

# The VAMPire, Part I

**Phil Schniter**



THE OHIO STATE UNIVERSITY



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

Supported in part by NSF grants IIP-1539960 and CCF-1527162.

Les Houches Workshop on  
Statistical Physics, Learning, Inference, and Networks—Feb 2017

# Overview

- 1 Linear Regression
- 2 EC, EP, and VAMP
- 3 EM-VAMP and Adaptive VAMP
- 4 Damped VAMP and Connections to ADMM
- 5 Plug-and-play VAMP & Whitening
- 6 VAMP as a Deep Neural Network

# Outline

- 1 Linear Regression
- 2 EC, EP, and VAMP
- 3 EM-VAMP and Adaptive VAMP
- 4 Damped VAMP and Connections to ADMM
- 5 Plug-and-play VAMP & Whitening
- 6 VAMP as a Deep Neural Network

## The Linear Regression Problem

In this talk, we consider the following linear regression problem:

- Observations:

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

- Prior:

$$p(\mathbf{x}; \boldsymbol{\theta}_1) \quad \text{with deterministic unknown parameters } \boldsymbol{\theta}_1.$$

- Likelihood function:

$$\ell(\mathbf{x}; \theta_2) \triangleq \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \mathbf{I}/\theta_2) \quad \text{with deterministic unknown precision } \theta_2.$$

**Goal:** infer  $\mathbf{x}$  and estimate  $\boldsymbol{\theta} \triangleq [\theta_1, \theta_2]$  (as needed).

(Sundeeep will discuss more complicated priors and likelihoods!)

# Outline

- 1 Linear Regression
- 2 EC, EP, and VAMP
- 3 EM-VAMP and Adaptive VAMP
- 4 Damped VAMP and Connections to ADMM
- 5 Plug-and-play VAMP & Whitening
- 6 VAMP as a Deep Neural Network

## Variational Inference

- For the moment, let's suppose that  $\theta$  is known.
- We would like to compute the **posterior density**

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x}; \theta_1)\ell(\mathbf{x}; \theta_2)}{Z(\theta)} \quad \text{for } Z(\theta) \triangleq \int p(\mathbf{x}; \theta_1)\ell(\mathbf{x}; \theta_2) d\mathbf{x},$$

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

- We can avoid computing  $Z(\theta)$  through **variational optimization**:

$$p(\mathbf{x}|\mathbf{y}) = \arg \min_b D(b(\mathbf{x})\|p(\mathbf{x}|\mathbf{y})) \quad \text{where } D(\cdot\|\cdot) \text{ is KL divergence}$$

$$= \arg \min_b \underbrace{D(b(\mathbf{x})\|p(\mathbf{x}; \theta_1)) + D(b(\mathbf{x})\|\ell(\mathbf{x}; \theta_2)) + H(b(\mathbf{x}))}_{\text{Gibbs free energy}}$$

$$= \arg \min_{b_1, b_2} \max_q \underbrace{D(b_1(\mathbf{x})\|p(\mathbf{x}; \theta_1)) + D(b_2(\mathbf{x})\|\ell(\mathbf{x}; \theta_2)) + H(q(\mathbf{x}))}_{\triangleq J(b_1, b_2, q; \theta)}$$

$$\text{such that } b_1 = b_2 = q, \quad \triangleq J(b_1, b_2, q; \theta)$$

but the density constraint keeps the problem difficult.

## Expectation Consistent Approximation

- In **Expectation consistent approximation (EC)**<sup>1</sup>, the density constraint is relaxed to moment-matching constraints:

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

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

for any affine linear diagonalization operator  $\text{diag} : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^N$ .

- The **stationary points** of EC are

$$\begin{aligned} b_1(\mathbf{x}) &\propto p(\mathbf{x}; \boldsymbol{\theta}_1) \mathcal{N}(\mathbf{x}; \mathbf{r}_1, \mathbf{1} \otimes \boldsymbol{\gamma}_1) \\ b_2(\mathbf{x}) &\propto \ell(\mathbf{x}; \theta_2) \mathcal{N}(\mathbf{x}; \mathbf{r}_2, \mathbf{1} \otimes \boldsymbol{\gamma}_2) \\ q(\mathbf{x}) &= \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \mathbf{1} \otimes \boldsymbol{\eta}) \end{aligned}$$

$$\text{s.t. } \begin{cases} \mathbb{E}\{\mathbf{x}|b_1\} = \mathbb{E}\{\mathbf{x}|b_2\} = \hat{\mathbf{x}} \\ \text{diag}(\text{Cov}\{\mathbf{x}|b_1\}) = \text{diag}(\text{Cov}\{\mathbf{x}|b_2\}) = \boldsymbol{\eta}. \end{cases}$$

which match the replica prediction of MMSE.<sup>2</sup>

<sup>1</sup>Opper, Winther'04, <sup>2</sup>Kabashima, Vehkaperä'14

## An EP Algorithm

An expectation-propagation (EP)<sup>3</sup> algorithm whose fixed points match EC is

Initialize $\mathbf{r}_1, \gamma_1$ .	
For $k = 1, 2, 3, \dots$	
$\boldsymbol{\eta}_1 \leftarrow \mathbf{1} \oslash \text{diag} [\text{Cov} \{ \mathbf{x} \mid \mathbf{r}_1 = \mathbf{x} + \mathcal{N}(\mathbf{0}, \text{Diag}(\boldsymbol{\gamma}_1)^{-1}) \}]$	posterior precision
$\hat{\mathbf{x}}_1 \leftarrow \mathbb{E} \{ \mathbf{x} \mid \mathbf{r}_1 = \mathbf{x} + \mathcal{N}(\mathbf{0}, \text{Diag}(\boldsymbol{\gamma}_1)^{-1}) \}$	denoising
$\boldsymbol{\gamma}_2 \leftarrow \boldsymbol{\eta}_1 - \boldsymbol{\gamma}_1$	extrinsic precision
$\mathbf{r}_2 \leftarrow (\hat{\mathbf{x}}_1 \odot \boldsymbol{\eta}_1 - \mathbf{r}_1 \odot \boldsymbol{\gamma}_1) \oslash \boldsymbol{\gamma}_2$	Onsager correction
<hr/>	
$\boldsymbol{\eta}_2 \leftarrow \mathbf{1} \oslash \text{diag} [(\theta_2 \mathbf{A}^\top \mathbf{A} + \text{Diag}(\boldsymbol{\gamma}_2))^{-1}]$	posterior precision
$\hat{\mathbf{x}}_2 \leftarrow (\theta_2 \mathbf{A}^\top \mathbf{A} + \text{Diag}(\boldsymbol{\gamma}_2))^{-1} (\theta_2 \mathbf{A}^\top \mathbf{y} + \boldsymbol{\gamma}_2 \odot \mathbf{r}_2)$	LMMSE estimation
$\boldsymbol{\gamma}_1 \leftarrow \boldsymbol{\eta}_2 - \boldsymbol{\gamma}_2$	extrinsic precision
$\mathbf{r}_1 \leftarrow (\hat{\mathbf{x}}_2 \odot \boldsymbol{\eta}_2 - \mathbf{r}_2 \odot \boldsymbol{\gamma}_2) \oslash \boldsymbol{\gamma}_1$	Onsager correction

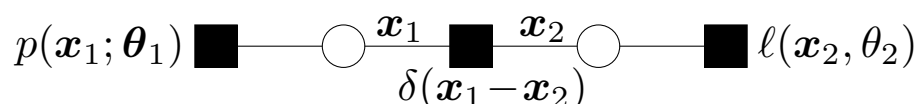
**Questions:** 1) Does it converge? Can we predict its trajectory?  
2) Can we avoid the matrix inverse?

<sup>3</sup>Minka'01



## Vector AMP (VAMP)

- Consider EP with the “uniform diagonalization”  $\text{diag}(\mathbf{C}) = \mathbf{1} \text{tr}(\mathbf{C})/N$ . By precomputing the SVD  $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ , we can<sup>4</sup>
  - 1) avoid matrix inversion and
  - 2) derive a rigorous state evolution. (Details in Sundeep’s talk!)
- We call the algorithm “vector AMP” because, like the AMP algorithm,<sup>56</sup> it has a rigorous state-evolution and can be derived using an approximation of message passing, now on a vector-valued factor graph:



- Importantly, VAMP is provably convergent under either
  - 1) iid prior  $p(\mathbf{x}; \boldsymbol{\theta}_1)$  and large, right-rotationally invariant  $\mathbf{A}$ ,
  - 2) strictly log-concave prior  $p(\mathbf{x}; \boldsymbol{\theta}_1)$  and arbitrary  $\mathbf{A}$ .

<sup>4</sup>Rangan, Schniter, Fletcher’16,    <sup>5</sup>Donoho, Maleki, Montanari’09,    <sup>6</sup>Bayati, Montanari’10

# The VAMP Algorithm

For any Lipschitz function  $\mathbf{g}_1(\cdot; \gamma_1, \boldsymbol{\theta}_1) : \mathbb{R}^N \rightarrow \mathbb{R}^N \dots$

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

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

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

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

$$\gamma_2 \leftarrow \eta_1 - \gamma_1 \quad \text{extrinsic precision}$$

$$\mathbf{r}_2 \leftarrow (\hat{\mathbf{x}}_1 \eta_1 - \mathbf{r}_1 \gamma_1) / \gamma_2 \quad \text{Onsager correction}$$

$$\eta_2 \leftarrow \frac{1}{N} \sum_{n=1}^N \frac{1}{\theta_2 s_n^2 + \gamma_2} \quad \text{posterior precision}$$

$$\hat{\mathbf{x}}_2 \leftarrow \mathbf{V} (\theta_2 \text{Diag}(\mathbf{s})^2 + \gamma_2 \mathbf{I})^{-1} (\theta_2 \text{Diag}(\mathbf{s}) \mathbf{U}^T \mathbf{y} + \gamma_2 \mathbf{V}^T \mathbf{r}_2) \quad \text{LMMSE estimation}$$

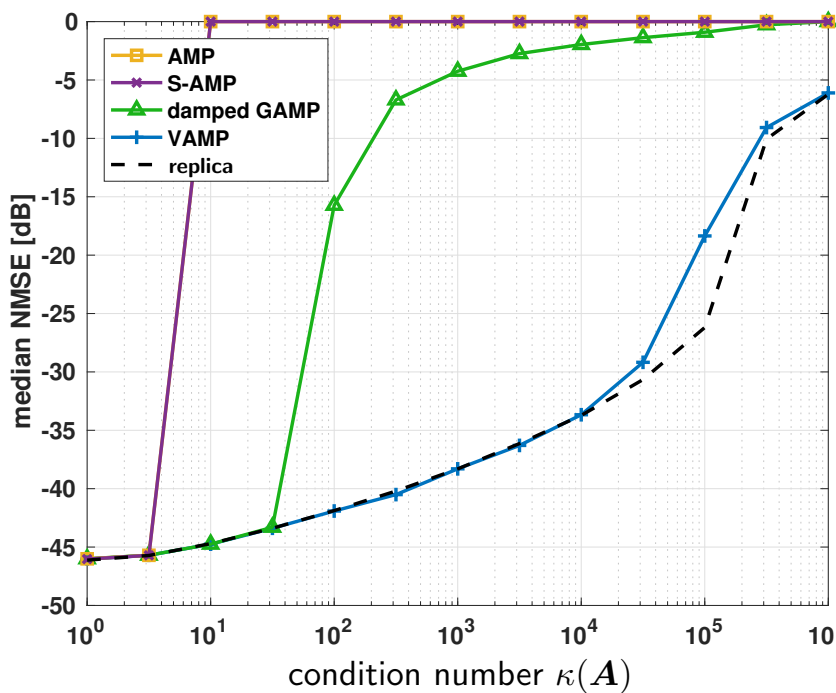
$$\gamma_1 \leftarrow \eta_2 - \gamma_2 \quad \text{extrinsic precision}$$

$$\mathbf{r}_1 \leftarrow (\hat{\mathbf{x}}_2 \eta_2 - \mathbf{r}_2 \gamma_2) / \gamma_1 \quad \text{Onsager correction}$$

Complexity is dominated by two matrix-vector multiplies per iteration.

## Experiment with Matched Priors

Comparison of several algorithms<sup>7</sup> with priors matched to data.



$N = 1024$   
 $M/N = 0.5$

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

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

SNR = 40dB

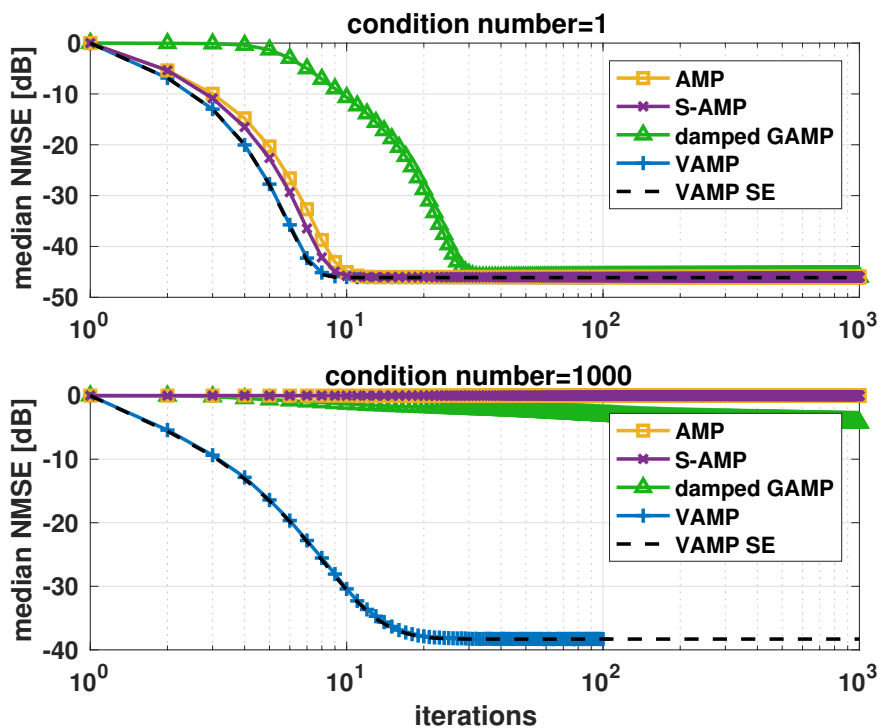
VAMP follows replica prediction<sup>8</sup> over a wide range of condition numbers.

<sup>7</sup>S-AMP: Cakmak, Fleury, Winther'14, AD-GAMP: Vila, Schniter, Rangan, Krzakala, Zdeborová'15

<sup>8</sup>Tulino, Caire, Verdú, Shamai'13

## Experiment with Matched Priors

Comparison of several algorithms with priors matched to data.



$$N = 1024$$

$$M/N = 0.5$$

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

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

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

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

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

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

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

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

# Outline

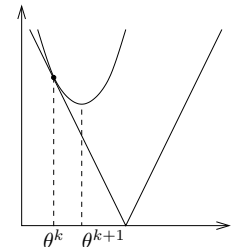
- 1 Linear Regression
- 2 EC, EP, and VAMP
- 3 EM-VAMP and Adaptive VAMP**
- 4 Damped VAMP and Connections to ADMM
- 5 Plug-and-play VAMP & Whitening
- 6 VAMP as a Deep Neural Network

## Expectation Maximization

- We now return to the problem of estimating  $\theta$ .
- The **EM algorithm**<sup>9</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  $b^k(\mathbf{x}) = p(\mathbf{x}|\mathbf{y}; \hat{\theta}^k)$



- Can also write EM in terms of the Gibbs free energy:<sup>10</sup>

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

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

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

<sup>9</sup>Dempster, Laird, Rubin'77,    <sup>10</sup>Neal, Hinton'98

# The EM-VAMP Algorithm

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

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

$$\eta_1 \leftarrow \gamma_1 / \langle \mathbf{g}'_1(\mathbf{r}_1; \gamma_1, \hat{\boldsymbol{\theta}}_1) \rangle \quad \text{posterior precision}$$

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

$$\gamma_2 \leftarrow \eta_1 - \gamma_1 \quad \text{extrinsic precision}$$

$$\mathbf{r}_2 \leftarrow (\hat{\mathbf{x}}_1 \eta_1 - \mathbf{r}_1 \gamma_1) / \gamma_2 \quad \text{Onsager correction}$$

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

$$\eta_2 \leftarrow \gamma_2 / \langle \mathbf{g}'_2(\mathbf{r}_2; \gamma_2, \hat{\boldsymbol{\theta}}_2) \rangle \quad \text{posterior precision}$$

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

$$\gamma_1 \leftarrow \eta_2 - \gamma_2 \quad \text{extrinsic precision}$$

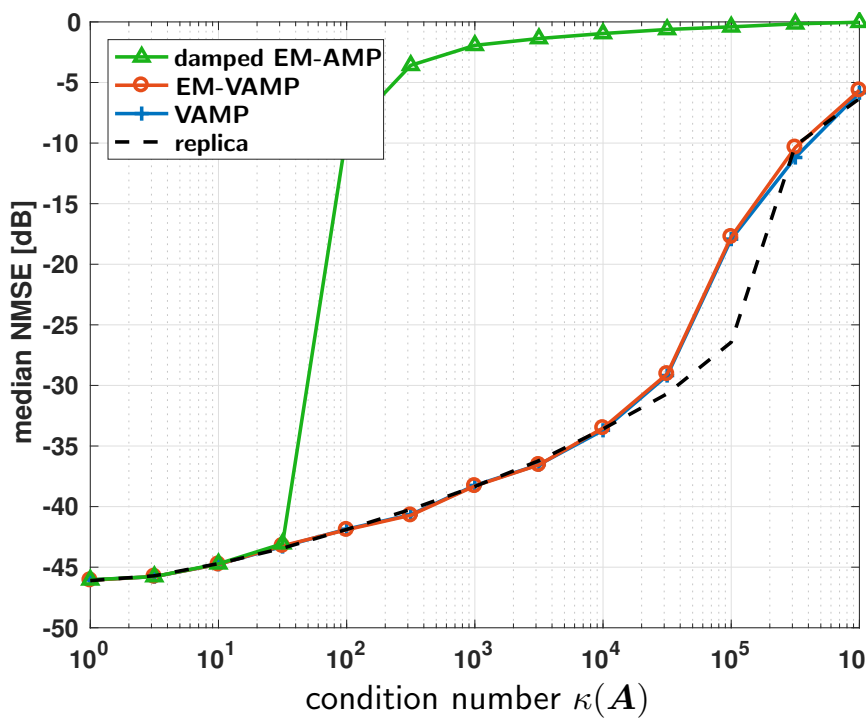
$$\mathbf{r}_1 \leftarrow (\hat{\mathbf{x}}_2 \eta_2 - \mathbf{r}_2 \gamma_2) / \gamma_1 \quad \text{Onsager correction}$$

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

Experiments suggest it helps to update  $\hat{\boldsymbol{\theta}}_2$  several times per VAMP iteration.

## Experiment with Learned $\theta$

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



$N = 1024$   
 $M/N = 0.5$

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

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

SNR = 40dB

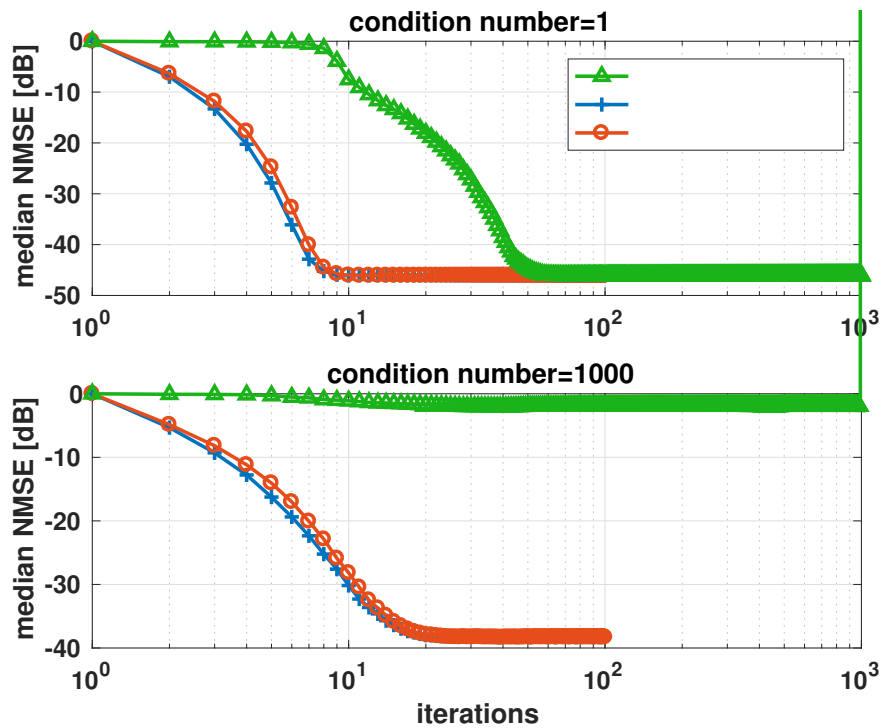
EM-VAMP achieves oracle performance at all condition numbers!<sup>11</sup>

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



## Experiment with Learned $\theta$

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



## State Evolution and Consistency

- Like VAMP, EM-VAMP has a rigorous state-evolution when the prior is i.i.d. and  $\mathbf{A}$  is large and right-rotationally invariant.<sup>12</sup>
- 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>12</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

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

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

*More details in Sundeep's talk!*

---

<sup>12</sup>Fletcher, Rangan, Schniter'17

# Outline

- 1 Linear Regression
- 2 EC, EP, and VAMP
- 3 EM-VAMP and Adaptive VAMP
- 4 Damped VAMP and Connections to ADMM**
- 5 Plug-and-play VAMP & Whitening
- 6 VAMP as a Deep Neural Network

## Damped VAMP

- Practical applications of VAMP may involve matrices  $\mathbf{A}$  that are *not* large and right-rotationally invariant and/or signals that are *not* i.i.d.!
- Can prove that that a **damped** version of VAMP **converges** when

$$\exists c_1, c_2 > 0 \text{ s.t. } \frac{\gamma}{\gamma + c_1} \mathbf{I} \leq \frac{\partial \mathbf{g}_1}{\partial \mathbf{r}}(\mathbf{r}, \gamma) \leq \frac{\gamma}{\gamma + c_2} \mathbf{I},$$

as occurs with MAP inference under strictly log-concave prior  $p(\mathbf{x}; \boldsymbol{\theta}_1)$ .

- To do the damping, we simply replace the  $\mathbf{r}_1$  update with

$$\mathbf{r}_1 \leftarrow \zeta [(\hat{\mathbf{x}}_2 \eta_2 - \mathbf{r}_2 \gamma_2) / \gamma_1] + (1 - \zeta) \mathbf{r}_1 \text{ for some } \zeta \in (0, 1].$$

- For convergence, it suffices that  $\zeta \leq \frac{2 \min\{\gamma_1, \gamma_2\}}{\gamma_1 + \gamma_2}$ . Thus
  - the damping factor can be adapted using  $\gamma_1, \gamma_2$ ,
  - no damping (i.e.,  $\zeta = 1$ ) suffices when  $\gamma_1 = \gamma_2$ .

## Connections to ADMM

- Consider the optimization problem

$$\arg \min_{\mathbf{x}} f_1(\mathbf{x}) + f_2(\mathbf{x}) \quad \text{with, e.g.,} \quad \begin{cases} f_1(\mathbf{x}) = -\log p(\mathbf{x}; \boldsymbol{\theta}_1) \\ f_2(\mathbf{x}) = \frac{\theta_2}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 \end{cases}$$

and define the augmented Lagrangian

$$L_\gamma(\mathbf{x}_1, \mathbf{x}_2, \mathbf{s}) = f_1(\mathbf{x}_1) + f_2(\mathbf{x}_2) + \mathbf{s}^\top (\mathbf{x}_1 - \mathbf{x}_2) + \frac{\gamma}{2} \|\mathbf{x}_1 - \mathbf{x}_2\|^2.$$

- An ADMM-like algorithm (known as [Peaceman-Rachford splitting](#)) is

$$\hat{\mathbf{x}}_1 \leftarrow \arg \min_{\mathbf{x}_1} L_\gamma(\mathbf{x}_1, \hat{\mathbf{x}}_2, \mathbf{s})$$

$$\mathbf{s} \leftarrow \mathbf{s} + \gamma(\hat{\mathbf{x}}_1 - \hat{\mathbf{x}}_2)$$

$$\hat{\mathbf{x}}_2 \leftarrow \arg \min_{\mathbf{x}_2} L_\gamma(\hat{\mathbf{x}}_1, \mathbf{x}_2, \mathbf{s})$$

$$\mathbf{s} \leftarrow \mathbf{s} + \gamma(\hat{\mathbf{x}}_1 - \hat{\mathbf{x}}_2)$$

- PRS is guaranteed to converge for strictly convex  $f_1(\cdot)$  and  $f_2(\cdot)$ .

## Connections to ADMM

- Now consider VAMP applied to the same MAP optimization problem, but with  $\gamma_1 = \gamma_2 \triangleq \gamma$  enforced at each iteration.
- Define  $\mathbf{s}_i \triangleq \gamma(\hat{\mathbf{x}}_i - \mathbf{r}_i)$  for  $i = 1, 2$ .
- The  $\gamma$ -forced VAMP manifests as<sup>13</sup>

$$\begin{aligned}\hat{\mathbf{x}}_1 &\leftarrow \arg \min_{\mathbf{x}_1} L_\gamma(\mathbf{x}_1, \hat{\mathbf{x}}_2, \mathbf{s}_1) \\ \mathbf{s}_2 &\leftarrow \mathbf{s}_1 + \gamma(\hat{\mathbf{x}}_1 - \hat{\mathbf{x}}_2) \\ \hat{\mathbf{x}}_2 &\leftarrow \arg \min_{\mathbf{x}_2} L_\gamma(\hat{\mathbf{x}}_1, \mathbf{x}_2, \mathbf{s}_2) \\ \mathbf{s}_1 &\leftarrow \mathbf{s}_2 + \gamma(\hat{\mathbf{x}}_1 - \hat{\mathbf{x}}_2)\end{aligned}$$

which is equivalent to Peaceman-Rachford splitting.

- This agrees with the earlier result that VAMP needs no damping when  $\gamma_1 = \gamma_2$  (for strictly convex  $f_1, f_2$ ).

<sup>13</sup>Fletcher, Sahraee, Rangan, Schniter'16

# Outline

- 1 Linear Regression
- 2 EC, EP, and VAMP
- 3 EM-VAMP and Adaptive VAMP
- 4 Damped VAMP and Connections to ADMM
- 5 Plug-and-play VAMP & Whitening**
- 6 VAMP as a Deep Neural Network

## Plug-and-play VAMP

- Recall that the nonlinear estimation step in VAMP (or AMP)

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

can be interpreted as “denoising” the pseudo-measurement  $\mathbf{r}_1$ .

- For certain signal classes, very sophisticated non-scalar denoising procedures have been developed (e.g., [BM3D](#) for images).
- Such denoising procedures can be “plugged into” signal recovery algorithms like ADMM<sup>14</sup>, AMP<sup>15</sup>, or VAMP<sup>16</sup>.
- For AMP and VAMP, the divergence can be approximated using Monte-Carlo:

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

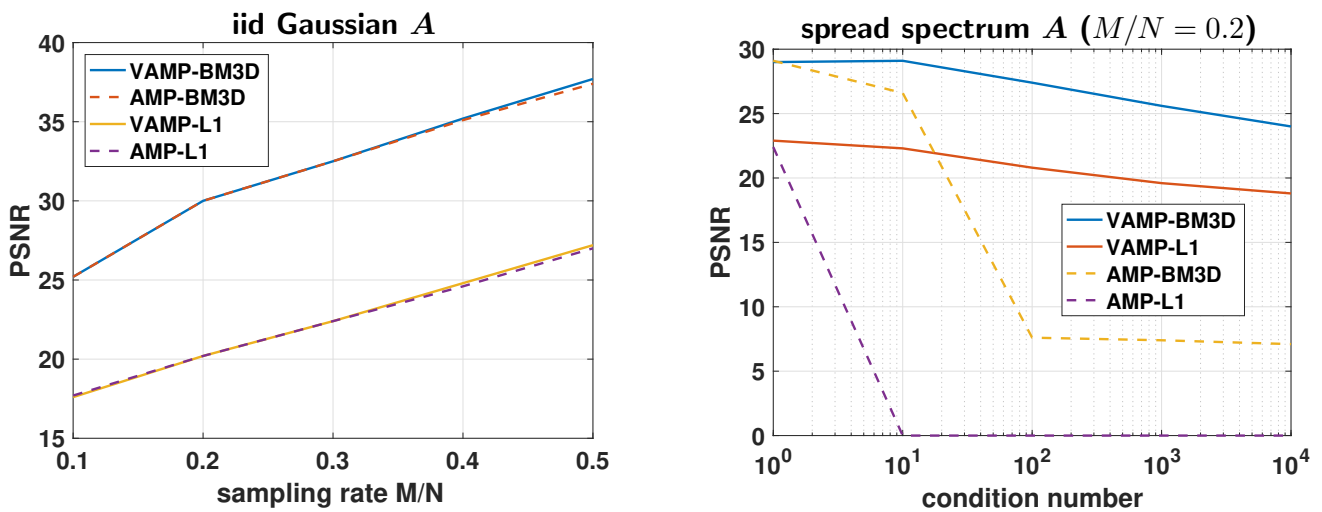
with random vectors  $\mathbf{p}_k \in \{\pm 1\}^N$  and small  $\epsilon > 0$ . Often,  $K = 1$  suffices.

<sup>14</sup>Bouman et al'13, <sup>15</sup>Metzler, Maleki, Baraniuk'14, <sup>16</sup>Schniter, Rangan, Fletcher'16



## Experiment: Image Recovery with Random Matrices

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



Results above are averaged over  $128 \times 128$  versions of

*lena, barbara, boat, fingerprint, house, peppers*

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

## Image Recovery with Non-Random Matrices

- Many imaging applications (e.g., MRI) use **low-frequency Fourier** measurements, in which case  $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top = \mathbf{I} [\mathbf{I} \mathbf{0}] \mathbf{F}$ .
- This causes problems for VAMP because the signal correlation structure interacts with  $\mathbf{F}$  in a way that VAMP is not designed to handle.
- Why? Say  $\mathbf{x}$  is a natural image, and consider  $\mathbf{q} = \mathbf{V}^\top \mathbf{x}$ .
  - If  $\mathbf{V}$  is large and Haar, then  $\mathbf{q}$  will be iid Gaussian.
  - If  $\mathbf{V}^\top = \mathbf{F}$ , the low-freq entries of  $\mathbf{q}$  will be much stronger than the others.

*VAMP treats  $\mathbf{V}^\top \mathbf{x}$  as iid Gaussian and thus diverges when  $\mathbf{V}^\top = \mathbf{F}$ !*

# Whitened VAMP for Image REcovery (VAMPire)

- To apply VAMP with non-random Fourier measurements, we propose to operate on the whitened signal:

$$\mathbf{y} = \underbrace{[\mathbf{I} \ \mathbf{0}] \mathbf{F} \mathbf{R}_x^{1/2}}_{\mathbf{A}} \mathbf{s} + \mathbf{w} \quad \text{for} \quad \begin{cases} \mathbf{R}_x = \mathbb{E}\{\mathbf{x}\mathbf{x}^\top\} \\ \mathbf{s} = \text{whitened signal coefficients} \end{cases}$$

and perform plug-and-play denoising from the whitened-coefficient space:

$$\hat{\mathbf{s}}_1 = \mathbf{g}_1(\mathbf{r}_1, \gamma_1) = \mathbf{R}_x^{-1/2} \text{denoise}(\mathbf{R}_x^{1/2} \mathbf{r}_1; \gamma_1 N / \text{tr}(\mathbf{R}_x)).$$

- In practice, we approximate  $\mathbf{R}_x \approx \mathbf{W}^\top \text{Diag}(\boldsymbol{\tau})^2 \mathbf{W}$ , where  $\mathbf{W}$  is a wavelet transform and  $\tau_i^2$  specifies the energy of the  $i$ th wavelet coefficient (which is easy to predict for natural images).

## Whitened VAMP for Image REcovery (VAMPire)

- The resulting matrix  $\mathbf{A} = [\mathbf{I} \ \mathbf{0}] \mathbf{F} \mathbf{W} \text{Diag}(\boldsymbol{\tau})$  does not yield a right singular vector matrix  $\mathbf{V}$  with a fast multiplication.
- But since  $\mathbf{A}$  has a fast implementation, the LMMSE stage can be computed via (preconditioned) LSQR:

$$\mathbf{g}_2(\mathbf{r}_2; \gamma_2) = (\gamma_w \mathbf{A}^\top \mathbf{A} + \gamma_2 \mathbf{I})^{-1} (\gamma_w \mathbf{A}^\top \mathbf{y} + \gamma_2 \mathbf{r}_2) = \begin{bmatrix} \sqrt{\gamma_w} \mathbf{A} \\ \sqrt{\gamma_2} \mathbf{I} \end{bmatrix}^+ \begin{bmatrix} \sqrt{\gamma_w} \mathbf{y} \\ \sqrt{\gamma_2} \mathbf{r}_2 \end{bmatrix}$$

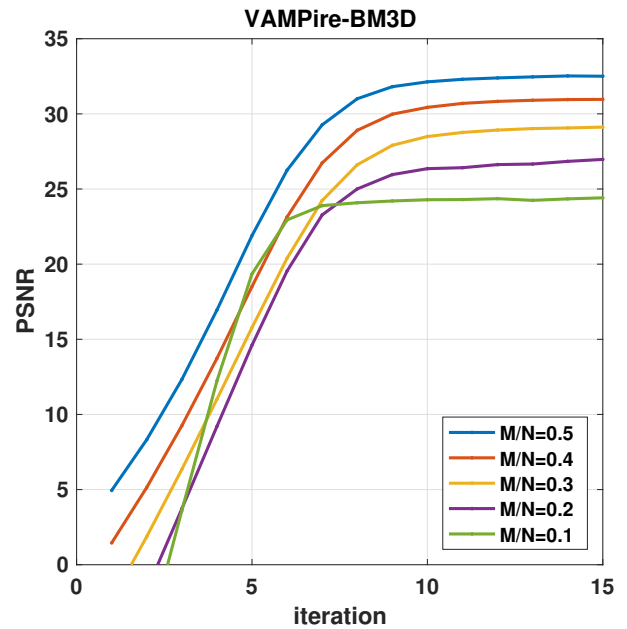
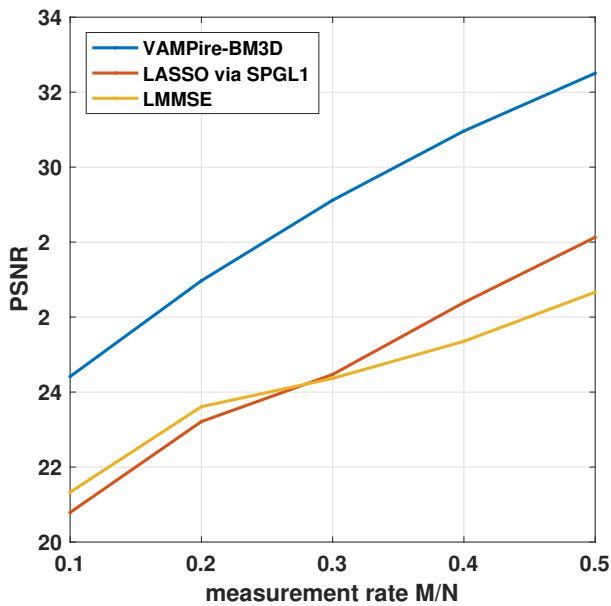
- The divergence  $\langle \mathbf{g}'_2(\mathbf{r}_2; \gamma_2) \rangle$  can be approximated using Monte-Carlo:

$$\langle \mathbf{g}'_2 \rangle = \frac{\gamma_2}{N} \text{tr} \left[ \left( \gamma_w \mathbf{A}^\text{H} \mathbf{A} + \gamma_2 \mathbf{I} \right)^{-1} \right] \approx \frac{1}{NK} \sum_{k=1}^K \mathbf{p}_k \begin{bmatrix} \sqrt{\gamma_w} \mathbf{A} \\ \sqrt{\gamma_2} \mathbf{I} \end{bmatrix}^+ \begin{bmatrix} \mathbf{0} \\ \sqrt{\gamma_2} \mathbf{p}_k \end{bmatrix},$$

where  $\mathbb{E}\{\mathbf{p}_k \mathbf{p}_k^\text{H}\} = \mathbf{I}$ . Here again, (preconditioned) LSQR can be used. In practice,  $K = 1$  suffices.

# Image Recovery Experiments

- Fourier measurements sampled at  $M$  lowest frequencies
- SNR=40dB
- $128 \times 128$  images  $\{lena, barbara, boat, fingerprint, house, peppers\}$
- db1 wavelet decomposition,  $D = 2$  levels



# Outline

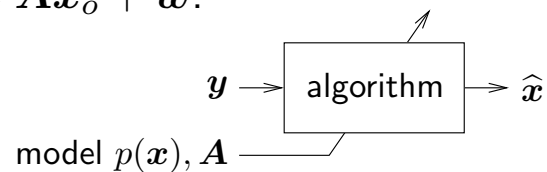
1 Linear Regression

2 EC, EP, and VAMP

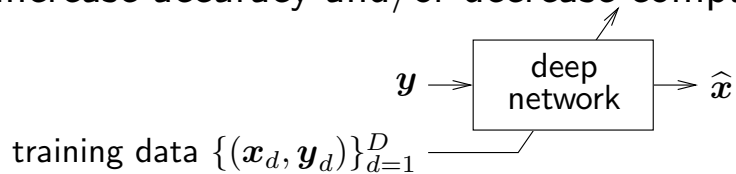


## Deep learning for sparse reconstruction

- 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?



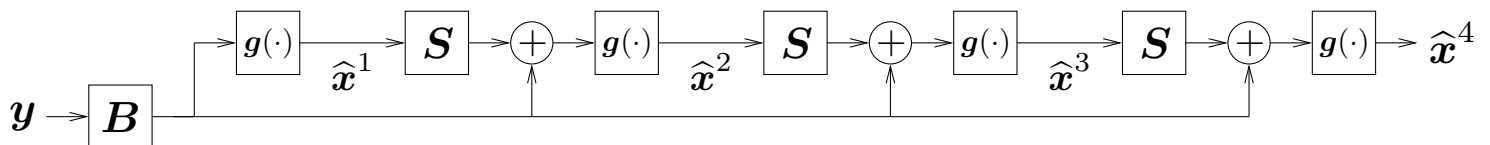
- Are there **connections** between these approaches?

## Unfolding Algorithms into Networks

Consider, e.g., the classical sparse-reconstruction algorithm, [ISTA](#).<sup>17</sup>

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

Gregor & LeCun<sup>18</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>19</sup>

<sup>17</sup>Daubechies,Defrise,DeMol'04.

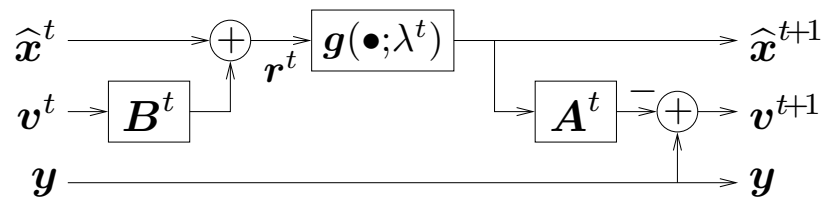
<sup>18</sup>Gregor,LeCun'10.

<sup>19</sup>Borgerding,Schniter'16.



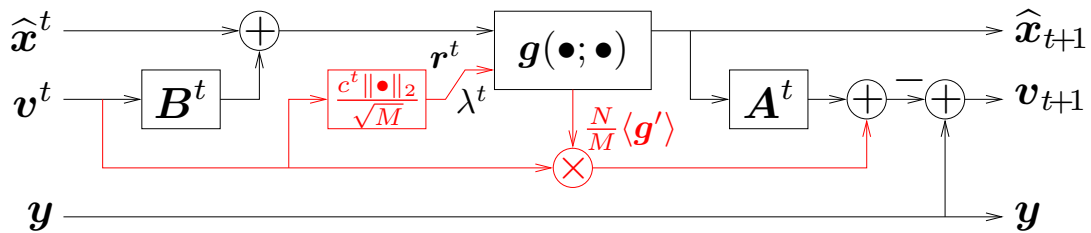
# Onsager-Corrected Deep Networks

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



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

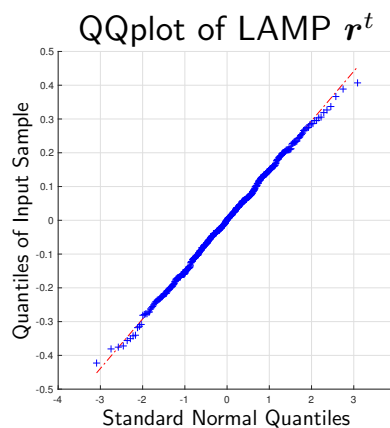
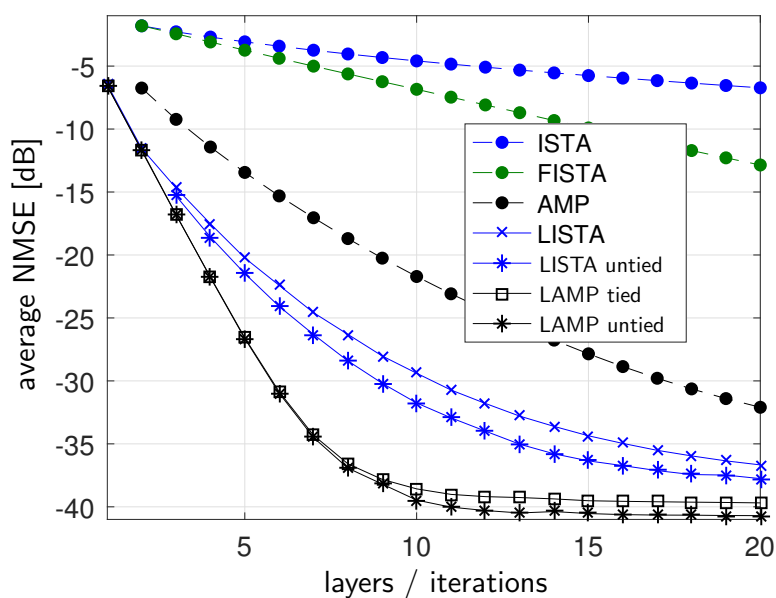
$t^{\text{th}}$  LAMP layer:



Onsager correction now aims to decouple errors across layers.

# LAMP performance with soft-threshold denoising

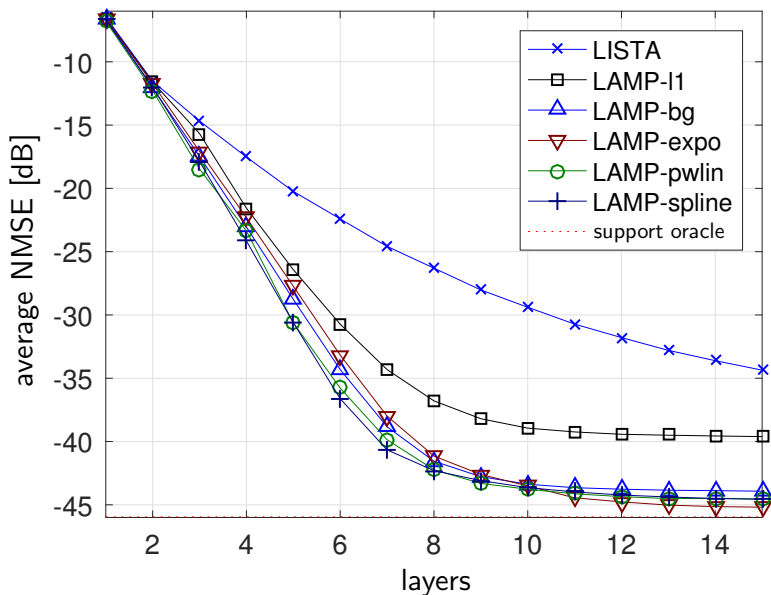
LISTA beats AMP, FISTA, ISTA in convergence speed and asymptotic MSE.  
 LAMP beats LISTA



## LAMP with more sophisticated scalar denoisers

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

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



Here we learned the parameters of these denoiser families:

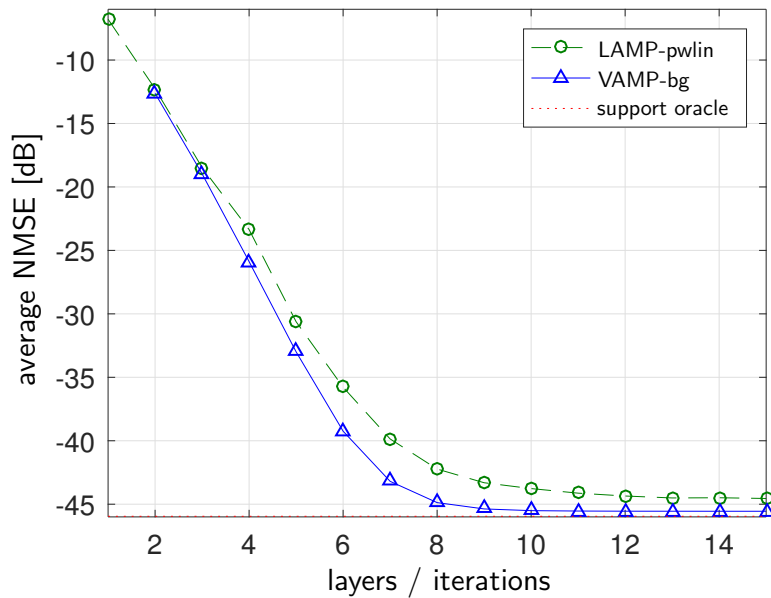
- scaled soft-thresholding
- conditional mean under BG
- Exponential kernel<sup>20</sup>
- Piecewise Linear<sup>20</sup>
- Spline<sup>21</sup>

Big improvement!

<sup>20</sup>Guo, Davies'15. <sup>21</sup>Kamilov, Mansour'16.

# LAMP versus VAMP

How does our best **Learned AMP** compare to (unlearned) **VAMP**?



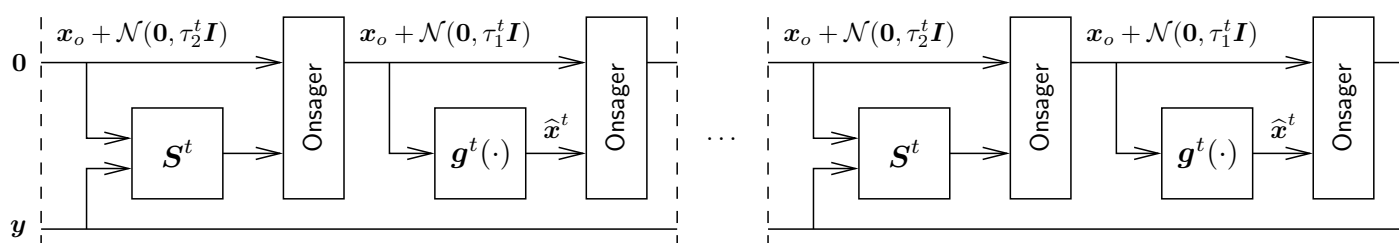
*VAMP wins!*

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.

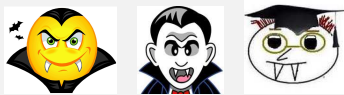


- Remarkably, **backpropagation does not improve matched VAMP!**

*VAMP is locally optimal*

- 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



- By merging classical algorithms (EC,EP) with SVDs and modern analysis tools (Bolthausen conditioning trick), we get “VAMP”: a **fast and robust algorithm** with a **rigorous SE analysis**.
- VAMP is like Peaceman-Rachford splitting with adaptive penalties  $\gamma_1, \gamma_2$ .
- Can combine with EM methods to handle priors/likelihood with unknown parameters, **rigorously**.
- Can apply sophisticated denoising algorithms (BM3D) via “plug and play.”
- Damping and whitening may be needed in practice, when  $\mathbf{A}$  is not random and  $\mathbf{x}$  is not iid.
- Can unroll VAMP into an deep network that is interpretable (and in some cases optimal).
- Lots more to do! (Generalized linear models, multilayer, bilinear models. . .)

**Stay tuned for Sundeep's talk!**