

# Recent Advances in Approximate Message Passing

**Phil Schniter**



THE OHIO STATE UNIVERSITY

---

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

Supported in part by NSF grants IIP-1539960, and CCF-1716388.

London Symposium on Information Theory— May 30, 2019

# Overview

- 1 Linear Regression and AMP
- 2 Vector AMP (VAMP)
- 3 VAMP for Optimization
- 4 Variational Interpretation and EM-VAMP
- 5 Plug-and-play VAMP
- 6 VAMP as a Deep Neural Network

# Outline

- 1 Linear Regression and AMP
- 2 Vector AMP (VAMP)
- 3 VAMP for Optimization
- 4 Variational Interpretation and EM-VAMP
- 5 Plug-and-play VAMP
- 6 VAMP as a Deep Neural Network

# The Linear Regression Problem

Consider the following linear regression problem:

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

Typical methodologies:

- 1 Optimization (or MAP estimation):

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \left\{ \frac{\theta_2}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + R(\mathbf{x}; \boldsymbol{\theta}_1) \right\}$$

- 2 Approximate MMSE:

$$\hat{\mathbf{x}} \approx \mathbb{E}\{\mathbf{x}|\mathbf{y}\} \quad \text{for } \mathbf{x} \sim p(\mathbf{x}; \boldsymbol{\theta}_1), \quad \mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{A}\mathbf{x}, \mathbf{I}/\theta_2)$$

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

# The AMP Methodology

- All of the aforementioned methodologies can be addressed using the **Approximate Message Passing (AMP)** framework.<sup>1</sup>
- AMP tackles these problems via **iterative denoising**.
- Each method defines the **denoiser**  $\mathbf{g}(\cdot; \gamma, \boldsymbol{\theta}_1) : \mathbb{R}^N \rightarrow \mathbb{R}^N$  differently:
  - Optimization:  $\mathbf{g}(\mathbf{r}; \gamma, \boldsymbol{\theta}_1) = \arg \min_{\mathbf{x}} \{R(\mathbf{x}; \boldsymbol{\theta}_1) + \frac{\gamma}{2} \|\mathbf{x} - \mathbf{r}\|_2^2\} \triangleq \text{“prox}_{R/\gamma}(\mathbf{r})”$
  - MMSE:  $\mathbf{g}(\mathbf{r}; \gamma, \boldsymbol{\theta}_1) = \mathbb{E} \{ \mathbf{x} \mid \mathbf{r} = \mathbf{x} + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma) \}$
  - Plug-and-play:<sup>2</sup>  $\mathbf{g}(\mathbf{r}; \gamma, \boldsymbol{\theta}_1) = \text{BM3D}(\mathbf{r}, 1/\gamma)$
  - Deep network:  $\mathbf{g}(\mathbf{r}; \gamma, \boldsymbol{\theta}_1)$  is learned from training data.

<sup>1</sup>Donoho, Maleki, Montanari'09,    <sup>2</sup>Venkatakrisnan, Bouman, Wohlberg'13

# The Original AMP Algorithm

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

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

$$\mathbf{v}^t = \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}^t + \frac{N}{M}\mathbf{v}^{t-1}\langle \mathbf{g}'(\hat{\mathbf{x}}^{t-1} + \mathbf{A}^\top \hat{\mathbf{v}}^{t-1}, \gamma^{t-1}) \rangle \quad \text{corrected residual}$$

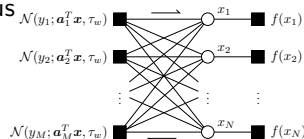
$$\hat{\mathbf{x}}^{t+1} = \mathbf{g}(\hat{\mathbf{x}}^t + \mathbf{A}^\top \mathbf{v}^t; \gamma^t) \quad \text{denoising}$$

where

$$\langle \mathbf{g}'(\mathbf{r}) \rangle \triangleq \frac{1}{N} \text{tr} \left[ \frac{\partial \mathbf{g}(\mathbf{r})}{\partial \mathbf{r}} \right] = \frac{1}{N} \sum_{j=1}^N \frac{\partial g_j(\mathbf{r})}{\partial r_j} \quad \text{“divergence.”}$$

Note:

- Can be recognized as iterative thresholding plus “Onsager correction.”
- Can be derived using Gaussian & Taylor-series approximations of belief-propagation.



# AMP's Denoising Property

## Assumption 1

- $\mathbf{A} \in \mathbb{R}^{M \times N}$  is i.i.d. sub-Gaussian
- $M, N \rightarrow \infty$  s.t.  $\frac{M}{N} \rightarrow \delta \in (0, \infty)$  ... “large-system limit”
- $[\mathbf{g}(\mathbf{r})]_j = g(r_j)$  with Lipschitz  $g(\cdot)$  ... “separable denoising”

Under Assumption 1, the elements of the denoiser's input  $\mathbf{r}^t \triangleq \hat{\mathbf{x}}^t + \mathbf{A}^\top \mathbf{v}^t$  obey<sup>34</sup>

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

- That is,  $\mathbf{r}^t$  is a Gaussian-noise corrupted version of the true signal  $\mathbf{x}_o$ .
- It is now clear why  $g(\cdot)$  is called a “denoiser.”

Furthermore, the noise variance can be consistently estimated:

$$\hat{\tau}_r^t \triangleq \frac{1}{M} \|\mathbf{v}^t\|^2 \longrightarrow \tau_r^t \quad \text{under Assumption 1.}$$

<sup>3</sup>Bayati, Montanari'11, <sup>4</sup>Bayati, Lelarge, Montanari'15

# AMP's State Evolution

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

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

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

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

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

- Then, under Assumption 1, AMP obeys the following scalar **state evolution**:

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

$$\tau_r^t = \tau_w + \frac{N}{M} \mathcal{E}^t$$

$$\mathcal{E}^{t+1} = \mathbb{E} \{ [g(X_o + \mathcal{N}(0, \tau_r^t); \gamma^t) - X_o]^2 \}$$



# Bayes Optimality of AMP

- Now suppose that Assumption 1 holds, and that

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

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

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

$$g(R; \gamma^t) = \mathbb{E} \{ X_o \mid R = X_o + \mathcal{N}(0, 1/\gamma^t) \} \quad \text{with} \quad \gamma^t = 1/\tau_r^t.$$

- Then, if the state evolution has a **unique** fixed point,  $\hat{\mathbf{x}}^t$  converges to the MMSE estimate<sup>5</sup> of  $\mathbf{x}_o$  as  $t \rightarrow \infty$ .

---

<sup>5</sup>Bayati, Montanari'11

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

## The good:

- With **large**<sup>6</sup> **i.i.d. sub-Gaussian**  $\mathbf{A}$ , AMP is rigorously characterized by a scalar **state-evolution** whose fixed points, when unique, are **Bayes optimal**.
- **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  $\mathbf{A}$ )




---

<sup>6</sup>Rush, Venkataramanan'16

# Outline

- 1 Linear Regression and AMP
- 2 Vector AMP (VAMP)
- 3 VAMP for Optimization
- 4 Variational Interpretation and EM-VAMP
- 5 Plug-and-play VAMP
- 6 VAMP as a Deep Neural Network

## Vector AMP (VAMP)

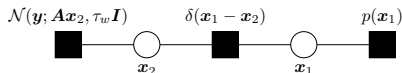


- VAMP is similar to AMP, but it supports a larger class of random matrices.
- As before, the goal is to recover  $\mathbf{x}_o$  from  $\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathcal{N}(\mathbf{0}, \tau_w \mathbf{I})$ .
- VAMP yields a precise analysis for **right-orthogonally invariant  $\mathbf{A}$** :

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

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

- Can be derived as a form of expectation propagation (EP).



# VAMP: The Algorithm

Take **SVD**  $\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s}) \mathbf{V}^\top$ , choose  $\zeta \in (0, 1]$  and Lipschitz  $\mathbf{g}_1(\cdot; \gamma_1, \boldsymbol{\theta}_1) : \mathbb{R}^N \rightarrow \mathbb{R}^N$ .

Initialize  $\mathbf{r}_1, \gamma_1$ .

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

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

$$\eta_1 \leftarrow \gamma_1 N / \text{tr} \left[ \frac{\partial \mathbf{g}_1(\mathbf{r}_1; \gamma_1, \boldsymbol{\theta}_1)}{\partial \mathbf{r}_1} \right]$$

$$\mathbf{r}_2 \leftarrow (\eta_1 \hat{\mathbf{x}}_1 - \gamma_1 \mathbf{r}_1) / (\eta_1 - \gamma_1) \quad \text{Onsager correction}$$

$$\gamma_2 \leftarrow \eta_1 - \gamma_1$$

$$\hat{\mathbf{x}}_2 \leftarrow \mathbf{g}_2(\mathbf{r}_2; \gamma_2, \boldsymbol{\theta}_2) \quad \begin{array}{l} \text{LMMSE estimate } \mathbf{x} \sim \mathcal{N}(\mathbf{r}_2, \mathbf{I}/\gamma_2) \\ \text{from } \mathbf{y} = \mathbf{A}\mathbf{x} + \mathcal{N}(\mathbf{0}, \mathbf{I}/\theta_2) \end{array}$$

$$\eta_2 \leftarrow \gamma_2 N / \text{tr} \left[ \frac{\partial \mathbf{g}_2(\mathbf{r}_2; \gamma_2, \boldsymbol{\theta}_2)}{\partial \mathbf{r}_2} \right]$$

$$\mathbf{r}_1 \leftarrow \zeta (\eta_2 \hat{\mathbf{x}}_2 - \gamma_2 \mathbf{r}_2) / (\eta_2 - \gamma_2) + (1 - \zeta) \mathbf{r}_1 \quad \text{Onsager correction}$$

$$\gamma_1 \leftarrow \zeta (\eta_2 - \gamma_2) + (1 - \zeta) \gamma_1 \quad \text{damping}$$

where  $\mathbf{g}_2(\mathbf{r}_2; \gamma_2, \boldsymbol{\theta}_2) = \mathbf{V} (\theta_2 \text{Diag}(\mathbf{s})^2 + \gamma_2 \mathbf{I})^{-1} (\theta_2 \text{Diag}(\mathbf{s}) \mathbf{U}^\top \mathbf{y} + \gamma_2 \mathbf{V}^\top \mathbf{r}_2)$   
 $\eta_2 = \frac{1}{N} \sum_{n=1}^N (\theta_2 s_n^2 + \gamma_2)^{-1}$  two mat-vec mults per iteration!

# VAMP's Denoising Property

## Assumption 2

- $\mathbf{A} \in \mathbb{R}^{M \times N}$  is right-orthogonally invariant
- $M, N \rightarrow \infty$  s.t.  $\frac{M}{N} \rightarrow \delta \in (0, \infty)$  ... “large-system limit”
- $[\mathbf{g}_1(\mathbf{r})]_j = g_1(r_j)$  with Lipschitz  $g_1(\cdot)$  ... “separable denoising”

Under Assumption 2, the elements of the denoiser's input  $\mathbf{r}_1^t$  obey<sup>7</sup>

$$r_{1,j}^t = x_{o,j} + \mathcal{N}(0, \tau_1^t)$$

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

<sup>7</sup>Rangan, Schniter, 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

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

Then under Assumption 2, the VAMP obeys the following state-evolution:

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

$$\mathcal{E}_1^t = \mathbb{E} \{ [g(X_o + \mathcal{N}(0, \tau_1^t); \gamma_1^t) - X_o]^2 \} \quad \text{MSE}$$

$$\alpha_1^t = \mathbb{E} \{ g'(X_o + \mathcal{N}(0, \tau_1^t); \gamma_1^t) \} \quad \text{divergence}$$

$$\gamma_2^t = \gamma_1^t \frac{1 - \alpha_1^t}{\alpha_1^t}, \quad \tau_2^t = \frac{1}{(1 - \alpha_1^t)^2} [\mathcal{E}_1^t - (\alpha_1^t)^2 \tau_1^t]$$

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

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

$$\gamma_1^{t+1} = \gamma_2^t \frac{1 - \alpha_2^t}{\alpha_2^t}, \quad \tau_1^{t+1} = \frac{1}{(1 - \alpha_2^t)^2} [\mathcal{E}_2^t - (\alpha_2^t)^2 \tau_2^t]$$

Note: Above assumes  $g_2(\cdot)$  uses matched noise variance  $\theta_2 = 1/\tau_w$ .

If not, there are more complicated expressions for  $\mathcal{E}_2^t$  and  $\alpha_2^t$ .

# Bayes Optimality of VAMP

- Now suppose that Assumption 2 holds, and that

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

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

- Suppose also that  $g_1(\cdot)$  is the MMSE denoiser, i.e.,

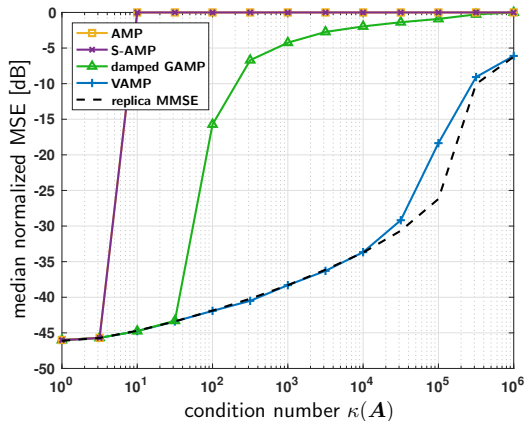
$$g_1(R_1; \gamma_1^t) = \mathbb{E} \{ X_o \mid R_1 = X_o + \mathcal{N}(0, 1/\gamma_1^t) \} \quad \text{with} \quad \gamma_1^t = 1/\tau_1^t.$$

- Then, if the state evolution has a **unique** fixed point, the MSE of  $\hat{\mathbf{x}}_1^t$  converges to the replica prediction of the MMSE as  $t \rightarrow \infty$ .
  - For right-orthogonally invariant  $\mathbf{A}$ , the replica prediction was derived by Tulino/Caire/Verdu/Shamai in 2013. It is conjectured to be correct.
  - For the special case of i.i.d. Gaussian  $\mathbf{A}$ , it was proven to be correct by Reeves/Pfister, and by Barbier/Dia/Macris/Krzakala, both in 2016.



# Experiment with MMSE Denoising

Comparison of several algorithms<sup>8</sup> with MMSE denoising.



$N = 1024$

$M/N = 0.5$

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

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

$$s_n/s_{n-1} = \phi \quad \forall n$$

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

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

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

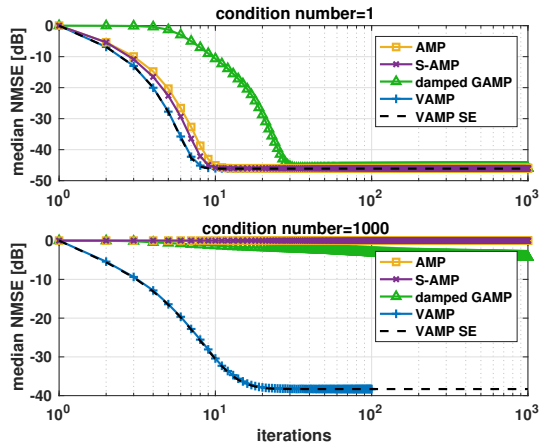
SNR = 40dB

VAMP achieves MMSE over a wide range of condition numbers.

<sup>8</sup>S-AMP: Cakmak, Fleury, Winther'14, damped GAMP: Vila, Schniter, Rangan, Krzakala, Zdeborová'15

# Experiment with MMSE Denoising

Comparison of several algorithms with priors matched to data.



$$N = 1024$$

$$M/N = 0.5$$

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

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

$$s_n/s_{n-1} = \phi \quad \forall n$$

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

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

$$\Pr\{X_0 \neq 0\} = 0.1$$

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

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

# Outline

- 1 Linear Regression and AMP
- 2 Vector AMP (VAMP)
- 3 VAMP for Optimization**
- 4 Variational Interpretation and EM-VAMP
- 5 Plug-and-play VAMP
- 6 VAMP as a Deep Neural Network

# VAMP for Optimization

- Consider the optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \left\{ \frac{\theta_2}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + R(\mathbf{x}; \boldsymbol{\theta}_1) \right\}$$

where  $R(\cdot)$  is **strictly convex**.

- If we choose the denoiser

$$\mathbf{g}_1(\mathbf{r}; \gamma, \boldsymbol{\theta}_1) = \arg \min_{\mathbf{x}} \left\{ R(\mathbf{x}; \boldsymbol{\theta}_1) + \frac{\gamma}{2} \|\mathbf{x} - \mathbf{r}\|^2 \right\} = \text{prox}_{R/\gamma}(\mathbf{r})$$

and the damping parameter

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

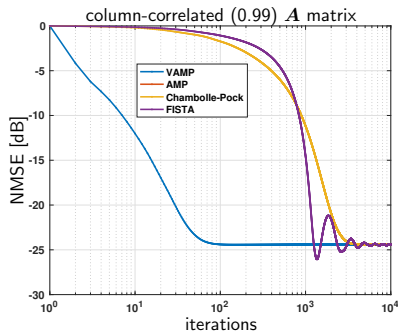
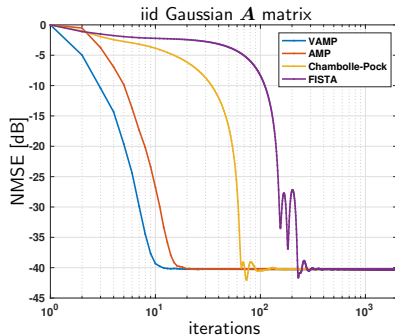
then a double-loop version of VAMP **converges<sup>9</sup> to the solution for any  $\mathbf{A}$** .

- Furthermore, if the  $\gamma_1$  and  $\gamma_2$  variables are fixed over the iterations, then VAMP reduces to the Peaceman-Rachford variant of ADMM.

---

<sup>9</sup>Fletcher, Sahraee, Rangan, Schniter'16

# Example of VAMP applied to the LASSO Problem



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

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

# Outline

- 1 Linear Regression and AMP
- 2 Vector AMP (VAMP)
- 3 VAMP for Optimization
- 4 Variational Interpretation and EM-VAMP**
- 5 Plug-and-play VAMP
- 6 VAMP as a Deep Neural Network

# Interpretation as Variational Inference

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

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x}; \boldsymbol{\theta}_1)\ell(\mathbf{x}; \theta_2)}{Z(\boldsymbol{\theta})} \quad \text{for } Z(\boldsymbol{\theta}) \triangleq \int p(\mathbf{x}; \boldsymbol{\theta}_1)\ell(\mathbf{x}; \theta_2) d\mathbf{x},$$

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

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

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \arg \min_b D(b(\mathbf{x}) \| p(\mathbf{x}|\mathbf{y})) \quad \text{where } D(\cdot \| \cdot) \text{ is KL divergence} \\ &= \arg \min_b \underbrace{D(b(\mathbf{x}) \| p(\mathbf{x}; \boldsymbol{\theta}_1)) + D(b(\mathbf{x}) \| \ell(\mathbf{x}; \theta_2)) + H(b(\mathbf{x}))}_{\text{Gibbs free energy}} \\ &= \arg \min_{b_1, b_2, q} \underbrace{D(b_1(\mathbf{x}) \| p(\mathbf{x}; \boldsymbol{\theta}_1)) + D(b_2(\mathbf{x}) \| \ell(\mathbf{x}; \theta_2)) + H(q(\mathbf{x}))}_{\triangleq J_{\text{Gibbs}}(b_1, b_2, q; \boldsymbol{\theta})} \\ &\quad \text{s.t. } b_1 = b_2 = q, \end{aligned}$$

but the density constraint keeps the problem difficult.

# Expectation Consistent Approximation

- In **expectation-consistent approximation (EC)**<sup>10</sup>, the density constraint is relaxed to moment-matching constraints:

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

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

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

$$\begin{aligned} b_1(\mathbf{x}) &\propto p(\mathbf{x}; \boldsymbol{\theta}_1) \mathcal{N}(\mathbf{x}; \mathbf{r}_1, \mathbf{I}/\gamma_1) \\ b_2(\mathbf{x}) &\propto \ell(\mathbf{x}; \theta_2) \mathcal{N}(\mathbf{x}; \mathbf{r}_2, \mathbf{I}/\gamma_2) \\ q(\mathbf{x}) &= \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \mathbf{I}/\eta) \end{aligned} \quad \text{s.t. } \begin{cases} \mathbb{E}\{\mathbf{x}|b_1\} = \mathbb{E}\{\mathbf{x}|b_2\} = \hat{\mathbf{x}} \\ \text{tr}(\text{Cov}\{\mathbf{x}|b_1\}) = \text{tr}(\text{Cov}\{\mathbf{x}|b_2\}) = N/\eta, \end{cases}$$

- VAMP iteratively solves for the quantities  $\mathbf{r}_1, \gamma_1, \mathbf{r}_2, \gamma_2, \hat{\mathbf{x}}, \eta$  above.
  - In this setting, VAMP is simply an instance of **expectation propagation (EP)**.

<sup>10</sup>Opper, Winther'04,

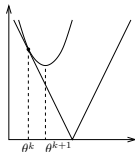


# Expectation Maximization

- What if the hyperparameters  $\theta$  of the prior & likelihood are unknown?
- The **EM algorithm**<sup>11</sup> is majorization-minimization approach to **ML estimation** that iteratively minimizes a tight upper bound on  $-\ln p(\mathbf{y}|\theta)$ :

$$\hat{\theta}^{k+1} = \arg \min_{\theta} \left\{ -\ln p(\mathbf{y}|\theta) + \underbrace{D(b^k(\mathbf{x}) \| p(\mathbf{x}|\mathbf{y}; \theta))}_{\geq 0} \right\}$$

with  $b^k(\mathbf{x}) = p(\mathbf{x}|\mathbf{y}; \hat{\theta}^k)$



- EM can also be written in terms of the **Gibbs free energy**:<sup>12</sup>

$$\hat{\theta}^{k+1} = \arg \min_{\theta} \underbrace{D(b^k(\mathbf{x}) \| p(\mathbf{x}; \theta_1)) + D(b^k(\mathbf{x}) \| \ell(\mathbf{x}; \theta_2)) + H(b^k(\mathbf{x}))}_{J_{\text{Gibbs}}(b^k, b^k, b^k; \theta)}$$

- Thus, we can **interleave EM and VAMP** to solve

$$\min_{\theta} \min_{b_1, b_2, q} J_{\text{Gibbs}}(b_1, b_2, q; \theta) \text{ s.t. } \begin{cases} \mathbb{E}\{\mathbf{x}|b_1\} = \mathbb{E}\{\mathbf{x}|b_2\} = \mathbb{E}\{\mathbf{x}|q\} \\ \text{tr}[\text{Cov}\{\mathbf{x}|b_1\}] = \text{tr}[\text{Cov}\{\mathbf{x}|b_2\}] = \text{tr}[\text{Cov}\{\mathbf{x}|q\}]. \end{cases}$$

<sup>11</sup>Dempster, Laird, Rubin'77,    <sup>12</sup>Neal, Hinton'98

# The EM-VAMP Algorithm

Input conditional-mean  $\mathbf{g}_1(\cdot)$  and  $\mathbf{g}_2(\cdot)$ , and initialize  $\mathbf{r}_1, \gamma_1, \hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_2$ .

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

$$\hat{\mathbf{x}}_1 \leftarrow \mathbf{g}_1(\mathbf{r}_1; \gamma_1, \hat{\boldsymbol{\theta}}_1) \quad \text{MMSE estimation}$$

$$\eta_1 \leftarrow \gamma_1 N / \text{tr} \left[ \partial \mathbf{g}_1(\mathbf{r}_1; \gamma_1, \hat{\boldsymbol{\theta}}_1) / \partial \mathbf{r}_1 \right]$$

$$\mathbf{r}_2 \leftarrow (\eta_1 \hat{\mathbf{x}}_1 - \gamma_1 \mathbf{r}_1) / (\eta_1 - \gamma_1)$$

$$\gamma_2 \leftarrow \eta_1 - \gamma_1$$

$$\hat{\boldsymbol{\theta}}_2 \leftarrow \arg \max_{\boldsymbol{\theta}_2} \mathbb{E} \{ \ln \ell(\mathbf{x}; \boldsymbol{\theta}_2) \mid \mathbf{r}_2; \gamma_2, \hat{\boldsymbol{\theta}}_2 \} \quad \text{EM update}$$

---


$$\hat{\mathbf{x}}_2 \leftarrow \mathbf{g}_2(\mathbf{r}_2; \gamma_2, \hat{\boldsymbol{\theta}}_2) \quad \text{LMMSE estimation}$$

$$\eta_2 \leftarrow \gamma_2 N / \text{tr} \left[ \partial \mathbf{g}_2(\mathbf{r}_2; \gamma_2, \hat{\boldsymbol{\theta}}_2) / \partial \mathbf{r}_2 \right]$$

$$\mathbf{r}_1 \leftarrow \zeta(\eta_2 \hat{\mathbf{x}}_2 - \gamma_2 \mathbf{r}_2) / (\eta_2 - \gamma_2) + (1 - \zeta) \mathbf{r}_1$$

$$\gamma_1 \leftarrow \zeta(\eta_2 - \gamma_2) + (1 - \zeta) \gamma_1$$

$$\hat{\boldsymbol{\theta}}_1 \leftarrow \arg \max_{\boldsymbol{\theta}_1} \mathbb{E} \{ \ln p(\mathbf{x}; \boldsymbol{\theta}_1) \mid \mathbf{r}_1; \gamma_1, \hat{\boldsymbol{\theta}}_1 \} \quad \text{EM update}$$

# State Evolution and Consistency

- EM-VAMP has a rigorous state-evolution<sup>13</sup> when the prior is i.i.d. and  $\mathbf{A}$  is large and right-orthogonally invariant.
- Furthermore, a variant known as “adaptive VAMP” can be shown to yield consistent parameter estimates with an i.i.d. prior in the exponential-family or with finite-cardinality  $\theta_1$ .<sup>13</sup>
- Essentially, adaptive VAMP replaces the EM update

$$\hat{\theta}_1 \leftarrow \arg \max_{\theta_1} \mathbb{E}\{\ln p(\mathbf{x}; \theta_1) \mid \mathbf{r}_1, \gamma_1, \hat{\theta}_1\}$$

with

$$(\hat{\theta}_1, \hat{\gamma}_1) \leftarrow \arg \max_{(\theta_1, \gamma_1)} \mathbb{E}\{\ln p(\mathbf{x}; \theta_1) \mid \mathbf{r}_1, \gamma_1, \hat{\theta}_1\},$$

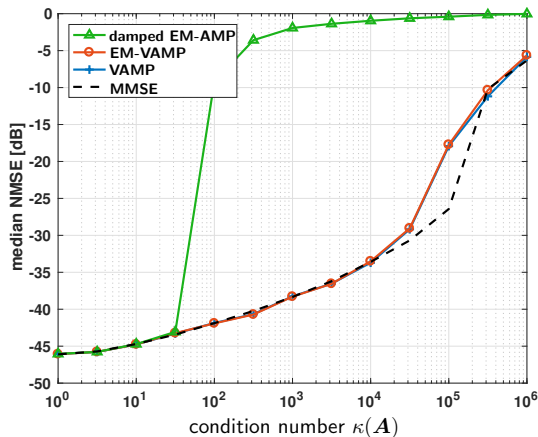
which re-estimates the precision  $\gamma_1$ . (And similar for  $\theta_2, \gamma_2$ .)

---

<sup>13</sup>Fletcher, Rangan, Schniter'17

# Experiment with Unknown Hyperparameters $\theta$

Learning both noise precision  $\theta_2$  and BG mean/variance/sparsity  $\theta_1$ :



$N = 1024$

$M/N = 0.5$

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

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

$$s_n/s_{n-1} = \phi \quad \forall n$$

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

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

$$\Pr\{X_0 \neq 0\} = 0.1$$

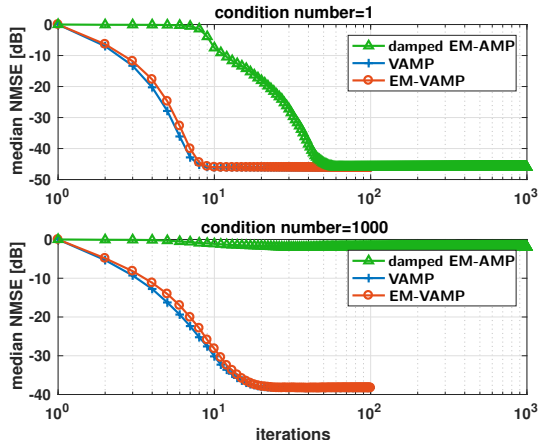
SNR = 40dB

EM-VAMP achieves oracle performance at all condition numbers!<sup>14</sup>

<sup>14</sup>EM-AMP proposed in Vila, Schniter'11 and Krzakala, Mézard, Sausset, Sun, Zdeborová'12

# Experiment with Unknown Hyperparameters $\theta$

Learning both noise precision  $\theta_2$  and BG mean/variance/sparsity  $\theta_1$ :



$$N = 1024$$

$$M/N = 0.5$$

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

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

$$s_n/s_{n-1} = \phi \quad \forall n$$

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

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

$$\Pr\{X_0 \neq 0\} = 0.1$$

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

EM-VAMP nearly as fast as VAMP and much faster than damped EM-GAMP.

# Outline

- 1 Linear Regression and AMP
- 2 Vector AMP (VAMP)
- 3 VAMP for Optimization
- 4 Variational Interpretation and EM-VAMP
- 5 Plug-and-play VAMP
- 6 VAMP as a Deep Neural Network

# Plug-and-play VAMP

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

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

- For certain signal classes (e.g., images), very sophisticated *non-scalar* denoising procedures have been developed (e.g., [BM3D](#), [DnCNN](#)).
- Such denoising procedures can be “[plugged into](#)” signal recovery algorithms like ADMM, AMP<sup>15</sup>, VAMP. Divergence can be approximated via

$$\frac{1}{N} \text{tr} \left[ \frac{\partial \mathbf{g}_1}{\partial \mathbf{r}_1} \right] \approx \frac{1}{K} \sum_{k=1}^K \frac{\mathbf{p}_k^T [\mathbf{g}_1(\mathbf{r} + \epsilon \mathbf{p}_k, \gamma_1) - \mathbf{g}_1(\mathbf{r}, \gamma_1)]}{N\epsilon}$$

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

- [Rigorous state-evolutions](#) established for plug-and-play AMP<sup>16</sup> and VAMP.<sup>17</sup>

<sup>15</sup>Metzler, Maleki, Baraniuk'14, <sup>16</sup>Berthier, Montanari, Nguyen'17, <sup>17</sup>Fletcher, Rangan, Sarkar, Schniter'18

# Bilinear estimation via Lifting

- As we now describe, non-scalar denoising facilitates **bilinear** recovery.
- Say the goal is to recover  $\mathbf{b} = [b_1, \dots, b_L]^\top$  and  $\mathbf{c}$  from measurements

$$\mathbf{y} = \left( \sum_{l=1}^L b_l \Phi_l \right) \mathbf{c} + \mathbf{w}$$

where  $\{\Phi_l\}$  are known. This arises in calibration problems.

- We can “lift”<sup>18</sup> this bilinear problem to the linear problem

$$\mathbf{y} = \underbrace{[\Phi_1 \quad \Phi_2 \quad \cdots \quad \Phi_L]}_A \underbrace{\text{vec}(\mathbf{c}\mathbf{b}^\top)}_{\mathbf{x}} + \mathbf{w}$$

and apply VAMP with an appropriate denoiser.

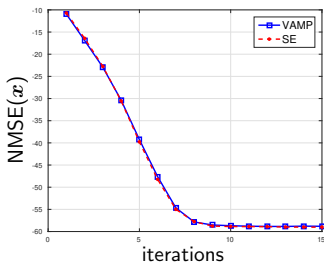
<sup>18</sup>Candes,Strohmer,Voroninski’13, Ahmed,Recht,Romberg’14



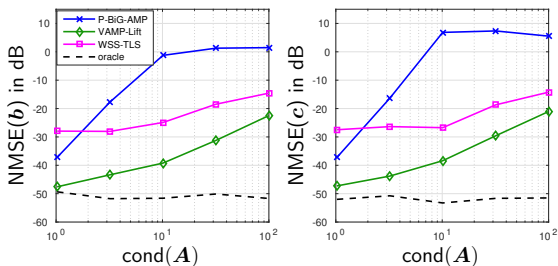
# Experiment: Compressed Sensing with Matrix Uncertainty

Goal: Recover<sup>19</sup>  $\mathbf{b}$  and sparse  $\mathbf{c}$  from  $\mathbf{y} = \left(\sum_{l=1}^L b_l \Phi_l\right) \mathbf{c} + \mathbf{w} = \mathbf{A} \mathbf{x} + \mathbf{w}$ .

State Evolution:



NMSE versus condition number of  $\mathbf{A}$ :



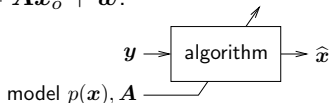
<sup>19</sup>WSS-TLS is from Zhu, Leus, Giannakis'11, P-BiG-AMP is from Parker, Schniter'16.

# Outline

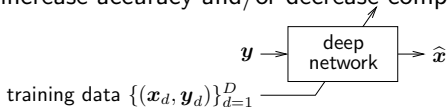
- 1 Linear Regression and AMP
- 2 Vector AMP (VAMP)
- 3 VAMP for Optimization
- 4 Variational Interpretation and EM-VAMP
- 5 Plug-and-play VAMP
- 6 VAMP as a Deep Neural Network

# Deep learning for sparse reconstruction

- Until now we've focused on **designing algorithms** to recover  $x_o \sim p(x)$  from measurements  $y = Ax_o + w$ .



- What about **training deep networks** to predict  $x_o$  from  $y$ ?  
Can we increase accuracy and/or decrease computation?



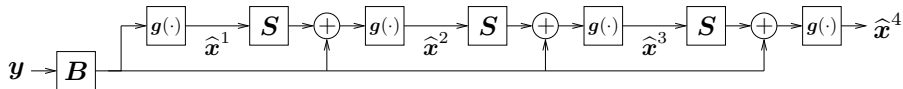
- Are there **connections** between these approaches?

# Unfolding Algorithms into Networks

Consider, e.g., the classical sparse-reconstruction algorithm, [ISTA](#).<sup>20</sup>

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

Gregor & LeCun<sup>21</sup> 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).<sup>22</sup>

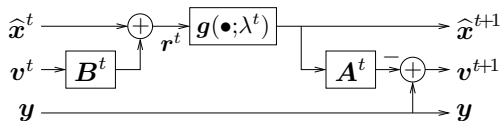
<sup>20</sup>Daubechies,Defrise,DeMol'04.

<sup>21</sup>Gregor,LeCun'10.

<sup>22</sup>Borgerding,Schniter'16.

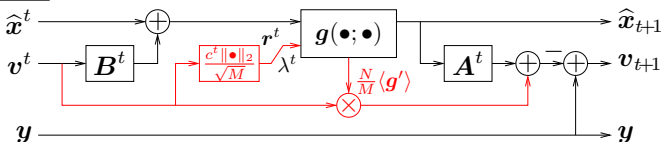
# Onsager-Corrected Deep Networks

$t^{\text{th}}$  LISTA layer:



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

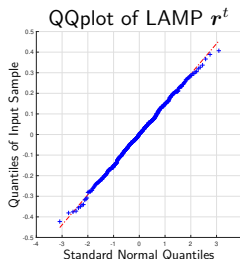
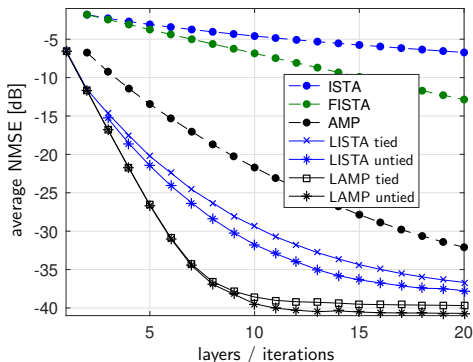
$t^{\text{th}}$  LAMP layer:



Onsager correction now aims to decouple errors across layers.

# LAMP performance with soft-threshold denoising

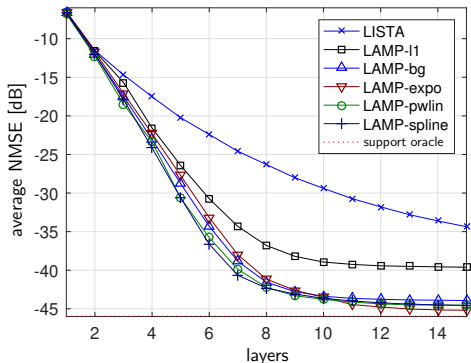
LISTA beats AMP, FISTA, ISTA in convergence speed and asymptotic MSE.  
LAMP beats LISTA



# 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 kernel<sup>23</sup>
- Piecewise Linear<sup>23</sup>
- Spline<sup>24</sup>

**Big improvement!**

<sup>23</sup>Guo,Davies'15. <sup>24</sup>Kamilov,Mansour'16.

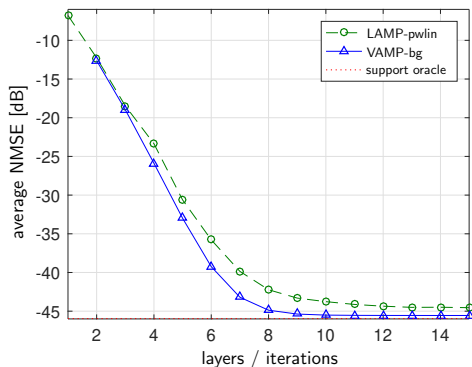
LAMP



versus VAMP



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



*VAMP wins!*

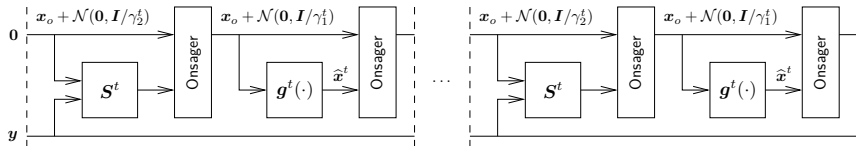
So what about “learned VAMP”?



# Learned VAMP



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



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

*Theory explains the deep network!*

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

# Conclusions

- VAMP is a computationally efficient algorithm for **linear** regression.
- For inference under large, right orthogonally-invariant  $\mathbf{A}$ , VAMP has a **rigorous state evolution** whose fixed-points, when unique, match the replica prediction of the **MMSE**.
- For **convex** optimization problems, VAMP is **provably convergent** for any  $\mathbf{A}$ .
- VAMP can be **combined with EM** to handle priors/likelihood with unknown parameters, again with a rigorous state evolution.
- VAMP supports **nonseparable (i.e., “plug-in”) denoisers**, with a rigorous state evolution.
- Can unfold VAMP into an **interpretable deep network**.
- Not discussed: **GLMs**, **multilayer** approaches, **bilinear** approaches.