

Recent Advances in Approximate Message Passing

Phil Schniter



THE OHIO STATE UNIVERSITY

Supported in part by NSF grant CCF-1716388.

July 5, 2019

Overview

- 1 Linear Regression
- 2 Approximate Message Passing (AMP)
- 3 Vector AMP (VAMP)
- 4 Unfolding AMP and VAMP into Deep Neural Networks
- 5 Extensions: GLMs, Parameter Learning, Bilinear Problems

Outline

- 1 Linear Regression
- 2 Approximate Message Passing (AMP)
- 3 Vector AMP (VAMP)
- 4 Unfolding AMP and VAMP into Deep Neural Networks
- 5 Extensions: GLMs, Parameter Learning, Bilinear Problems

The Linear Regression Problem

Consider the following linear regression problem:

Recover \mathbf{x}_o from	with	$\begin{cases} \mathbf{x}_o \in \mathbb{R}^n & \text{unknown signal} \\ \mathbf{A} \in \mathbb{R}^{m \times n} & \text{known linear operator} \\ \mathbf{w} \in \mathbb{R}^m & \text{white Gaussian noise.} \end{cases}$
$\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathbf{w}$		

Typical methodologies:

- 1 Optimization (or MAP estimation):

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + R(\mathbf{x}) \right\}$$

- 2 Approximate MMSE:

$$\hat{\mathbf{x}} \approx \mathbb{E}\{\mathbf{x}|\mathbf{y}\} \quad \text{for } \mathbf{x} \sim p(\mathbf{x}), \quad \mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{A}\mathbf{x}, \nu_w \mathbf{I})$$

- 3 Plug-and-play:¹ iteratively apply a denoising algorithm like BM3D
- 4 Train a deep network to recover \mathbf{x}_o from \mathbf{y} .

¹Venkatkrishnan, Bouman, Wohlberg'13

Outline

- 1 Linear Regression
- 2 Approximate Message Passing (AMP)**
- 3 Vector AMP (VAMP)
- 4 Unfolding AMP and VAMP into Deep Neural Networks
- 5 Extensions: GLMs, Parameter Learning, Bilinear Problems

The AMP Methodology

- All of the aforementioned methodologies can be addressed using the **Approximate Message Passing (AMP)** framework.
- AMP tackles these problems via **iterative denoising**.
 - We will write the iteration- t **denoiser** as $\boldsymbol{\eta}^t(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^n$.
- Each method defines the denoiser $\boldsymbol{\eta}^t(\cdot)$ differently:
 - Optimization: $\boldsymbol{\eta}^t(\mathbf{r}) = \arg \min_{\mathbf{x}} \{R(\mathbf{x}) + \frac{1}{2\nu^t} \|\mathbf{x} - \mathbf{r}\|_2^2\} \triangleq \text{“prox}_{R\nu^t}(\mathbf{r})\text{”}$
 - MMSE: $\boldsymbol{\eta}^t(\mathbf{r}) = \mathbb{E} \{ \mathbf{x} \mid \mathbf{r} = \mathbf{x} + \mathcal{N}(\mathbf{0}, \nu^t) \}$
 - Plug-and-play: $\boldsymbol{\eta}^t(\mathbf{r}) = \text{BM3D}(\mathbf{r}, \nu^t)$
 - Deep network: $\boldsymbol{\eta}^t(\mathbf{r})$ is learned from training data.

The AMP Algorithm

initialize $\hat{\mathbf{x}}^0 = \mathbf{0}$, $\mathbf{v}^{-1} = \mathbf{0}$

for $t = 0, 1, 2, \dots$

$$\mathbf{v}^t = \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}^t + \frac{N}{M}\mathbf{v}^{t-1} \operatorname{div}(\boldsymbol{\eta}^{t-1}(\hat{\mathbf{x}}^{t-1} + \mathbf{A}^\top \hat{\mathbf{v}}^{t-1})) \quad \text{corrected residual}$$

$$\hat{\mathbf{x}}^{t+1} = \boldsymbol{\eta}^t(\hat{\mathbf{x}}^t + \mathbf{A}^\top \mathbf{v}^t) \quad \text{denoising}$$

where

$$\operatorname{div}(\boldsymbol{\eta}^t(\mathbf{r})) \triangleq \frac{1}{n} \operatorname{tr} \left(\frac{\partial \boldsymbol{\eta}^t(\mathbf{r})}{\partial \mathbf{r}} \right) \quad \text{“divergence.”}$$

Note:

- Original version proposed by [Donoho, Maleki, and Montanari in 2009](#).
 - They considered “scalar” denoisers, such that $[\boldsymbol{\eta}^t(\mathbf{r})]_j = \eta^t(r_j) \forall j$
 - For scalar denoisers, $\operatorname{div}(\boldsymbol{\eta}^t(\mathbf{r})) = \frac{1}{n} \sum_{j=1}^n \eta^{t'}(r_j)$
- Can be recognized as iterative shrinkage/thresholding² plus “[Onsager correction](#).”
- Can be derived using Gaussian & Taylor-series approximations of loopy belief-propagation (hence “AMP”).

²Chambolle, DeVore, Lee, Lucier'98

AMP's Denoising Property

Original AMP Assumptions

- $\mathbf{A} \in \mathbb{R}^{m \times n}$ is drawn i.i.d. Gaussian
- $m, n \rightarrow \infty$ s.t. $\frac{m}{n} \rightarrow \delta \in (0, \infty)$... “large-system limit”
- $[\boldsymbol{\eta}^t(\mathbf{r})]_j = \eta^t(r_j)$ with Lipschitz $\eta(\cdot)$... “scalar denoising”

Under these assumptions, the denoiser's input $\mathbf{r}^t \triangleq \widehat{\mathbf{x}}^t + \mathbf{A}^\top \mathbf{v}^t$ obeys³

$$r_j^t = x_{o,j} + \mathcal{N}(0, \nu_r^t)$$

- That is, r^t is a Gaussian-noise corrupted version of the true signal x_o .
- It should now be clear why we think of $\eta^t(\cdot)$ as a “denoiser.”

Furthermore, the effective noise variance can be consistently estimated:

$$\widehat{\nu}_r^t \triangleq \frac{1}{m} \|\mathbf{v}^t\|^2 \longrightarrow \nu_r^t.$$

³Bayati, Montanari'11

AMP's State Evolution

- Assume that the measurements \mathbf{y} were generated via

$$\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathcal{N}(\mathbf{0}, \nu_w \mathbf{I})$$

where \mathbf{x}_o empirically converges to some random variable X_o as $n \rightarrow \infty$.

- Define the iteration- t mean-squared error (MSE)

$$\mathcal{E}^t \triangleq \frac{1}{n} \|\hat{\mathbf{x}}^t - \mathbf{x}_o\|^2.$$

- Under above assumptions, AMP obeys the following state evolution (SE):⁴

for $t = 0, 1, 2, \dots$

$$\nu_r^t = \nu_w + \frac{n}{m} \mathcal{E}^t$$

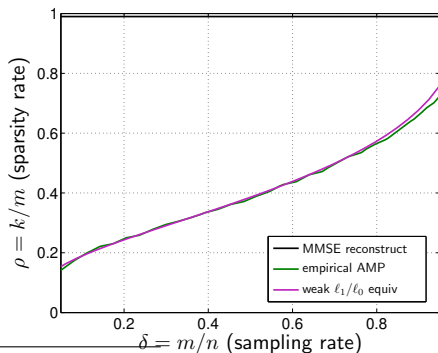
$$\mathcal{E}^{t+1} = \mathbb{E} \left\{ \left[\eta^t(X_o + \mathcal{N}(0, \nu_r^t)) - X_o \right]^2 \right\}$$

⁴Bayati, Montanari'11

Achievability Analysis via the AMP SE

- AMP's SE can be applied to analyze achievability in various problems.
- E.g., it yields a **closed-form expression**⁵ for the sparsity/sampling region where ℓ_1 -penalized regression is equivalent to ℓ_0 -penalized regression:

$$\rho(\delta) = \max_{c>0} \frac{1 - 2\delta^{-1}[(1 + c^2)\Phi(-c) - c\phi(c)]}{1 + c^2 - 2[(1 + c^2)\Phi(-c) - c\phi(c)]},$$



⁵Donoho, Maleki, Montanari'09

MMSE Optimality of AMP

- Now suppose that the AMP Assumptions hold, and that

$$\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathcal{N}(\mathbf{0}, \nu_w \mathbf{I}),$$

where the elements of \mathbf{x}_o are **i.i.d. draws** of some random variable X_o .

- Suppose also that $\eta^t(\cdot)$ is the **MMSE denoiser**, i.e.,

$$\eta^t(R) = \mathbb{E} \{X_o \mid R = X_o + \mathcal{N}(0, \nu_r^t)\}$$

- Then, if the state evolution has a **unique** fixed point, the MSE of $\hat{\mathbf{x}}^t$ converges⁶ to the **replica prediction of the MMSE** as $t \rightarrow \infty$.
- Under the AMP Assumptions, the replica prediction of the MMSE was shown to be **correct**.⁷⁸

⁶Bayati, Montanari'11, ⁷Reeves, Pfister'16, ⁸Barbier, Dia, Macris, Krzakala'16

Universality of AMP State Evolution

- Until now, it was assumed that \mathbf{A} is drawn i.i.d. Gaussian.
- The state evolution also holds when \mathbf{A} is drawn from i.i.d. A_{ij} such that

$$\mathbb{E}\{A_{ij}\} = 0$$

$$\mathbb{E}\{A_{ij}^2\} = 1/m$$

$$\mathbb{E}\{A_{ij}^6\} = C/m \text{ for some fixed } C > 0.$$

often abbreviated as “sub-Gaussian A_{ij} .”

- The proof⁹ assumes polynomial scalar denoising $\eta^t(\cdot)$ of bounded order.

⁹Bayati, Lelarge, Montanari'15

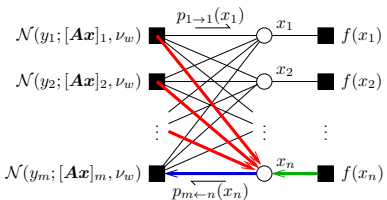
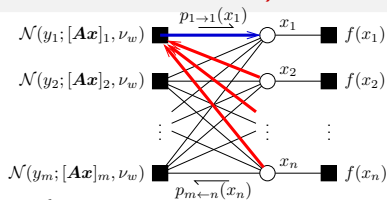
Deriving AMP via Loopy BP (e.g., sum-product alg)

- 1 Message from y_i node to x_j node:

$$\begin{aligned}
 p_{i \rightarrow j}(x_j) &\approx \mathcal{N} \text{ via CLT} \\
 p_{i \rightarrow j}(x_j) &\propto \int_{\{x_l\}_{l \neq j}} \mathcal{N}(y_i; \underbrace{\sum_l a_{il} x_l}_{z_i}, \nu_w) \prod_{l \neq j} p_{i \leftarrow l}(x_l) \\
 &\approx \int_{z_i} \mathcal{N}(y_i; z_i, \nu_w) \mathcal{N}(z_i; \hat{z}_i(x_j), \nu_i^z(x_j)) \sim \mathcal{N}
 \end{aligned}$$

To compute $\hat{z}_i(x_j), \nu_i^z(x_j)$, the means and variances of $\{p_{i \leftarrow l}\}_{l \neq j}$ suffice, implying Gaussian message passing, similar to expectation-propagation. Remaining problem: we have $2mn$ messages to compute (too many!).

- 2 Exploiting similarity among the messages $\{p_{i \leftarrow j}\}_{i=1}^m$, AMP employs a **Taylor-series approximation** of their difference whose error vanishes as $m \rightarrow \infty$ for dense \mathbf{A} (and similar for $\{p_{i \leftarrow j}\}_{j=1}^n$ as $n \rightarrow \infty$). Finally, need to compute only $O(m+n)$ messages!



Understanding AMP

- The belief-propagation derivation of AMP provides very little insight!
 - Loopy BP is suboptimal, even if implemented exactly
 - The i.i.d. property of \mathcal{A} is never used in the derivation
- And the rigorous proofs of AMP's state evolution are very technical!
- As a middle ground, we suggest an alternate derivation that gives insight into how and why AMP works.
 - Based on the idea of “[first-order cancellation](#)”
 - We will assume equiprobable Bernoulli $a_{ij} \in \pm 1/\sqrt{m}$ and polynomial $\eta(\cdot)$

AMP as First-Order Cancellation

Recall the AMP recursion:

$$\begin{aligned} \mathbf{v}^t &= \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}^t + \frac{n}{m} \mathbf{v}^{t-1} \operatorname{div}(\eta(\mathbf{r}^{t-1})) \\ \hat{\mathbf{x}}^{t+1} &= \eta(\underbrace{\hat{\mathbf{x}}^t + \mathbf{A}^\top \mathbf{v}^t}_{\triangleq \mathbf{r}^t}) \end{aligned}$$

Notice that

$$\begin{aligned} [\mathbf{A}\hat{\mathbf{x}}^t]_i &= \mathbf{a}_i^\top \eta(\hat{\mathbf{x}}^{t-1} + \sum_l \mathbf{a}_l v_l^{t-1}) \quad \text{where } \mathbf{a}_i^\top \text{ is the } i\text{th row of } \mathbf{A} \\ &= \mathbf{a}_i^\top \eta(\underbrace{\hat{\mathbf{x}}^{t-1} + \sum_{l \neq i} \mathbf{a}_l v_l^{t-1}}_{\triangleq \mathbf{r}_i^{t-1}} + \mathbf{a}_i v_i^{t-1}) \\ &\quad \triangleq \mathbf{r}_i^{t-1}, \text{ which removes the direct contribution of } \mathbf{a}_i \text{ from } \mathbf{r}^{t-1} \\ &= \mathbf{a}_i^\top [\eta(\mathbf{r}_i^{t-1}) + \frac{\partial \eta}{\partial \mathbf{r}}(\mathbf{r}_i^{t-1}) \mathbf{a}_i v_i^{t-1} + O(1/m)] \quad \text{using a Taylor expansion} \\ &= \mathbf{a}_i^\top \eta(\mathbf{r}_i^{t-1}) + v_i^{t-1} \sum_j \mathbf{a}_{ij}^2 \eta'(\mathbf{r}_{ij}^{t-1}) + O(1/\sqrt{m}) \\ &= \mathbf{a}_i^\top \eta(\mathbf{r}_i^{t-1}) + \frac{n}{m} v_i^{t-1} \underbrace{\frac{1}{n} \sum_j \eta'(\mathbf{r}_{ij}^{t-1})}_{\operatorname{div}(\eta(\mathbf{r}_i^{t-1}))} + O(1/\sqrt{m}) \quad \text{since } \mathbf{a}_{ij}^2 = 1/m \quad \forall ij \end{aligned}$$

which uncovers the Onsager correction.

AMP as First-Order Cancellation (cont.)

Now use $[\mathbf{A}\hat{\mathbf{x}}^t]_i$ to study j th component of denoiser input error $\mathbf{e}^t \triangleq \mathbf{r}^t - \mathbf{x}_o$:

$$\begin{aligned} e_j^t &= \sum_i a_{ij} \sum_{l \neq j} a_{il} [x_{o,l} - \eta(r_{il}^{t-1})] + \sum_i a_{ij} w_i \\ &+ \sum_i a_{ij} \left[\frac{n}{m} v_i^{t-1} \operatorname{div}(\eta(\mathbf{r}^{t-1})) - \frac{n}{m} v_i^{t-1} \operatorname{div}(\eta(\mathbf{r}_i^{t-1})) \right] + O(1/\sqrt{m}) \end{aligned}$$

where the divergence difference can be absorbed into the $O(1/\sqrt{m})$ term...

$$\begin{aligned} &= \underbrace{\sum_i a_{ij} \sum_{l \neq j} a_{il} \underbrace{[x_{o,l} - \eta(r_{il}^{t-1})]}_{\triangleq \epsilon_{il}}}_{\sim \mathcal{N}(0, \frac{1}{m^2} \sum_i \sum_{l \neq j} (\epsilon_{il}^t)^2)} + \underbrace{\sum_i a_{ij} w_i}_i + O(1/\sqrt{m}) \\ &\sim \mathcal{N}(0, \frac{1}{m^2} \sum_i \sum_{l \neq j} (\epsilon_{il}^t)^2) \sim \mathcal{N}(0, \frac{1}{m} \sum_i w_i^2) \end{aligned}$$

using the CLT and assuming independence of $\{a_{il}\}_{l=1}^n$ and $\{r_{il}^{t-1}\}_{l=1}^n$

$$\sim \mathcal{N}(0, \frac{n}{m} \mathcal{E}^{(t)} + \nu_w) + O(1/\sqrt{m}) \quad \dots \text{the AMP state evolution}$$

$$\text{where } \mathcal{E}^{(t)} \triangleq \frac{1}{n} \sum_{j=1}^n [x_{o,j} - \hat{x}_j^{(t)}]^2 \text{ and } \nu_w \triangleq \frac{1}{m} \sum_{i=1}^m w_i^2$$

AMP with Non-Separable Denoisers

- Until now, we have focused on **separable** denoisers, i.e., $[\eta^t(\mathbf{r})]_j = \eta^t(r_j) \forall j$
- Can we use sophisticated **non-separable** $\eta(\cdot)$ with AMP?
- Yes! Many examples. . .
 - Markov chain,¹⁰ Markov field,¹² Markov tree,¹² denoisers in 2010
 - Blockwise & TV denoising considered by Donoho, Johnstone, Montanari in 2011
 - BM3D denoising considered by Metzler, Maleki, Baraniuk in 2015
- Rigorous state-evolution proven by Berthier, Montanari, Nguyen in 2017.
 - Assumes \mathbf{A} drawn i.i.d. Gaussian
 - Assumes η is Lipschitz and “convergent under Gaussian inputs”

¹⁰S'10, ¹¹Som,S'11, ¹²Som,S'12

AMP at Large but Finite Dimensions

- Until now, we have focused on the large-system limit $m, n \rightarrow \infty$ with $m/n \rightarrow \delta \in (0, \infty)$
- The **non-asymptotic** case was analyzed by Rush and Venkataramanan.¹³
- They showed that probability of ϵ -deviation between the finite and limiting SE falls exponentially in m , as long as the number of iterations $t < o(\frac{\log n}{\log \log n})$

¹³Rush, Venkataramanan'18

AMP Summary: The good, the bad, and the ugly

The good:

- With **large i.i.d. sub-Gaussian \mathbf{A}** , AMP is rigorously characterized by a scalar **state-evolution** whose fixed points, when unique, are **MMSE optimal** under proper choice of denoiser.
- **Empirically, AMP behaves well** with many other “sufficiently random” \mathbf{A} (e.g., randomly sub-sampled Fourier \mathbf{A} & i.i.d. sparse x).

The bad:

- With **general \mathbf{A}** , AMP gives **no guarantees**.

The ugly:

- With **some \mathbf{A}** , AMP may **fail to converge!** (e.g., ill-conditioned or non-zero-mean \mathbf{A})



Outline

- 1 Linear Regression
- 2 Approximate Message Passing (AMP)
- 3 Vector AMP (VAMP)**
- 4 Unfolding AMP and VAMP into Deep Neural Networks
- 5 Extensions: GLMs, Parameter Learning, Bilinear Problems

Vector AMP (VAMP)

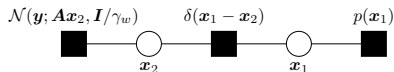
- Recall goal is linear regression: Recover \mathbf{x}_o from $\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w)$.
 - Now it will be easier to work with inverse variances, i.e., **precisions**
- VAMP is like AMP in many ways, but **supports a larger class of random matrices**.

- VAMP yields a precise analysis for **right-orthogonally invariant \mathbf{A}** :

$$\text{svd}(\mathbf{A}) = \mathbf{U}\mathbf{S}\mathbf{V}^T \text{ for } \begin{cases} \mathbf{U}: \text{deterministic orthogonal} \\ \mathbf{S}: \text{deterministic diagonal} \\ \mathbf{V}: \text{"Haar;" uniform on set of orthogonal matrices} \end{cases}$$

of which i.i.d. Gaussian is a special case.

- Can be derived as a form of message passing on a vector-valued factor graph.



VAMP: The Algorithm

With SVD $\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^T$, damping $\zeta \in (0, 1]$, and Lipschitz $\boldsymbol{\eta}_1^t(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^n$.

Initialize \mathbf{r}_1, γ_1 .

For $t = 1, 2, 3, \dots$

$$\hat{\mathbf{x}}_1 \leftarrow \boldsymbol{\eta}_1^t(\mathbf{r}_1) \quad \text{denoising of } \mathbf{r}_1 = \mathbf{x}_o + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_1)$$

$$\xi_1 \leftarrow \gamma_1 / \text{div}(\boldsymbol{\eta}_1^t(\mathbf{r}_1))$$

$$\mathbf{r}_2 \leftarrow (\xi_1 \hat{\mathbf{x}}_1 - \gamma_1 \mathbf{r}_1) / (\xi_1 - \gamma_1) \quad \text{ Onsager correction}$$

$$\gamma_2 \leftarrow \xi_1 - \gamma_1$$

$$\hat{\mathbf{x}}_2 \leftarrow \boldsymbol{\eta}_2(\mathbf{r}_2; \gamma_2) \quad \text{LMMSE estimate of } \mathbf{x} \sim \mathcal{N}(\mathbf{r}_2, \mathbf{I}/\gamma_2)$$

$$\xi_2 \leftarrow \gamma_2 / \text{div}(\boldsymbol{\eta}_2(\mathbf{r}_2; \gamma_2)) \quad \text{from } \mathbf{y} = \mathbf{A}\mathbf{x} + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w)$$

$$\mathbf{r}_1 \leftarrow \zeta(\xi_2 \hat{\mathbf{x}}_2 - \gamma_2 \mathbf{r}_2) / (\xi_2 - \gamma_2) + (1 - \zeta) \mathbf{r}_1 \quad \text{ Onsager correction}$$

$$\gamma_1 \leftarrow \zeta(\xi_2 - \gamma_2) + (1 - \zeta) \gamma_1 \quad \text{damping}$$

where $\boldsymbol{\eta}_2(\mathbf{r}_2; \gamma_2) = (\gamma_w \mathbf{A}^T \mathbf{A} + \gamma_2 \mathbf{I})^{-1} (\gamma_w \mathbf{A}^T \mathbf{y} + \gamma_2 \mathbf{r}_2)$
 $= \mathbf{V} (\gamma_w \text{Diag}(\mathbf{s})^2 + \gamma_2 \mathbf{I})^{-1} (\gamma_w \text{Diag}(\mathbf{s}) \mathbf{U}^T \mathbf{y} + \gamma_2 \mathbf{V}^T \mathbf{r}_2)$
 $\xi_2 = [\frac{1}{n} \sum_{j=1}^n (\gamma_w s_j^2 + \gamma_2)^{-1}]^{-1}$ **two mat-vec mults per iteration!**

VAMP's Denoising Property

Original VAMP Assumptions

- $\mathbf{A} \in \mathbb{R}^{m \times n}$ is right-orthogonally invariant
- $m, n \rightarrow \infty$ s.t. $m/n \rightarrow \delta \in (0, \infty)$... “large-system limit”
- $[\boldsymbol{\eta}_1^t(\mathbf{r})]_j = \eta_1^t(r_j)$ with Lipschitz $\eta_1^t(\cdot)$... “separable denoising”

Under Assumption 2, the elements of the denoiser's input \mathbf{r}_1^t obey¹⁴

$$r_{1,j}^t = x_{o,j} + \mathcal{N}(0, \nu_1^t)$$

- That is, \mathbf{r}_1^t is a Gaussian-noise corrupted version of the true signal \mathbf{x}_o .
- As with AMP, we can interpret $\boldsymbol{\eta}_1(\cdot)$ as a “denoiser.”

¹⁴Rangan, S, Fletcher'16

VAMP's State Evolution

Assume empirical convergence of $\{s_j\} \rightarrow S$ and $\{(r_{1,j}^0, x_{o,j})\} \rightarrow (R_1^0, X_o)$, and define

$$\mathcal{E}_i^t \triangleq \frac{1}{n} \|\widehat{\mathbf{x}}_i^t - \mathbf{x}_o\|^2 \text{ for } i = 1, 2.$$

Then under the VAMP Assumptions, VAMP obeys the following state-evolution:

for $t = 0, 1, 2, \dots$

$$\mathcal{E}_1^t = \mathbb{E} \left\{ \left[\eta_1^t (X_o + \mathcal{N}(0, \nu_1^t)) - X_o \right]^2 \right\} \quad \text{MSE}$$

$$\alpha_1^t = \mathbb{E} \left\{ \eta_1^{t'} (X_o + \mathcal{N}(0, \nu_1^t)) \right\} \quad \text{divergence}$$

$$\gamma_2^t = \gamma_1^t \frac{1 - \alpha_1^t}{\alpha_1^t}, \quad \nu_2^t = \frac{1}{(1 - \alpha_1^t)^2} \left[\mathcal{E}_1^t - (\alpha_1^t)^2 \nu_1^t \right]$$

$$\mathcal{E}_2^t = \mathbb{E} \left\{ \left[\gamma_w S^2 + \gamma_2^t \right]^{-1} \right\} \quad \text{MSE}$$

$$\alpha_2^t = \gamma_2^t \mathbb{E} \left\{ \left[\gamma_w S^2 + \gamma_2^t \right]^{-1} \right\} \quad \text{divergence}$$

$$\gamma_1^{t+1} = \gamma_2^t \frac{1 - \alpha_2^t}{\alpha_2^t}, \quad \nu_1^{t+1} = \frac{1}{(1 - \alpha_2^t)^2} \left[\mathcal{E}_2^t - (\alpha_2^t)^2 \nu_2^t \right]$$

Note: Above equations assume $\eta_2(\cdot)$ uses true noise precision γ_w .
If not, there are more complicated expressions for \mathcal{E}_2^t and α_2^t .

MMSE Optimality of VAMP

- Now suppose that the VAMP Assumptions hold, and that

$$\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_w),$$

where the elements of \mathbf{x}_o are i.i.d. draws of some random variable X_o .

- Suppose also that $\eta_1^t(\cdot)$ is the MMSE denoiser, i.e.,

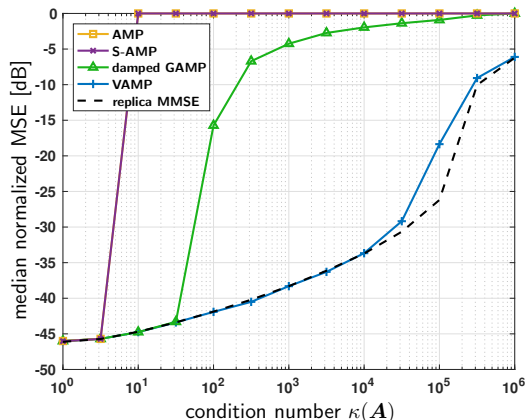
$$\eta_1^t(R_1) = \mathbb{E} \{ X_o \mid R_1 = X_o + \mathcal{N}(0, \nu_1^t) \}$$

- Then, if the state evolution has a **unique** fixed point, the MSE of $\hat{\mathbf{x}}_1^t$ converges¹⁵ to the **replica prediction**¹⁶ of the **MMSE** as $t \rightarrow \infty$.

¹⁵Rangan, S, Fletcher'16, ¹⁶Tulino, Caire, Verdu, Shamai'13

Experiment with MMSE Denoising

Comparison of several algorithms¹⁷ with MMSE denoising.



$n = 1024$
 $m/n = 0.5$

$\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^T$
 $\mathbf{U}, \mathbf{V} \sim \text{Haar}$
 $s_j/s_{j-1} = \phi \forall j$
 ϕ determines $\kappa(\mathbf{A})$

$X_o \sim \text{Bernoulli-Gaussian}$
 $\Pr\{X_o \neq 0\} = 0.1$

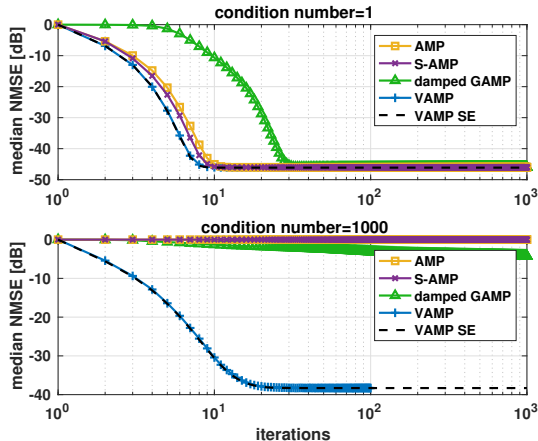
SNR = 40dB

VAMP achieves the replica MMSE over a wide range of condition numbers.

¹⁷S-AMP: Cakmak, Fleury, Winther'14, damped GAMP: Vila, S, Rangan, Krzakala, Zdeborová'15

Experiment with MMSE Denoising (cont.)

Comparison of several algorithms with priors matched to data.



$$n = 1024$$

$$m/n = 0.5$$

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

$$\mathbf{U}, \mathbf{V} \sim \text{Haar}$$

$$s_j/s_{j-1} = \phi \quad \forall j$$

$$\phi \text{ determines } \kappa(\mathbf{A})$$

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

$$\Pr\{X_o \neq 0\} = 0.1$$

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

VAMP is relative fast even when \mathbf{A} is ill-conditioned.

VAMP for Optimization

- Consider the optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + R(\mathbf{x}) \right\}$$

where $R(\cdot)$ is strictly convex and \mathbf{A} is arbitrary (e.g., not necessarily RRI).

- If we choose the denoiser

$$\eta_1^t(\mathbf{r}) = \arg \min_{\mathbf{x}} \left\{ R(\mathbf{x}) + \frac{\gamma_1^t}{2} \|\mathbf{x} - \mathbf{r}\|^2 \right\} = \text{prox}_{R/\gamma_1^t}(\mathbf{r})$$

and the damping parameter

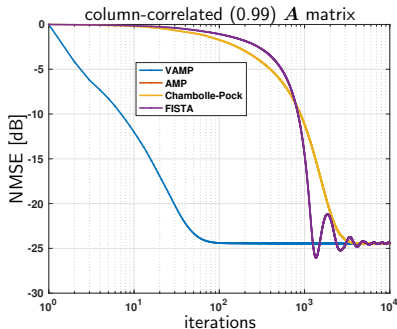
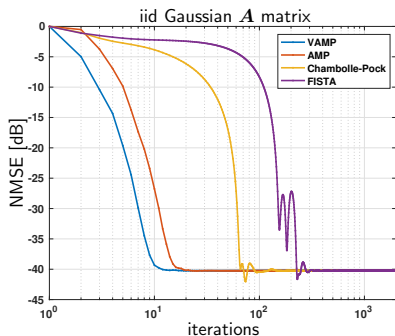
$$\zeta \leq \frac{2 \min\{\gamma_1, \gamma_2\}}{\gamma_1 + \gamma_2},$$

then a double-loop version of VAMP converges¹⁸ to $\hat{\mathbf{x}}$ from above.

- Furthermore, if the γ_1 and γ_2 variables are fixed over the iterations, then VAMP reduces to the Peaceman-Rachford variant of ADMM.

¹⁸Fletcher, Sahraee, Rangan, S'16

Example of AMP & VAMP on the LASSO Problem



Solving LASSO to reconstruct 40-sparse $\mathbf{x} \in \mathbb{R}^{1000}$ from noisy $\mathbf{y} \in \mathbb{R}^{400}$.

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\}.$$

Deriving VAMP from EC

- Ideally, we would like to compute the exact **posterior density**

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x})\ell(\mathbf{x}; \mathbf{y})}{Z(\mathbf{y})} \quad \text{for } Z(\mathbf{y}) \triangleq \int p(\mathbf{x})\ell(\mathbf{x}; \mathbf{y}) d\mathbf{x},$$

but the high-dimensional integral in $Z(\mathbf{y})$ is difficult to compute.

- We might try to circumvent $Z(\mathbf{y})$ through **variational optimization**:

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \arg \min_b D(b(\mathbf{x})\|p(\mathbf{x}|\mathbf{y})) \quad \text{where } D(\cdot\|\cdot) \text{ is KL divergence} \\ &= \arg \min_b \underbrace{D(b(\mathbf{x})\|p(\mathbf{x})) + D(b(\mathbf{x})\|\ell(\mathbf{x}; \mathbf{y})) + H(b(\mathbf{x}))}_{\text{Gibbs free energy}} \\ &= \arg \min_{b_1, b_2, q} \underbrace{D(b_1(\mathbf{x})\|p(\mathbf{x})) + D(b_2(\mathbf{x})\|\ell(\mathbf{x}; \mathbf{y})) + H(q(\mathbf{x}))}_{\triangleq J_{\text{Gibbs}}(b_1, b_2, q)} \\ &\quad \text{s.t. } b_1 = b_2 = q, \end{aligned}$$

but the density constraint keeps the problem difficult.

Deriving VAMP from EC (cont.)

- In **expectation-consistent approximation (EC)**¹⁹, the density constraint is relaxed to moment-matching constraints:

$$p(\mathbf{x}|\mathbf{y}) \approx \arg \min_{b_1, b_2, q} J_{\text{Gibbs}}(b_1, b_2, q)$$

$$\text{s.t.} \quad \begin{cases} \mathbb{E}\{\mathbf{x}|b_1\} = \mathbb{E}\{\mathbf{x}|b_2\} = \mathbb{E}\{\mathbf{x}|q\} \\ \text{tr}(\text{Cov}\{\mathbf{x}|b_1\}) = \text{tr}(\text{Cov}\{\mathbf{x}|b_2\}) = \text{tr}(\text{Cov}\{\mathbf{x}|q\}). \end{cases}$$

- The **stationary points** of EC are the densities

$$\begin{aligned} b_1(\mathbf{x}) &\propto p(\mathbf{x})\mathcal{N}(\mathbf{x}; \mathbf{r}_1, \mathbf{I}/\gamma_1) \\ b_2(\mathbf{x}) &\propto \ell(\mathbf{x}; \mathbf{y})\mathcal{N}(\mathbf{x}; \mathbf{r}_2, \mathbf{I}/\gamma_2) \\ q(\mathbf{x}) &= \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \mathbf{I}/\xi) \end{aligned} \quad \text{s.t.} \quad \begin{cases} \mathbb{E}\{\mathbf{x}|b_1\} = \mathbb{E}\{\mathbf{x}|b_2\} = \hat{\mathbf{x}} \\ \frac{1}{n}\text{tr}(\text{Cov}\{\mathbf{x}|b_1\}) = \frac{1}{n}\text{tr}(\text{Cov}\{\mathbf{x}|b_2\}) = \frac{1}{\xi} \end{cases}$$

- VAMP iteratively solves for the quantities $\mathbf{r}_1, \gamma_1, \mathbf{r}_2, \gamma_2, \hat{\mathbf{x}}, \xi$ above.
 - Leads to $\eta_1^t(\cdot)$ being the MMSE denoiser of $\mathbf{r}_1 = \mathbf{x}_o + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_1^t)$
 - In this setting, VAMP is simply an instance of **expectation propagation (EP)**²⁰.
 - But VAMP is more general than EP, in that it allows non-MMSE denoisers η_1 .

¹⁹Opper, Winther'04, ²⁰Minka'01

Plug-and-play VAMP

- Recall the scalar denoising step of VAMP (or AMP):

$$\hat{\mathbf{x}}_1 = \eta_1^t(\mathbf{r}_1) \quad \text{where } \mathbf{r}_1 = \mathbf{x}_o + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_1^t)$$

- For many signal classes (e.g., images), very sophisticated *non-separable* denoisers $\eta_1(\cdot)$ have been developed (e.g., [BM3D](#), [DnCNN](#)).
- These non-separable denoisers can be “[plugged into](#)” VAMP!
- Their divergence can be approximated via Monte Carlo²¹

$$\text{div}(\eta_1^t(\mathbf{r})) \approx \frac{1}{K} \sum_{k=1}^K \frac{\mathbf{p}_k^\top [\eta_1^t(\mathbf{r} + \epsilon \mathbf{p}_k) - \eta_1^t(\mathbf{r})]}{n\epsilon}$$

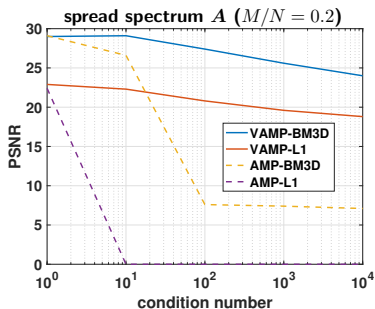
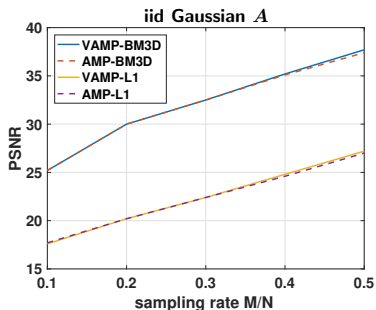
with random vectors $\mathbf{p}_k \in \{\pm 1\}^n$ and small $\epsilon > 0$. Empirically, $K=1$ suffices.

- A [rigorous state-evolution](#) has been established for plug-and-play VAMP.²²

²¹Ramani, Blu, Unser'08, ²²Fletcher, Rangan, Sarker, S'18

Experiment: Compressive Image Recovery with BM3D

Plug-and-play versions of VAMP and AMP behave similarly with i.i.d. Gaussian \mathbf{A} is i.i.d., but VAMP can handle a larger class of random matrices \mathbf{A} .



Results above are averaged over 128×128 versions of

lena, barbara, boat, fingerprint, house, peppers

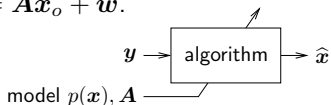
and 10 random realizations of \mathbf{A}, \mathbf{w} .

Outline

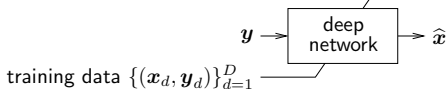
- 1 Linear Regression
- 2 Approximate Message Passing (AMP)
- 3 Vector AMP (VAMP)
- 4 Unfolding AMP and VAMP into Deep Neural Networks
- 5 Extensions: GLMs, Parameter Learning, Bilinear Problems

Deep learning for sparse reconstruction

- Until now we've focused on **designing algorithms** to recover $\mathbf{x}_o \sim p(\mathbf{x})$ from measurements $\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathbf{w}$.



- What about **training deep networks** to predict \mathbf{x}_o from \mathbf{y} ?
Can we increase accuracy and/or decrease computation?



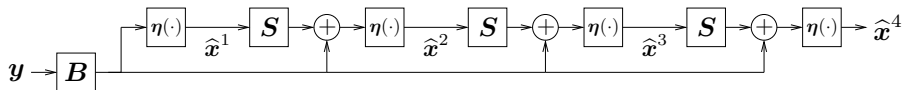
- Are there **connections** between these approaches?

Unfolding Algorithms into Networks

Consider, e.g., the classical sparse-reconstruction algorithm, [ISTA](#).²³

$$\boxed{\begin{aligned} \mathbf{v}^t &= \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}^t \\ \hat{\mathbf{x}}^{t+1} &= \eta(\hat{\mathbf{x}}^t + \mathbf{A}^\top \mathbf{v}^t) \end{aligned}} \Leftrightarrow \boxed{\hat{\mathbf{x}}^{t+1} = \eta(\mathbf{S}\hat{\mathbf{x}}^t + \mathbf{B}\mathbf{y}) \text{ with } \begin{aligned} \mathbf{S} &\triangleq \mathbf{I} - \mathbf{A}^\top \mathbf{A} \\ \mathbf{B} &\triangleq \mathbf{A}^\top \end{aligned}}$$

Gregor & LeCun²⁴ proposed to “[unfold](#)” it into a deep net and “[learn](#)” improved parameters using training data, yielding “[learned ISTA](#)” (LISTA):



The same “[unfolding & learning](#)” idea can be used to improve AMP, yielding “[learned AMP](#)” (LAMP).²⁵

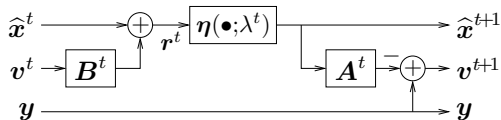
²³Chambolle, DeVore, Lee, Lucier’98.

²⁴Gregor, LeCun’10.

²⁵Borgerding, S’16.

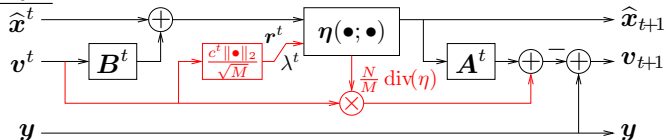
Onsager-Corrected Deep Networks

t^{th} LISTA layer:



to exploit low-rank $B^t A^t$ in linear stage $S^t = I - B^t A^t$.

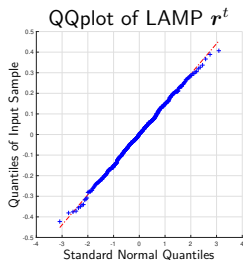
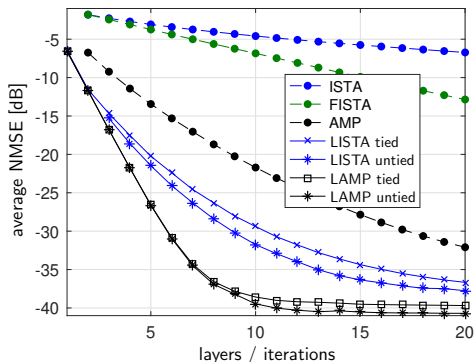
t^{th} LAMP layer:



Onsager correction now aims to decouple errors across layers.

LAMP performance with soft-threshold denoising

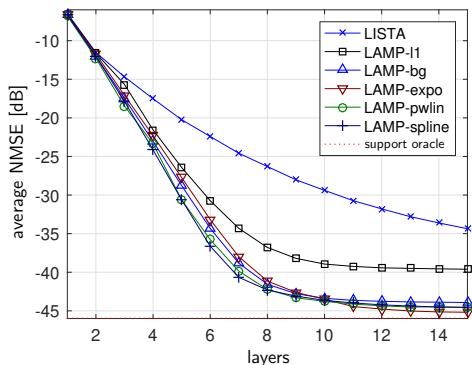
LISTA beats AMP, FISTA, ISTA
LAMP beats LISTA in convergence speed and asymptotic MSE.



LAMP beyond soft-thresholding

So far, we used [soft-thresholding](#) to isolate the effects of Onsager correction.

What happens with [more sophisticated \(learned\) denoisers](#)?



Here we learned the parameters of these denoiser families:

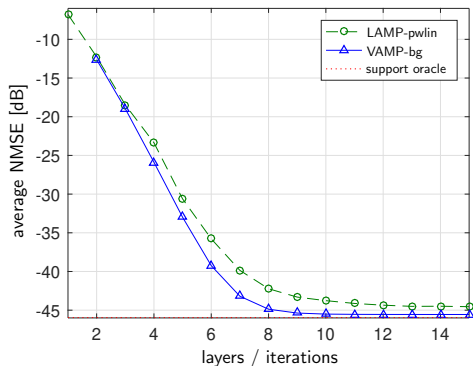
- scaled soft-thresholding
- conditional mean under BG
- Exponential kernel²⁶
- Piecewise Linear²⁶
- Spline²⁷

Big improvement!

²⁶Guo, Davies'15. ²⁷Kamilov, Mansour'16.

LAMP versus VAMP

How does our best **Learned AMP** compare to MMSE **VAMP**?



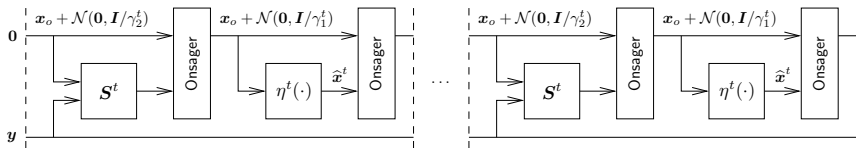
VAMP wins!

So what about “learned VAMP”?

Learned VAMP



- Suppose we **unfold** VAMP and **learn (via backprop)** the parameters $\{\mathbf{S}^t, \eta^t\}_{t=1}^T$ that minimize the training MSE.



- Remarkably, **backpropagation learns the parameters prescribed by VAMP!**

Theory explains the deep network!

- Onsager correction **decouples** the design of $\{\mathbf{S}^t, \eta^t(\cdot)\}_{t=1}^T$:
 Layer-wise optimal $\mathbf{S}^t, \eta^t(\cdot) \Rightarrow$ Network optimal $\{\mathbf{S}^t, \eta^t(\cdot)\}_{t=1}^T$

Outline

- 1 Linear Regression
- 2 Approximate Message Passing (AMP)
- 3 Vector AMP (VAMP)
- 4 Unfolding AMP and VAMP into Deep Neural Networks
- 5 Extensions: GLMs, Parameter Learning, Bilinear Problems

Generalized linear models

- Until now we have considered the **standard linear model**: $\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathbf{w}$.
- One may also consider the **generalized linear model** (GLM), where

$$\mathbf{y} \sim p(\mathbf{y}|\mathbf{z}) \text{ with hidden } \mathbf{z} = \mathbf{A}\mathbf{x}_o$$

which supports, e.g.,

- $y_i = z_i + w_i$: additive, possibly **non-Gaussian noise**
 - $y_i = Q(z_i + w_i)$: **quantization**
 - $y_i = \text{sgn}(z_i + w_i)$: **binary classification**
 - $y_i = |z_i + w_i|$: **phase retrieval**
 - Poisson y_i : **photon-limited imaging**
- For this, there is a **Generalized AMP**²⁹ with a rigorous state evolution.³⁰
 - There is also a **Generalized VAMP**³¹ with a rigorous state evolution.³²

²⁹Rangan'11, ³⁰Javanmard, Montanari'12, ³¹S, Fletcher, Rangan'16. ³²Fletcher, Rangan, S'18.

Parameter learning

- Consider inference under prior $p(\mathbf{x}; \boldsymbol{\theta}_1)$ and likelihood $\ell(\mathbf{x}; \mathbf{y}, \boldsymbol{\theta}_2)$, where the hyperparameters $\boldsymbol{\theta} \triangleq [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2]$ are unknown.
 - $\boldsymbol{\theta}_1$ might specify sparsity rate, or all parameters of a GMM
 - $\boldsymbol{\theta}_2$ might specify the measurement noise variance, or forward model \mathbf{A}

- EM-inspired extensions of (G)AMP and (G)VAMP that simultaneously estimate \mathbf{x} and learn $\boldsymbol{\theta}$ from \mathbf{y} have been developed.
 - Have rigorous state evolutions³³³⁴
 - “Adaptive VAMP” yields asymptotically consistent³⁴ estimates of $\boldsymbol{\theta}$

- SURE-based auto-tuning AMP algorithms have also been proposed
 - for LASSO by Mousavi, Maleki, and Baraniuk
 - for parametric separable denoisers by Guo and Davies

³³Kamilov, Rangan, Fletcher, Unser'12, ³⁴Fletcher, Sahraee, Rangan, S'17

Bilinear problems

- So far we have considered (generalized) linear models.
- AMP has also been applied to (generalized) *bilinear models*.
- The typical problem is to recover $\mathbf{B} \in \mathbb{R}^{m \times k}$ and $\mathbf{C} \in \mathbb{R}^{k \times n}$ from ...
 - $\begin{cases} \mathbf{Y} = \mathbf{BC} + \mathbf{W} & \text{(standard bilinear model)} \\ \mathbf{Y} \sim p(\mathbf{Y}|\mathbf{Z}) \text{ for } \mathbf{Z} = \mathbf{BC} & \text{(generalized bilinear model)} \end{cases}$
 - The case where $m, n \rightarrow \infty$ for fixed k is well understood.³⁵ (See Jean's talk)
 - With $m, n, k \rightarrow \infty$, algorithms work (e.g., BiGAMP³⁶) but are not well understood.
- A more general bilinear problem is to recover $\mathbf{b} \in \mathbb{R}^k$ and $\mathbf{c} \in \mathbb{R}^n$ from
 - $\begin{cases} y_i = \mathbf{b}^\top \mathbf{A}_i \mathbf{c} + w_i, & i = 1 \dots m \\ y_i \sim p(y_i|z_i) \text{ for } z_i = \mathbf{b}^\top \mathbf{A}_i \mathbf{c}, & i = 1 \dots m \end{cases}$ where $\{\mathbf{A}_i\}$ are known matrices
 - Algorithms³⁷ and replica analyses³⁸ (for $m, n, k \rightarrow \infty$ and i.i.d. \mathbf{A}_i) exist.

³⁵Montanari, Venkataramanan'17, ³⁶Parker, S, Cevher'14, ³⁷Parker, S'16, ³⁸Schulke, S, Zdeborova'16

Conclusions

- AMP and VAMP are a computationally efficient algorithms for (generalized) linear regression.
- With large random \mathbf{A} , the ensemble behaviors of AMP and VAMP obey rigorous state evolutions whose fixed-points, when unique, agree with the replica predictions of the MMSE.
- AMP and VAMP support nonseparable (i.e., “plug-in”) denoisers, also with rigorous state evolutions.
- For convex optimization problems, VAMP is provably convergent for any \mathbf{A} .
- Extensions of AMP and VAMP cover ...
 - unfolded deep networks
 - the learning of unknown prior/likelihood parameters
 - bilinear problems
- Not discussed: multilayer versions of AMP & VAMP.

References I



S. V. Venkatakrisnan, C. A. Bouman, and B. Wohlberg, "Plug-and-play priors for model based reconstruction," in *Proc. IEEE Global Conf. Signal Info. Process.*, pp. 945–948, 2013.



D. L. Donoho, A. Maleki, and A. Montanari, "Message passing algorithms for compressed sensing," *Proc. Nat. Acad. Sci.*, vol. 106, pp. 18914–18919, Nov. 2009.



A. Chambolle, R. A. DeVore, N. Lee, and B. J. Lucier, "Nonlinear wavelet image processing: Variational problems, compression, and noise removal through wavelet shrinkage," *IEEE Trans. Image Process.*, vol. 7, pp. 319–335, Mar. 1998.



M. Bayati and A. Montanari, "The dynamics of message passing on dense graphs, with applications to compressed sensing," *IEEE Trans. Inform. Theory*, vol. 57, pp. 764–785, Feb. 2011.



G. Reeves and H. D. Pfister, "The replica-symmetric prediction for compressed sensing with Gaussian matrices is exact," in *Proc. IEEE Int. Symp. Inform. Thy.*, 2016.



J. Barbier, M. Dia, N. Macris, and F. Krzakala, "The mutual information in random linear estimation," in *Proc. Allerton Conf. Commun. Control Comput.*, pp. 625–632, 2016.



M. Bayati, M. Lelarge, and A. Montanari, "Universality in polytope phase transitions and message passing algorithms," *Ann. App. Prob.*, vol. 25, no. 2, pp. 753–822, 2015.

References II



P. Schniter, “Turbo reconstruction of structured sparse signals,” in *Proc. Conf. Inform. Science & Syst.*, (Princeton, NJ), pp. 1–6, Mar. 2010.



S. Som and P. Schniter, “Approximate message passing for recovery of sparse signals with Markov-random-field support structure.” *Internat. Conf. Mach. Learning—Workshop on Structured Sparsity: Learning and Inference*, (Bellevue, WA), July 2011.



S. Som and P. Schniter, “Compressive imaging using approximate message passing and a Markov-tree prior,” *IEEE Trans. Signal Process.*, vol. 60, pp. 3439–3448, July 2012.



D. L. Donoho, I. M. Johnstone, and A. Montanari, “Accurate prediction of phase transitions in compressed sensing via a connection to minimax denoising,” *IEEE Trans. Inform. Theory*, vol. 59, June 2013.










C. A. Metzler, A. Maleki, and R. G. Baraniuk, “BM3D-AMP: A new image recovery algorithm based on BM3D denoising,” in *Proc. IEEE Int. Conf. Image Process.*, pp. 3116–3120, 2015.









R. Berthier, A. Montanari, and P.-M. Nguyen, “State evolution for approximate message passing with non-separable functions,” *Inform. Inference*, 2019.

References III

-  C. Rush and R. Venkataramanan, "Finite-sample analysis of approximate message passing algorithms," *IEEE Trans. Inform. Theory*, vol. 64, no. 11, pp. 7264–7286, 2018.
-  S. Rangan, P. Schniter, and A. K. Fletcher, "Vector approximate message passing," *IEEE Trans. Inform. Theory*, to appear (see also arXiv:1610.03082).
-  A. M. Tulino, G. Caire, S. Verdú, and S. Shamai (Shitz), "Support recovery with sparsely sampled free random matrices," *IEEE Trans. Inform. Theory*, vol. 59, pp. 4243–4271, July 2013.
-  B. Çakmak, O. Winther, and B. H. Fleury, "S-AMP: Approximate message passing for general matrix ensembles," in *Proc. Inform. Theory Workshop*, pp. 192–196, 2014.
-  J. Vila, P. Schniter, S. Rangan, F. Krzakala, and L. Zdeborová, "Adaptive damping and mean removal for the generalized approximate message passing algorithm," in *Proc. IEEE Int. Conf. Acoust. Speech & Signal Process.*, pp. 2021–2025, 2015.
-  A. K. Fletcher, M. Sahraee-Ardakan, S. Rangan, and P. Schniter, "Expectation consistent approximate inference: Generalizations and convergence," in *Proc. IEEE Int. Symp. Inform. Thy.*, pp. 190–194, 2016.
-  M. Opper and O. Winther, "Expectation consistent approximate inference," *J. Mach. Learn. Res.*, vol. 1, pp. 2177–2204, 2005.

References IV

-  T. Minka, *A Family of Approximate Algorithms for Bayesian Inference*. PhD thesis, Dept. Comp. Sci. Eng., MIT, Cambridge, MA, Jan. 2001.
-  S. Ramani, T. Blu, and M. Unser, “Monte-Carlo SURE: A black-box optimization of regularization parameters for general denoising algorithms,” *IEEE Trans. Image Process.*, vol. 17, no. 9, pp. 1540–1554, 2008.
-  A. K. Fletcher, P. Pandit, S. Rangan, S. Sarkar, and P. Schniter, “Plug-in estimation in high-dimensional linear inverse problems: A rigorous analysis,” in *Proc. Neural Inform. Process. Syst. Conf.*, pp. 7440–7449, 2018.
-  M. Borgerding, P. Schniter, and S. Rangan, “AMP-inspired deep networks for sparse linear inverse problems,” *IEEE Trans. Signal Process.*, vol. 65, no. 15, pp. 4293–4308, 2017.
-  C. Guo and M. E. Davies, “Near optimal compressed sensing without priors: Parametric SURE approximate message passing,” *IEEE Trans. Signal Process.*, vol. 63, pp. 2130–2141, 2015.
-  U. Kamilov and H. Mansour, “Learning optimal nonlinearities for iterative thresholding algorithms,” *IEEE Signal Process. Lett.*, vol. 23, pp. 747–751, May 2016.

References V



S. Rangan, "Generalized approximate message passing for estimation with random linear mixing," in *Proc. IEEE Int. Symp. Inform. Thy.*, pp. 2168–2172, Aug. 2011.

(full version at [arXiv:1010.5141](https://arxiv.org/abs/1010.5141)).



A. Javanmard and A. Montanari, "State evolution for general approximate message passing algorithms, with applications to spatial coupling," *Inform. Inference*, vol. 2, no. 2, pp. 115–144, 2013.



P. Schniter, S. Rangan, and A. K. Fletcher, "Vector approximate message passing for the generalized linear model," in *Proc. Asilomar Conf. Signals Syst. Comput.*, pp. 1525–1529, 2016.



A. K. Fletcher, S. Rangan, and P. Schniter, "Inference in deep networks in high dimensions," in *Proc. IEEE Int. Symp. Inform. Thy.*, 2018.









U. S. Kamilov, S. Rangan, A. K. Fletcher, and M. Unser, "Approximate message passing with consistent parameter estimation and applications to sparse learning," *IEEE Trans. Inform. Theory*, vol. 60, pp. 2969–2985, May 2014.



A. K. Fletcher, M. Sahraee-Ardakan, S. Rangan, and P. Schniter, "Rigorous dynamics and consistent estimation in arbitrarily conditioned linear systems," in *Proc. Neural Inform. Process. Syst. Conf.*, pp. 2542–2551, 2017.

References VI

-  A. Mousavi, A. Maleki, and R. G. Baraniuk, “Consistent parameter estimation for LASSO and approximate message passing,” *Ann. Statist.*, vol. 45, no. 6, pp. 2427–2454, 2017.
-  A. Montanari and R. Venkataramanan, “Estimation of low-rank matrices via approximate message passing,” *arXiv:1711.01682*, 2017.
-  J. T. Parker, P. Schniter, and V. Cevher, “Bilinear generalized approximate message passing—Part I: Derivation,” *IEEE Trans. Signal Process.*, vol. 62, pp. 5839–5853, Nov. 2014.
-  J. T. Parker, P. Schniter, and V. Cevher, “Bilinear generalized approximate message passing—Part II: Applications,” *IEEE Trans. Signal Process.*, vol. 62, pp. 5854–5867, Nov. 2014.
-  J. T. Parker and P. Schniter, “Parametric bilinear generalized approximate message passing,” *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 4, pp. 795–808, 2016.
-  C. Schülke, P. Schniter, and L. Zdeborová, “Phase diagram of matrix compressed sensing,” *Physical Rev. E*, vol. 94, pp. 062136(1–16), Dec. 2016.