

Vector Approximate Message Passing for the Generalized Linear Model

Phil Schniter



THE OHIO STATE UNIVERSITY

Duke
UNIVERSITY



Collaborators: **Sundeeep Rangan** (NYU), **Alyson Fletcher** (UCLA)

Supported in part by NSF grant CCF-1527162.

Asilomar Conference — Nov 8, 2016

Signal Recovery

- We consider problems where we want to
 - recover a “structured” **signal** $\mathbf{x} \in \mathbb{C}^N$
 - from “corrupted” **measurements** $\mathbf{y} \in \mathbb{C}^M$
 - of hidden **linear-transform outputs** $\mathbf{z} = \mathbf{A}\mathbf{x} \in \mathbb{C}^M$.
- The measurement corruption mechanism might be
 - **additive**: $y_i = z_i + w_i$, but possibly non-Gaussian
 - **quantized**: $y_i = \text{sgn}(z_i + w_i)$, such as in classification & one-bit CS
 - **phase-less**: $y_i = |z_i + w_i|$, such as in phase retrieval
 - **Poisson**, such as in photon-limited imaging, etc...
- The signal \mathbf{x} might be
 - **(approximately) sparse**, such as in compressive sensing
 - **finite alphabet**, such as in communications
 - **constant modulus**, etc...

Generalized Linear Model (GLM)

We take a statistical approach to signal recovery:

- corruption modeled using a **likelihood fn** $p(\mathbf{y}|\mathbf{z})$ with $\mathbf{z} = \mathbf{A}\mathbf{x}$
- signal modeled using a **prior distribution** $p(\mathbf{x})$

The **posterior** tells all we can learn about \mathbf{x} , but it's not computable:

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x}) p(\mathbf{y}|\mathbf{A}\mathbf{x})}{p(\mathbf{y})}.$$

Instead, we usually settle for **point estimates** of \mathbf{x} like the

- **MAP** estimate: $\hat{\mathbf{x}}_{\text{MAP}} = \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y})$
 - **MMSE** estimate: $\hat{\mathbf{x}}_{\text{MMSE}} = \mathbb{E}\{\mathbf{x}|\mathbf{y}\} = \int_{\mathbb{C}^N} \mathbf{x} p(\mathbf{x}|\mathbf{y}) d\mathbf{x}$
- and perhaps **marginal uncertainty** information like $\text{var}\{x_j|\mathbf{y}\}$.

Assumptions

In this talk, we assume a

- separable prior: $p(\mathbf{x}) = \prod_{j=1}^N p(x_j)$
- separable likelihood: $p(\mathbf{y}|\mathbf{z}) = \prod_{i=1}^M p(y_i|z_i)$

Then MAP estimation reduces to a familiar optimization problem:

$$\begin{aligned}\hat{\mathbf{x}}_{\text{MAP}} &= \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \\ &= \arg \max_{\mathbf{x}} \ln p(\mathbf{x}|\mathbf{y}) \\ &= \arg \max_{\mathbf{x}} \underbrace{\sum_{i=1}^M \ln p(y_i | [\mathbf{A}\mathbf{x}]_i)}_{\text{data fidelity}} + \underbrace{\sum_{j=1}^N \ln p(x_j)}_{\text{regularization}}.\end{aligned}$$

E.g., AWGN & Laplace $\Rightarrow \hat{\mathbf{x}}_{\text{MAP}} = \arg \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1.$

But often the prior and/or likelihood are **not log-concave!**

Existing Methods

- 1 Convex optimization
 - MAP only 😞
 - need log-concave prior & likelihood 😞
- 2 Sparse Bayesian Learning (SBL) & Expectation Propagation (EP)
 - posterior must be log-concave 😞
 - additional constraints on prior & likelihood 😞
 - per-iteration matrix inverse (slow) 😞
- 3 MCMC
 - slow, convergence difficult to assess 😞
- 4 Generalized Approximate Message Passing (GAMP)
 - any prior & likelihood 😊
 - no matrix inverses (fast) 😊
 - guaranteed only under large, i.i.d. Gaussian \mathbf{A} 😞

Proposed Method

We propose to ...

- 1 Rewrite $\mathbf{z} = \mathbf{A}\mathbf{x}$ as $\mathbf{0} = [\mathbf{A} \quad -\mathbf{I}] \begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} \triangleq \overline{\mathbf{A}}\overline{\mathbf{x}}$, thereby converting the GLM problem to a **standard linear regression** problem:

Recover $\overline{\mathbf{x}}$ from $\overline{\mathbf{y}} = \overline{\mathbf{A}}\overline{\mathbf{x}} + \overline{\mathbf{w}}$ with $\overline{\mathbf{w}} \sim \mathcal{N}(\mathbf{0}, \epsilon\mathbf{I})$, where now $\overline{\mathbf{y}} = \mathbf{0}$ and $\epsilon \rightarrow 0$.

- 2 Apply the recently proposed “**Vector AMP**” algorithm,¹ tracking separate divergences on \mathbf{x} and \mathbf{z} .

¹Rangan, Schniter, Fletcher—arXiv:1610.03082

Vector AMP for Standard Linear Regression

To recover \mathbf{x} from $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ with $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \mathbf{I})$ and i.i.d. $x_j \sim p(x_j) \dots$

Initialize $\tilde{\mathbf{r}}_1 = \mathbf{0}$ and $\tilde{\sigma}_1 = \infty$.

Repeat for $t = 1, 2, 3, \dots$

$$\tilde{\mathbf{x}}_t = \tilde{\boldsymbol{\eta}}(\tilde{\mathbf{r}}_t; \tilde{\sigma}_t) \quad \text{LMMSE estimation} \quad (1a)$$

$$\tilde{\tau}_t = \langle \tilde{\boldsymbol{\eta}}'(\tilde{\mathbf{r}}_t; \tilde{\sigma}_t) \rangle \quad \text{divergence} \quad (1b)$$

$$\mathbf{r}_t = (\tilde{\mathbf{x}}_t - \tilde{\tau}_t \tilde{\mathbf{r}}_t) / (1 - \tilde{\tau}_t) \quad \text{Onsager correction} \quad (1c)$$

$$\sigma_t^2 = \tilde{\sigma}_t^2 \tilde{\tau}_t / (1 - \tilde{\tau}_t) \quad \text{variance update} \quad (1d)$$

$$\hat{\mathbf{x}}_t = \boldsymbol{\eta}(\mathbf{r}_t; \sigma_t) \quad \text{denoising} \quad (2a)$$

$$\tau_t = \langle \boldsymbol{\eta}'(\mathbf{r}_t; \sigma_t) \rangle \quad \text{divergence} \quad (2b)$$

$$\tilde{\mathbf{r}}_{t+1} = (\hat{\mathbf{x}}_t - \tau_t \mathbf{r}_t) / (1 - \tau_t) \quad \text{Onsager correction} \quad (2c)$$

$$\tilde{\sigma}_{t+1}^2 = \sigma_t^2 \tau_t / (1 - \tau_t) \quad \text{variance update} \quad (2d)$$

where

$$[\boldsymbol{\eta}(\mathbf{r}_t; \sigma_t)]_j = \begin{cases} \int x_j p(x_j | r_{tj}) dx_j & \text{MMSE} \\ \arg \max_{x_j} p(x_j | r_{tj}) & \text{MAP} \end{cases} \quad \text{with } p(x_j | r_{tj}) \propto p(x_j) \mathcal{N}(x_j; r_{tj}, \sigma_t^2)$$

$$\tilde{\boldsymbol{\eta}}(\tilde{\mathbf{r}}_t; \tilde{\sigma}_t) = \mathbf{V} \left(\text{Diag}(\mathbf{s})^2 + \frac{\sigma_w^2}{\tilde{\sigma}_t^2} \mathbf{I}_R \right)^{-1} \left(\text{Diag}(\mathbf{s}) \mathbf{U}^\top \mathbf{y} + \frac{\sigma_w^2}{\tilde{\sigma}_t^2} \mathbf{V}^\top \tilde{\mathbf{r}}_t \right)$$

with SVD $\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^\text{H}$

Why call this “Vector AMP”?

- 1) Can be derived using an **approximation of message passing** on a factor graph, now with **vector-valued** variable nodes.
- 2) Performance rigorously characterized by a scalar **state-evolution**² under certain large random \mathbf{A} :

$$SVD \mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

- \mathbf{U} is deterministic
- \mathbf{S} is deterministic
- \mathbf{V} is uniformly distributed on the group of orthogonal matrices

“ \mathbf{A} is **right rotationally invariant**.”

Thus the VAMP state evolution holds for “almost any \mathbf{A} .”

²Rangan, Fletcher, Schniter–16

Connections to the Replica Prediction

- The **replica method** from statistical physics is often used to characterize the average behavior of large disordered systems.
- Although not fully rigorous, replica predictions are usually correct.
- For estimation of i.i.d. \mathbf{x} from measurements $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathcal{N}(\mathbf{0}; \sigma_w^2 \mathbf{I})$ under large right-rotationally invariant \mathbf{A} :

The MMSE $\mathcal{E}(\sigma_t^2)$ should satisfy the fixed-point equation³

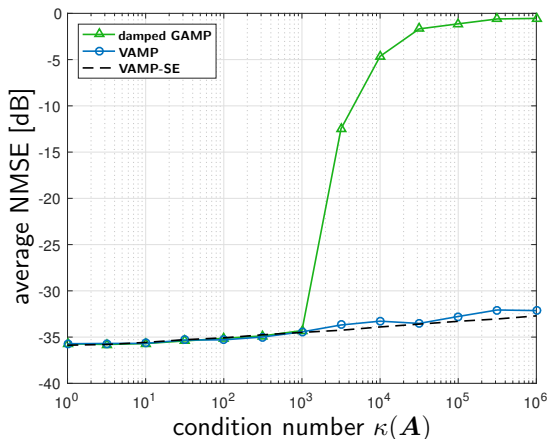
$$1/\sigma_t^2 = R_{\mathbf{A}^\top \mathbf{A} / \sigma_w^2}(-\mathcal{E}(\sigma_t^2)),$$

where $R_C(\cdot)$ denotes the R-transform of matrix C and $\mathcal{E}(\sigma_t^2) \triangleq \mathbb{E} \{ [\eta(x_j + \mathcal{N}(0, \sigma_t^2); \sigma_t^2) - x_j]^2 \}$.

- It can be shown that **VAMP's SE fixed-points obey the above equation**.
- Thus, assuming that the replica prediction is correct, VAMP will generate **MMSE estimates** whenever these fixed-points are unique.

³Tulino, Caire, Verdu, Shamai—TIT'13

Numerical Results: 1-Bit Compressive Sensing



$$N = 512$$

$$M/N = 4$$

$$\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^T$$

\mathbf{U}, \mathbf{V} drawn uniform

$$s_i/s_{i-1} = \rho \quad \forall i$$

ρ determines $\kappa(\mathbf{A})$

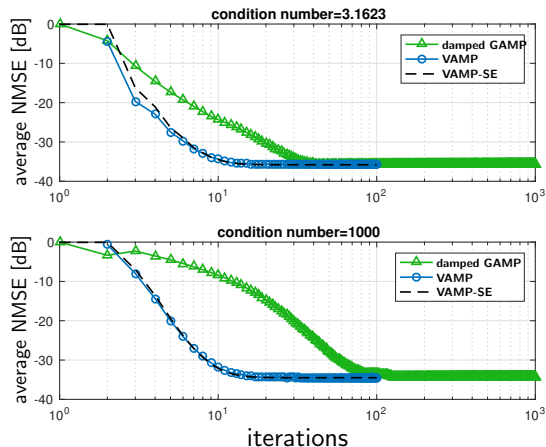
$x_j \sim \text{Bernoulli-Gaussian}$

$$\Pr\{x_j \neq 0\} = 1/32$$

$$\text{SNR} = 40\text{dB}$$

VAMP is robust to ill-conditioned \mathbf{A} ; GAMP is not.

Numerical Results: 1-Bit Compressive Sensing



$$N = 512$$

$$M/N = 4$$

$$\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^T$$

\mathbf{U}, \mathbf{V} drawn uniform

$$s_i/s_{i-1} = \rho \quad \forall i$$

ρ determines $\kappa(\mathbf{A})$

$x_j \sim \text{Bernoulli-Gaussian}$

$$\Pr\{x_j \neq 0\} = 1/32$$

$$\text{SNR} = 40\text{dB}$$

VAMP is much faster than damped GAMP.

Non-parametric Estimation

- So far we have considered estimating x from

$$\mathbf{y} \sim p(\mathbf{y}|\mathbf{z}; \boldsymbol{\theta}_z) \text{ where } \mathbf{z} = \mathbf{A}\mathbf{x} \text{ and } \mathbf{x} \sim p(\mathbf{x}; \boldsymbol{\theta}_x),$$

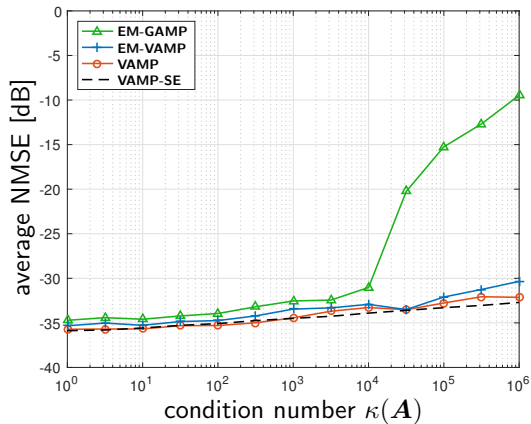
where $\boldsymbol{\theta}_z$ and $\boldsymbol{\theta}_x$ are parameters of the likelihood and prior.

- What if $\boldsymbol{\theta}_z$ and $\boldsymbol{\theta}_x$ are **unknown**? Can we learn them from \mathbf{y} ?
- Yes! The “**EM-VAMP**” approach⁴ can be directly applied.

⁴Fletcher, Schniter—arXiv:1602.08207

Numerical Results: Nonparametric 1-Bit CS

Learning both σ_w^2 and BG parameters:



$N = 512$

$M/N = 4$

$$\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^T$$

\mathbf{U}, \mathbf{V} drawn uniform

$$s_i/s_{i-1} = \rho \quad \forall i$$

ρ determines $\kappa(\mathbf{A})$

$x_j \sim \text{Bernoulli-Gaussian}$

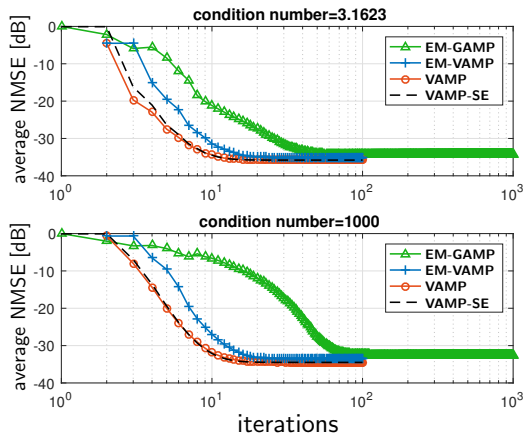
$$\Pr\{x_j \neq 0\} = 1/32$$

SNR = 40dB

EM-VAMP performs near oracle VAMP even with ill-conditioned \mathbf{A} .

Numerical Results: Nonparametric 1-Bit CS

Learning both σ_w^2 and BG parameters:



$N = 512$

$M/N = 4$

$\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^T$

\mathbf{U}, \mathbf{V} drawn uniform

$s_i/s_{i-1} = \rho \forall i$

ρ determines $\kappa(\mathbf{A})$

$x_j \sim \text{Bernoulli-Gaussian}$

$\Pr\{x_j \neq 0\} = 1/32$

SNR = 40dB

EM-VAMP slightly slower than VAMP but much faster than EM-GAMP.

Conclusions

- We proposed a new approach for inference under **generalized linear models** (GLMs).
- Applications include **1-bit compressive sensing**, **binary classification**, **(compressive) phase retrieval**, **photon-limited imaging**, etc.
- Our approach builds on the recently proposed “**vector AMP**” algorithm, which (unlike AMP) is **robust** to the choice of measurement operator \mathbf{A} .
- After an initial **SVD**, our approach consumes only **two matrix-vector multiplications per iteration** and converges in ~ 10 iterations.
- Our approach can be easily extended to the nonparametric case, where the likelihood and/or prior have **unknown parameters**, via EM-VAMP.
- In the future, we hope to rigorously prove the **state evolution** of VAMP-GLM and analyze the performance of EM-VAMP-GLM.