{\tilde{\mathbf{z}}}', \mathbf{C}_{\tilde{\mathbf{z}}}'), \quad (16)$$

where

$$\begin{aligned} \mathbf{C}_{\tilde{\mathbf{z}}}' &= (\mathbf{D}_{\tau_{\mathbf{z}}}^{-1} + (\mathbf{A}\mathbf{D}_{\tau_{\mathbf{x}}}\mathbf{A}^T)^{-1})^{-1} = \mathbf{A}(\mathbf{D}_{\tau_{\mathbf{z}}}^{-1} + \mathbf{A}^T\mathbf{D}_{\tau_{\mathbf{x}}}^{-1}\mathbf{A})^{-1}\mathbf{A}^T; \\ \mathbf{m}_{\tilde{\mathbf{z}}}' &= \mathbf{C}_{\tilde{\mathbf{z}}}'(\mathbf{D}_{\tau_{\mathbf{z}}}^{-1}\mathbf{m}_{\mathbf{z}} + (\mathbf{A}\mathbf{D}_{\tau_{\mathbf{x}}}\mathbf{A}^T)^{-1}\mathbf{A}\mathbf{m}_{\mathbf{x}}) \\ &= \mathbf{A}(\mathbf{D}_{\tau_{\mathbf{z}}}^{-1} + \mathbf{A}^T\mathbf{D}_{\tau_{\mathbf{x}}}^{-1}\mathbf{A})^{-1}(\mathbf{D}_{\tau_{\mathbf{z}}}^{-1}\mathbf{m}_{\mathbf{z}} + \mathbf{A}^T\mathbf{D}_{\tau_{\mathbf{x}}}^{-1}\mathbf{m}_{\mathbf{x}}). \end{aligned} \quad (17)$$

Denote

$$\tau_{\tilde{\mathbf{z}}}' = \text{diag}[\mathbf{C}_{\tilde{\mathbf{z}}}'], \quad (18)$$

we then obtain the extrinsic/cavity distribution $q_{\mathbf{y}_{\bar{j}}, z_j}(z_j)$ by multivariate Gaussian marginalization

$$\mathcal{N}(z_j | m_{p_j}, \tau_{p_j}) \propto \frac{\int \mathcal{N}(\mathbf{z} | \mathbf{m}_{\tilde{\mathbf{z}}}', \mathbf{C}_{\tilde{\mathbf{z}}}') d\mathbf{z}_{\bar{j}}}{q_{y_j | z_j}(z_j)} \propto \frac{\mathcal{N}(z_j | m_{z_j}', \tau_{z_j}')}{\mathcal{N}(z_j | m_{z_j}, \tau_{z_j})}, \quad (19)$$

where

$$\begin{aligned} \frac{1}{\tau_{p_j}} &= \frac{1}{\tau_{z_j}'} - \frac{1}{\tau_{z_j}}; \\ m_{p_j} &= \tau_{p_j} \left(\frac{m_{z_j}'}{\tau_{z_j}'} - \frac{m_{z_j}}{\tau_{z_j}} \right). \end{aligned} \quad (20)$$

Thus, the tilted posterior is proportional to

$$q_{z_j|\mathbf{y}}(z_j) \propto p_{y_j|z_j}(z_j)\mathcal{N}(z_j|m_{p_j}, \tau_{p_j}). \quad (21)$$

Matching the first and second order moment of the two parameters in the first KLD in (8), we obtain

$$q_{y_j|z_j}(z_j)q_{\mathbf{y}_{\bar{j}}, z_j}(z_j) \propto \mathcal{N}(z_j|m_{\hat{z}_j}, \tau_{\hat{z}_j}), \quad (22)$$

where \hat{z}_j and $\tau_{\hat{z}_j}$ are calculated as the mean and variances of the tilted posterior

$$\begin{aligned} m_{\hat{z}_j} &= \mathbb{E}_{q_{z_j|\mathbf{y}}}[z_j] = g_z(m_{p_j}, \tau_{p_j}); \\ \tau_{\hat{z}_j} &= \mathbb{E}_{q_{z_j|\mathbf{y}}}[(z_j - m_{\hat{z}_j})^2] = \tau_{p_j} g'_z(m_{p_j}, \tau_{p_j}), \end{aligned} \quad (23)$$

where we denote

$$g_z(m_{p_j}, \tau_{p_j}) = \frac{\int z_j p_{y_j|z_j}(z_j)\mathcal{N}(z_j|m_{p_j}, \tau_{p_j})dz_j}{\int p_{y_j|z_j}(z_j)\mathcal{N}(z_j|m_{p_j}, \tau_{p_j})dz_j}, \quad (24)$$

and the derivative is with respect to the first parameter m_{p_j} . The new approximated likelihood is then derived as

$$q_{y_j|z_j}(z_j) \propto \frac{\mathcal{N}(z_j|m_{\hat{z}_j}, \tau_{\hat{z}_j})}{q_{\mathbf{y}_{\bar{j}}, z_j}(z_j)} = \frac{\mathcal{N}(z_j|m_{\hat{z}_j}, \tau_{\hat{z}_j})}{\mathcal{N}(z_j|m_{p_j}, \tau_{p_j})}. \quad (25)$$

It is a quotient of two Gaussian, and thus, the approximated likelihood is updated by

$$q_{y_j|z_j}(z_j) = \mathcal{N}(z_j|m_{z_j}^+, \tau_{z_j}^+), \quad (26)$$

where

$$\begin{aligned} \frac{1}{\tau_{z_j}^+} &= \frac{1}{\tau_{\hat{z}_j}} - \frac{1}{\tau_{p_j}}; \\ m_{z_j}^+ &= \tau_{z_j}^+ \left(\frac{m_{\hat{z}_j}}{\tau_{\hat{z}_j}} - \frac{m_{p_j}}{\tau_{p_j}} \right). \end{aligned} \quad (27)$$

If sequential update is used, the j -th entry of \mathbf{m}_z and τ_z is replaced with $m_{z_j}^+$ and $\tau_{z_j}^+$ before the proceeding to the update of next variable node z_{j+1} . The computation of (17) can be simplified as a rank one update. To simplify the computation of the covariance, we define

$$\frac{1}{\Delta\tau_{z_j}} = \frac{1}{\tau_{z_j}^+} - \frac{1}{\tau_{z_j}}. \quad (28)$$

The updated covariance according to matrix inverse lemma is

$$\mathbf{C}_{\hat{z}'}^+ = \left(\mathbf{C}_{\hat{z}'}^{-1} + e_j \Delta\tau_{z_j}^{-1} e_j^T \right)^{-1} = \mathbf{C}_{\hat{z}'} - \frac{\mathbf{C}_{\hat{z}'} e_j e_j^T \mathbf{C}_{\hat{z}'}}{\Delta\tau_{z_j} + e_j^T \mathbf{C}_{\hat{z}'} e_j}. \quad (29)$$

Similarly, we define

$$\Delta m_{z_j} = \Delta\tau_{z_j} \left(\frac{m_{z_j}^+}{\tau_{z_j}^+} - \frac{m_{z_j}}{\tau_{z_j}} \right), \quad (30)$$

and the mean update is obtained as

$$\mathbf{m}_{\hat{z}'}^+ = \mathbf{m}_{\hat{z}'} + \frac{\Delta m_{z_j} - e_j^T \mathbf{m}_{\hat{z}'}}{\Delta\tau_{z_j} + e_j^T \mathbf{C}_{\hat{z}'} e_j} \mathbf{C}_{\hat{z}'} e_j. \quad (31)$$

Thus, the complexity for updating each z_j approximated prior is $O(N^2)$ due to (29). The complexity for approximation at the output node is dominated by $O(N^3)$ happening at the start of a sweep where we need to compute $\mathbf{C}_{\hat{z}'}$ according to (17) for the first time. Moreover, there are $O(N)$ variable nodes to be updated. In other words, the sweep over vector \mathbf{z} has complexity $O(N^3) + O(N)O(N^2) = O(N^3)$.

B. Input Node

At the input node, we optimize the second expression in (8) with respect to $q_{x_i}(x_i)$. This KLD's optimization is equivalent to matching the first and second order moment of the tilted posterior since q_{x_i} belongs to the Gaussian family.

With the notations defined in (12), the tilted posterior is given by

$$\begin{aligned} q_{x_i|\mathbf{y}}(x_i) &\propto \int \int \frac{p_{x_i}(x_i)}{q_{x_i}(x_i)} q_{\mathbf{x}}(\mathbf{x}) q_{\mathbf{y}|\mathbf{z}}(\mathbf{z}) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) d\mathbf{z} d\mathbf{x}_{\bar{i}} \\ &= \frac{p_{x_i}(x_i)}{q_{x_i}(x_i)} \int q_{\mathbf{x}}(\mathbf{x}) \int \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) \mathcal{N}(\mathbf{z}|\mathbf{m}_z, \mathbf{D}_{\tau_z}) d\mathbf{z} d\mathbf{x}_{\bar{i}} \\ &= \frac{p_{x_i}(x_i)}{q_{x_i}(x_i)} \int \mathcal{N}(\mathbf{x}|\mathbf{m}_x, \mathbf{D}_{\tau_x}) \mathcal{N}(\mathbf{A}\mathbf{x}|\mathbf{m}_z, \mathbf{D}_{\tau_z}) d\mathbf{x}_{\bar{i}}. \end{aligned} \quad (32)$$

Again, the last integral is proportional to a multi-variate Gaussian marginalization of the Gaussian joint posterior (with \mathbf{z} marginalized out), and hence, we compute the Gaussian joint posterior as

$$\mathcal{N}(\mathbf{x}|\mathbf{m}_x, \mathbf{D}_{\tau_x}) \mathcal{N}(\mathbf{A}\mathbf{x}|\mathbf{m}_z, \mathbf{D}_{\tau_z}) \propto \mathcal{N}(\mathbf{x}|\mathbf{m}_{\hat{x}'}, \mathbf{C}_{\hat{x}'}), \quad (33)$$

where

$$\begin{aligned} \mathbf{C}_{\hat{x}'} &= (\mathbf{D}_{\tau_x}^{-1} + \mathbf{A}^T \mathbf{D}_{\tau_z}^{-1} \mathbf{A})^{-1}; \\ \mathbf{m}_{\hat{x}'} &= \mathbf{C}_{\hat{x}'} (\mathbf{D}_{\tau_x}^{-1} \mathbf{m}_x + \mathbf{A}^T \mathbf{D}_{\tau_z}^{-1} \mathbf{m}_z). \end{aligned} \quad (34)$$

To obtain the marginal distribution, we denote

$$\tau_{\hat{x}'} = \text{diag}[\mathbf{C}_{\hat{x}'}]. \quad (35)$$

The extrinsic/cavity distribution $q_{\mathbf{y}|x_i}(x_i)$ can be computed as a quotient of two Gaussian distributions

$$\mathcal{N}(x_i|m_{r_i}, \tau_{r_i}) \propto \frac{\int \mathcal{N}(\mathbf{x}|\mathbf{m}_{\hat{x}'}, \mathbf{C}_{\hat{x}'}) d\mathbf{x}_{\bar{i}}}{q_{x_i}(x_i)} \propto \frac{\mathcal{N}(x_i|m_{\hat{x}'}, \tau_{\hat{x}'})}{\mathcal{N}(x_i|m_{x_i}, \tau_{x_i})} \quad (36)$$

where

$$\begin{aligned} \frac{1}{\tau_{r_i}} &= \frac{1}{\tau_{\hat{x}'}} - \frac{1}{\tau_{x_i}}; \\ m_{r_i} &= \tau_{r_i} \left(\frac{m_{\hat{x}'}}{\tau_{\hat{x}'}} - \frac{m_{x_i}}{\tau_{x_i}} \right). \end{aligned} \quad (37)$$

The tilted posterior is proportional to the product of Gaussian extrinsic and the prior of x_i ,

$$q_{x_i|\mathbf{y}}(x_i) \propto p_{x_i}(x_i) \mathcal{N}(x_i|m_{r_i}, \tau_{r_i}). \quad (38)$$

Analog to the discussion (22) to (24), the approximated posterior is obtained by

$$q_{x_i}(x_i) q_{\mathbf{y}|x_i}(x_i) \propto \mathcal{N}(x_i|m_{\hat{x}_i}, \tau_{\hat{x}_i}), \quad (39)$$

where

$$\begin{aligned} m_{\hat{x}_i} &= \mathbb{E}_{q_{x_i|\mathbf{y}}}[x_i] = g_x(m_{r_i}, \tau_{r_i}); \\ \tau_{\hat{x}_i} &= \mathbb{E}_{q_{x_i|\mathbf{y}}}[(x_i - m_{\hat{x}_i})^2] = \tau_{r_i} g'_x(m_{r_i}, \tau_{r_i}), \end{aligned} \quad (40)$$

with $g_x(m_{r_i}, \tau_{r_i})$ defined by

$$g_x(m_{r_i}, \tau_{r_i}) = \frac{\int x_i p_{x_i}(x_i) \mathcal{N}(x_i|m_{r_i}, \tau_{r_i}) dx_i}{\int p_{x_i}(x_i) \mathcal{N}(x_i|m_{r_i}, \tau_{r_i}) dx_i}, \quad (41)$$

and the derivative $g'_x(m_{r_i}, \tau_{r_i})$ is once again with respect to the first parameter m_{r_i} . Hence, the approximated prior can be updated by

$$q_{x_i}(x_i) \propto \frac{\mathcal{N}(x_i|m_{\hat{x}_i}, \tau_{x_i})}{q_{\mathbf{y}|x_i}(x_i)} = \frac{\mathcal{N}(x_i|m_{\hat{x}_i}, \tau_{\hat{x}_i})}{\mathcal{N}(x_i|m_{r_i}, \tau_{r_i})}. \quad (42)$$

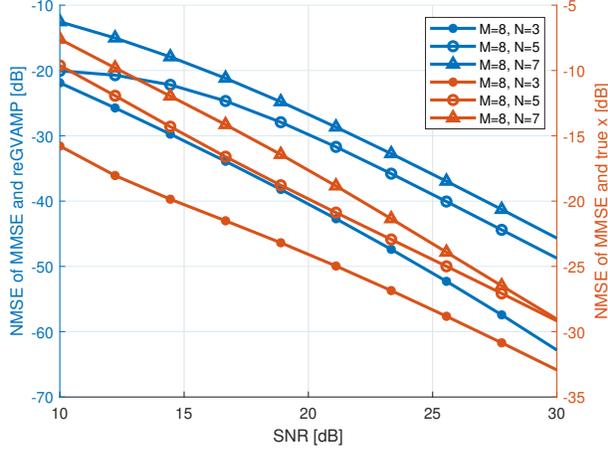


Fig. 1. Performance of GLM-reVAMP algorithm with different variable length N .

This Gaussian quotient can be calculated by

$$q_{x_i}(x_i) = \mathcal{N}(x_i | m_{x_i}^+, \tau_{x_i}^+) \propto \frac{\mathcal{N}(x_i | m_{\hat{x}_i}, \tau_{\hat{x}_i})}{\mathcal{N}(x_i | m_{r_i}, \tau_{r_i})}, \quad (43)$$

where

$$\begin{aligned} \frac{1}{\tau_{x_i}^+} &= \frac{1}{\tau_{\hat{x}_i}} - \frac{1}{\tau_{r_i}}; \\ m_{x_i}^+ &= \tau_{x_i}^+ \left(\frac{m_{\hat{x}_i}}{\tau_{\hat{x}_i}} - \frac{m_{r_i}}{\tau_{r_i}} \right). \end{aligned} \quad (44)$$

We use a similar method as described in (28) till (31) to simplify the LMMSE step (34). Define

$$\frac{1}{\Delta_{\tau_{x_i}}} = \frac{1}{\tau_{x_i}^+} - \frac{1}{\tau_{x_i}}; \Delta_{m_{x_i}} = \Delta_{\tau_{x_i}} \left(\frac{m_{x_i}^+}{\tau_{x_i}^+} - \frac{m_{x_i}}{\tau_{x_i}} \right). \quad (45)$$

The update of (34) with the new estimated prior is

$$\begin{aligned} \mathbf{C}_{\hat{\mathbf{x}}'}^+ &= \mathbf{C}_{\hat{\mathbf{x}}'} - \frac{\mathbf{C}_{\hat{\mathbf{x}}'} \mathbf{e}_i \mathbf{e}_i^T \mathbf{C}_{\hat{\mathbf{x}}'}}{\Delta_{\tau_{x_i}} + \mathbf{e}_i^T \mathbf{C}_{\hat{\mathbf{x}}'} \mathbf{e}_i}; \\ \mathbf{m}_{\hat{\mathbf{x}}'}^+ &= \mathbf{m}_{\hat{\mathbf{x}}'} + \frac{\Delta_{m_{x_i}} - \mathbf{e}_i^T \mathbf{m}_{\hat{\mathbf{x}}'}}{\Delta_{\tau_{x_i}} + \mathbf{e}_i^T \mathbf{C}_{\hat{\mathbf{x}}'} \mathbf{e}_i} \mathbf{C}_{\hat{\mathbf{x}}'} \mathbf{e}_i. \end{aligned} \quad (46)$$

Therefore, at the start of a sweep, the algorithm has a complexity of $O(N^3)$ due to the matrix inverse operation. After that, during the sweep, the algorithm still has a complexity of $O(N^3)$ because there are $O(N)$ messages each of which has a complexity of $O(N^2)$. This algorithm displays a nice symmetry between the input and output nodes. Considering both sweeps for the input and output node, the algorithm has an overall complexity of $O(N^3)$. We conclude the discussion in Algorithm reGVAMP.

III. SIMULATION RESULTS

We present the results of numerical experiments with varying parameters. In this experiment, we set M to 8 and N to $\{3, 5, 7\}$, respectively. Both the prior distribution $p_{\mathbf{x}}(\mathbf{x})$ and the conditional distribution $p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})$ are assumed to follow mixture Gaussian distributions. Specifically, they can be expressed as:

$$p_{x_i}(x_i) = 0.5\mathcal{N}(x_i | 0, \sigma_{x_{1i}}^2) + 0.5\mathcal{N}(x_i | 0, \sigma_{x_{2i}}^2), \quad (47)$$

$$p_{y_j|z_j}(z_j) = 0.5\mathcal{N}(y_j | z_j, \sigma_{v_{1j}}^2) + 0.5\mathcal{N}(y_j | z_j, \sigma_{v_{2j}}^2), \quad (48)$$

Algorithm 1 reGVAMP under EP Framework

Require: $\mathbf{y}, \mathbf{z} = \mathbf{A}\mathbf{x}, p_{\mathbf{x}}(\mathbf{x}), p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z})$

- 1: Initialize: $\mathbf{m}_{\mathbf{z}}, \tau_{\mathbf{z}}, \mathbf{m}_{\mathbf{x}}, \tau_{\mathbf{x}}$
- 2: **repeat**[For iteration step l]
- 3: [Update the Output Nodes]
- 4: $\mathbf{C}_{\hat{\mathbf{z}}'} = \mathbf{A}(\mathbf{D}_{\tau_{\mathbf{x}}}^{-1} + \mathbf{A}^T \mathbf{D}_{\tau_{\mathbf{z}}}^{-1} \mathbf{A})^{-1} \mathbf{A}^T$
- 5: $\mathbf{m}_{\hat{\mathbf{z}}}' = \mathbf{C}_{\hat{\mathbf{z}}}'(\mathbf{D}_{\tau_{\mathbf{z}}}^{-1} \mathbf{m}_{\mathbf{z}} + (\mathbf{A} \mathbf{D}_{\tau_{\mathbf{x}}} \mathbf{A}^T)^{-1} \mathbf{A} \mathbf{m}_{\mathbf{x}})$
- 6: **repeat**[For $j = 1$]
- 7: $\tau_{\hat{z}_j}' = [\mathbf{C}_{\hat{\mathbf{z}}'}]_{jj}$
- 8: $\tau_{p_j} = \frac{\tau_{z_j} \tau_{\hat{z}_j}'}{\tau_{z_j} - \tau_{\hat{z}_j}'}; m_{p_j} = \frac{m_{z_j} \tau_{z_j} - m_{z_j} \tau_{\hat{z}_j}'}{\tau_{z_j} - \tau_{\hat{z}_j}'}$
- 9: $m_{\hat{z}_j} = g_z(m_{p_j}, \tau_{p_j}); \tau_{\hat{z}_j} = \tau_{p_j} g_z'(m_{p_j}, \tau_{p_j})$
- 10: $\tau_{z_j}^+ = \frac{\tau_{p_j} \tau_{\hat{z}_j}}{\tau_{p_j} - \tau_{\hat{z}_j}}; m_{z_j}^+ = \frac{m_{z_j} \tau_{p_j} - m_{p_j} \tau_{\hat{z}_j}}{\tau_{p_j} - \tau_{\hat{z}_j}}$
- 11: Update $\mathbf{C}_{\hat{\mathbf{z}}}'^+$ and $\mathbf{m}_{\hat{\mathbf{z}}}'^+$ according to (29), (31)
- 12: **until** $j = M$
- 13: [Update the Input Nodes]
- 14: $\mathbf{C}_{\hat{\mathbf{x}}}' = (\mathbf{D}_{\tau_{\mathbf{x}}}^{-1} + \mathbf{A}^T \mathbf{D}_{\tau_{\mathbf{z}}}^{-1} \mathbf{A})^{-1}$
- 15: $\mathbf{m}_{\hat{\mathbf{x}}}' = \mathbf{C}_{\hat{\mathbf{x}}}'(\mathbf{D}_{\tau_{\mathbf{x}}}^{-1} \mathbf{m}_{\mathbf{x}} + \mathbf{A}^T \mathbf{D}_{\tau_{\mathbf{z}}}^{-1} \mathbf{m}_{\mathbf{z}})$
- 16: **repeat**[For $i = 1$]
- 17: $\tau_{\hat{x}_i}' = [\mathbf{C}_{\hat{\mathbf{x}}'}]_{ii}$
- 18: $\tau_{r_i} = \frac{\tau_{x_i} \tau_{\hat{x}_i}'}{\tau_{x_i} - \tau_{\hat{x}_i}'}; m_{r_i} = \frac{m_{x_i} \tau_{x_i} - m_{x_i} \tau_{\hat{x}_i}'}{\tau_{x_i} - \tau_{\hat{x}_i}'}$
- 19: $m_{\hat{x}_i} = g_x(m_{r_i}, \tau_{r_i}); \tau_{\hat{x}_i} = \tau_{r_i} g_x'(m_{r_i}, \tau_{r_i})$
- 20: $\tau_{x_i}^+ = \frac{\tau_{\hat{x}_i} \tau_{r_i}}{\tau_{\hat{x}_i} - \tau_{r_i}}; m_{x_i}^+ = \frac{m_{x_i} \tau_{r_i} - m_{r_i} \tau_{\hat{x}_i}}{\tau_{\hat{x}_i} - \tau_{r_i}}$
- 21: Update $\mathbf{C}_{\hat{\mathbf{x}}}'^+$ and $\mathbf{m}_{\hat{\mathbf{x}}}'^+$ according to (46)
- 22: **until** $i = N$
- 23: **until** Convergence

where $\sigma_{v_{1j}}, \sigma_{v_{2j}}, \sigma_{x_{1i}}$ and $\sigma_{x_{2i}}$ are independently and uniformly drawn from the interval $(0, 1]$. The elements of the measurement matrix \mathbf{A} are drawn independently from a Gaussian distribution $\mathcal{N}(0, \gamma)$. By adjusting γ , we can control the signal-to-noise ratio (SNR) of the system. We conducted simulations with SNR values ranging from 10 to 30 dB. Additionally, the recovery performance was evaluated using the normalized mean-squared error (NMSE) defined as $\|\hat{\mathbf{x}} - \mathbf{x}\|^2 / \|\mathbf{x}\|^2$ of one realization. The NMSE of $\hat{\mathbf{x}}_{mmse}$ with respect to \mathbf{x} and $\hat{\mathbf{x}}_{reGVAMP}$ were chosen as performance metrics.

The simulation results are presented in Figure 1. It is evident from the figure that as SNR increases, the bias decreases. Furthermore, as N approaches the value of M , the estimation performance deteriorates, although the overall recovery remains satisfactory. These results confirm the effectiveness of the proposed algorithm.

IV. CONCLUDING REMARKS

In this paper, we proposed an iterative Bayesian estimation method for estimating the GLM input signal which enables us to obtain the posterior mean $\mathbf{m}_{\hat{\mathbf{x}}}'$ and covariance matrix $\mathbf{C}_{\hat{\mathbf{x}}}'$ at the cost of higher complexity. It also gives Gaussian approximation for the prior and likelihood as a byproduct. Further study is required to study its convergence behavior.

Acknowledgements EURECOM's research is partially supported by its industrial members: ORANGE, BMW, SAP, iABG, Norton LifeLock, and by the Franco-German projects CellFree6G and 5G-OPERA.

REFERENCES

- [1] J. A. Nelder and R. W. Wedderburn, "Generalized linear models," *J. R. Stat. Soc. Ser. A. Stat. Soc.*, vol. 135, no. 3, pp. 370–384, 1972.
- [2] J. D. Blume, L. Su, R. M. Olveda, and S. T. McGarvey, "Statistical evidence for glm regression parameters: a robust likelihood approach," *STAT. MED.*, vol. 26, no. 15, pp. 2919–2936, 2007.
- [3] J. Boutros and G. Caire, "Iterative multiuser joint decoding: Unified framework and asymptotic analysis," *IEEE Trans. Inf. Theory*, vol. 48, no. 7, pp. 1772–1793, 2002.
- [4] Y. Wang, R. Wang, S. S. Du, and A. Krishnamurthy, "Optimism in reinforcement learning with generalized linear function approximation," *arXiv:1912.04136*, 2019.
- [5] S. Rangan, "Generalized approximate message passing for estimation with random linear mixing," in *IEEE Int. Symp. Inf. Theory - Proc.*, 2011, pp. 2168–2172.
- [6] D. L. Donoho, A. Maleki, and A. Montanari, "Message-passing algorithms for compressed sensing," *PNAS*, vol. 106, no. 45, pp. 18 914–18 919, 2009.
- [7] S. Rangan, P. Schniter, and A. K. Fletcher, "Vector approximate message passing," *IEEE Trans. Inf. Theory*, vol. 65, no. 10, pp. 6664–6684, 2019.
- [8] P. Schniter, S. Rangan, and A. K. Fletcher, "Vector approximate message passing for the generalized linear model," in *IEEE ASILOMAR*, 2016, pp. 1525–1529.
- [9] T. P. Minka, "Expectation propagation for approximate bayesian inference," *arXiv:1301.2294*, 2013.
- [10] Z. Zhao, F. Xiao, and D. Slock, "Approximate Message Passing for not so Large NIID Generalized Linear Models," in *SPAWC 2023, 24th IEEE International Workshop on Signal Processing Advances in Wireless Communications*.