

# Onsager Correction, Deep Learning, Sparse Reconstruction and VAMP

**Mark Borgerding and Phil Schniter**



**THE OHIO STATE UNIVERSITY**

---

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

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

IEEE GlobalSIP — Dec 8, 2016

# Sparse Reconstruction

## Goal:

Recover  $\mathbf{x}_o \in \mathbb{R}^N$  from measurements  $\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathbf{w} \in \mathbb{R}^M$

## Assumptions:

- $\mathbf{x}_o$  is sparse
- $\mathbf{A}$  is known and high dimensional
- often  $M \ll N$
- $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \tau_w \mathbf{I})$

# Regularized loss minimization

Popular approach:

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

where

- $f(\mathbf{x})$  is a **regularizer**, e.g.,  $\|\mathbf{x}\|_1$  in LASSO or BPDN
- $\lambda > 0$  is a **tuning parameter**

# The iterative soft thresholding algorithm (ISTA)

ISTA:

initialize  $\hat{\mathbf{x}}^0 = \mathbf{0}$

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

$$\mathbf{v}^t = \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}^t \quad \text{residual error}$$

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

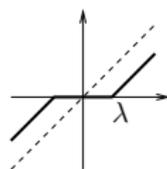
where

$$\mathbf{g}(\mathbf{r}) = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{r} - \mathbf{x}\|_2^2 + \lambda f(\mathbf{x}) \triangleq \text{prox}_{\lambda f}(\mathbf{r})$$

$\|\mathbf{A}\|_2^2 < 1$  ensures convergence<sup>1</sup> with convex  $f(\cdot)$ .

When  $f(\mathbf{x}) = \|\mathbf{x}\|_1$  we get “soft thresholding”

$$[\mathbf{g}(\mathbf{r})]_j = \text{sgn}(r_j) \max\{0, |r_j| - \lambda\}$$



<sup>1</sup>Daubechies, Defrise, DeMol-CPAM'04

# Approximate Message Passing (AMP)

Donoho, Maleki, and Montanari<sup>2</sup> proposed:

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

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

$$\mathbf{v}^t = \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}^t + \frac{N}{M}\mathbf{v}^{t-1}\langle \mathbf{g}^{t-1}'(\hat{\mathbf{x}}^{t-1} + \mathbf{A}^T\hat{\mathbf{v}}^{t-1}) \rangle$$

corrected residual

$$\hat{\mathbf{x}}^{t+1} = \mathbf{g}^t(\hat{\mathbf{x}}^t + \mathbf{A}^T\mathbf{v}^t)$$

thresholding

where

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

Note:

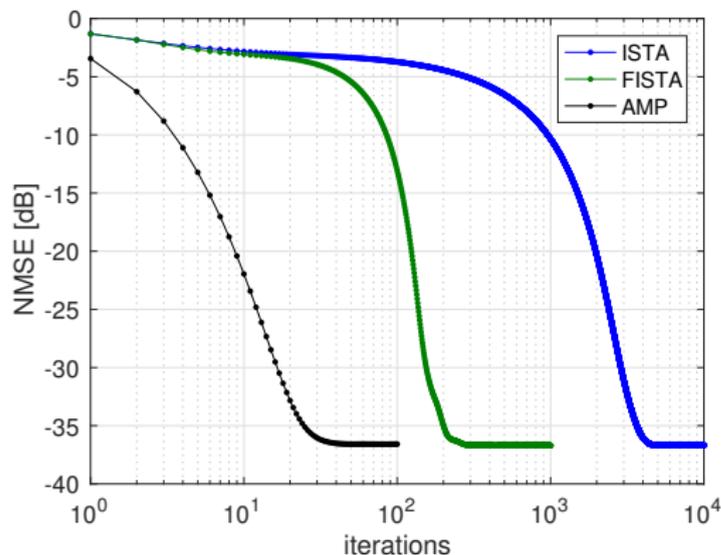
- “Onsager correction” aims to decouple the errors across iterations.
- The thresholding  $\mathbf{g}^t(\cdot)$  can vary with iteration  $t$ .

---

<sup>2</sup>Donoho, Maleki, Montanari–PNAS’09

# AMP vs ISTA (and FISTA)

Example: LASSO problem with i.i.d. Gaussian  $\mathbf{A}$ :



- $M = 250, N = 500$
- $\Pr\{x_n \neq 0\} = 0.1$
- SNR= 40dB
- Convergence to  $-35\text{dB}$ :
  - ISTA: 2407 iterations
  - FISTA:<sup>3</sup>174 iterations
  - AMP: 25 iterations

<sup>3</sup>Beck, Teboulle–JIS'09

# The limitations of AMP

- AMP's good performance is guaranteed only for large i.i.d. zero-mean sub-Gaussian  $\mathbf{A}$ .
- Deviations from this condition can cause AMP to diverge.
- Can we extend AMP to a larger class of matrices?

# Vector AMP (VAMP)



The recent **vector AMP**<sup>4</sup> algorithm for linear regression is

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

$$\hat{\mathbf{x}}_1^t = \mathbf{g}(\mathbf{r}_1^t; \gamma_1^t) \quad \text{thresholding}$$

$$\alpha_1^t = \frac{1}{N} \sum_j \frac{\partial g_r}{\partial r_j}(\mathbf{r}_1^t; \gamma_1^t) \quad \text{divergence}$$

$$\mathbf{r}_2^t = \frac{1}{1-\alpha_1^t} (\hat{\mathbf{x}}_1^t - \alpha_1^t \mathbf{r}_1^t) \quad \text{Onsager correction}$$

$$\gamma_2^t = \gamma_1^t \frac{1-\alpha_1^t}{\alpha_1^t} \quad \text{precision of } \mathbf{r}_2^t$$

$$\hat{\mathbf{x}}_2^t = (\mathbf{A}^\top \mathbf{A} / \hat{\tau}_w + \gamma_2^t \mathbf{I})^{-1} (\mathbf{A}^\top \mathbf{y} / \hat{\tau}_w + \gamma_2^t \mathbf{r}_2^t) \quad \text{LMMSE}$$

$$\alpha_2^t = \frac{\gamma_2^t}{N} \text{Tr} [(\mathbf{A}^\top \mathbf{A} / \hat{\tau}_w + \gamma_2^t \mathbf{I})^{-1}] \quad \text{divergence}$$

$$\mathbf{r}_1^{t+1} = \frac{1}{1-\alpha_2^t} (\hat{\mathbf{x}}_2^t - \alpha_2^t \mathbf{r}_2^t) \quad \text{Onsager correction}$$

$$\gamma_1^{t+1} = \gamma_2^t \frac{1-\alpha_2^t}{\alpha_2^t} \quad \text{precision of } \mathbf{r}_1^{t+1}$$

<sup>4</sup>Rangan, Schniter, Fletcher – arxiv:1610.03082

# VAMP without matrix inverses



LMMSE **matrix inverse** step is easy if the **SVD**  $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$  is known.  
Note if  $\kappa(\mathbf{A}) = 1$ <sup>5</sup>, the SVD is trivial  $\mathbf{U}\mathbf{S} \propto \mathbf{I}, \mathbf{V} \propto \mathbf{A}^T$ .

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

$$\hat{\mathbf{x}}^t = \mathbf{g}(\mathbf{r}_1^t; \gamma_1^t) \quad \text{thresholding}$$

$$\alpha_1^t = \frac{1}{N} \sum_j \frac{\partial g}{\partial r_j}(\mathbf{r}_1^t; \gamma_1^t) \quad \text{divergence}$$

$$\mathbf{r}_2^t = \frac{1}{1-\alpha_1^t} (\hat{\mathbf{x}}^t - \alpha_1^t \mathbf{r}_1^t) \quad \text{Onsager}$$

$$\gamma_2^t = \gamma_1^t \frac{1-\alpha_1^t}{\alpha_1^t} \quad \text{precision}$$

$$\alpha_2^t = \frac{1}{N} \sum_j \gamma_2^t / (s_j^2 / \hat{\tau}_w + \gamma_2^t) \quad \text{divergence}$$

$$\mathbf{r}_1^{t+1} = \mathbf{r}_2^t + \frac{1}{1-\alpha_2^t} \mathbf{V} (\mathbf{S}^2 + \hat{\tau}_w \gamma_2^t \mathbf{I})^{-1} \mathbf{S} (\mathbf{U}^T \mathbf{y} - \mathbf{S}\mathbf{V}^T \mathbf{r}_2^t) \quad \text{2 mat-vec}$$

$$\gamma_1^{t+1} = \gamma_2^t \frac{1-\alpha_2^t}{\alpha_2^t} \quad \text{precision}$$

<sup>5</sup>e.g. punctured DFT, DCT, DWT

## Why call this “Vector AMP”?

- 1) Can be derived using an **approximation of message passing** on a factor graph, now with **vector-valued** variable nodes.
- 2) Performance characterized by a rigorous **state-evolution**<sup>6</sup> under certain large random  $A$ :

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

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

*“ $\mathbf{A}$  is rotationally invariant”*

---

<sup>6</sup>Rangan, Fletcher, Schniter–16

# Message-passing derivation of VAMP

- Write **joint density** as  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x}) = p(\mathbf{x})\mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \tau_w \mathbf{I})$

$$p(\mathbf{x}) \blacksquare \text{---} \circ \overset{\mathbf{x}}{\text{---}} \blacksquare \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \tau_w \mathbf{I})$$

- Variable **splitting**:  $p(\mathbf{x}_1, \mathbf{x}_2, \mathbf{y}) = p(\mathbf{x}_1)\delta(\mathbf{x}_1 - \mathbf{x}_2)\mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}_2, \tau_w \mathbf{I})$

$$p(\mathbf{x}_1) \blacksquare \text{---} \circ \overset{\mathbf{x}_1}{\text{---}} \blacksquare \underset{\delta(\mathbf{x}_1 - \mathbf{x}_2)}{\text{---}} \circ \overset{\mathbf{x}_2}{\text{---}} \blacksquare \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}_2, \tau_w \mathbf{I})$$

- Perform message-passing with messages approximated as  $\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ .
  - An instance of **expectation-propagation**<sup>7</sup> (EP).
  - Also derivable through **expectation-consistent approximation**<sup>8</sup> (EC).

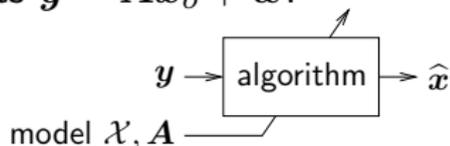
---

<sup>7</sup> Minka–Dissertation’01

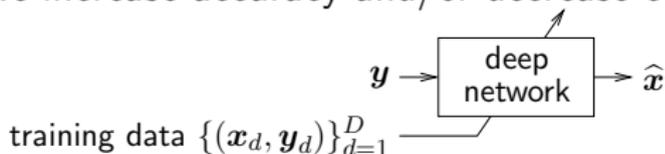
<sup>8</sup> Oppor, Winther–NIPS’04, Fletcher, Rangan, Schniter–ISIT’16

# Deep learning for sparse reconstruction

- Until now we've focused on **designing algorithms** to recover  $x_o \in \mathcal{X}$  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?



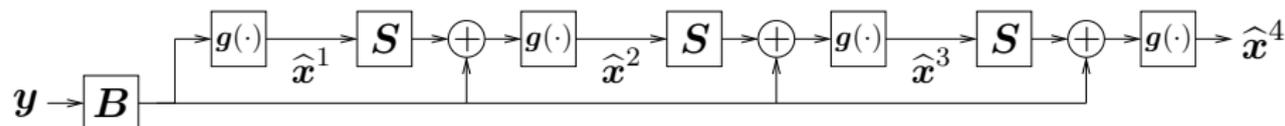
- Are there **connections** between these approaches?

# Unrolling ISTA

First, rewrite ISTA as

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

Then “unroll” into a network:

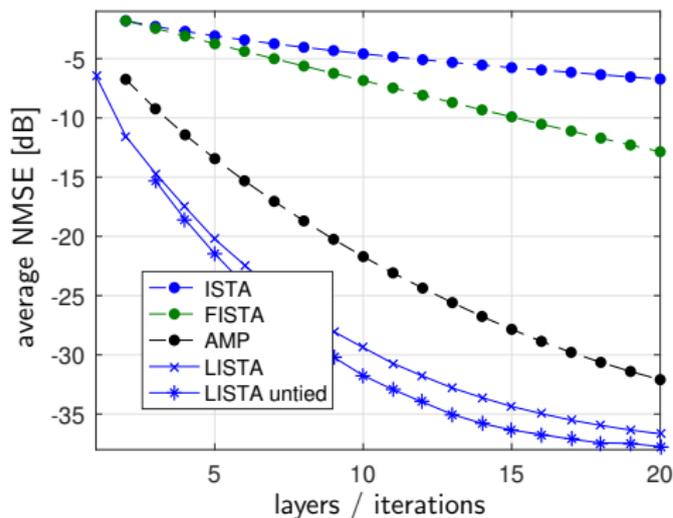


Note cascade of linear “ $\mathbf{S}$ ,” bias “ $\mathbf{B}\mathbf{y}$ ,” & separable non-linearity “ $\mathbf{g}(\cdot)$ .”

ISTA algorithm  $\Leftrightarrow$  deep neural network

# Learned ISTA (LISTA)

Gregor and LeCun<sup>9</sup> proposed to **learn (via backpropagation)** the linear transform  $\mathcal{S}$  and soft thresholds  $\{\lambda^t\}_{t=1}^T$  that minimize training MSE



$$\arg \min_{\Theta} \sum_{d=1}^D \left\| \hat{\mathbf{x}}(\mathbf{y}_d; \Theta) - \mathbf{x}_d \right\|^2.$$

The resulting “LISTA” beats LASSO-AMP in convergence speed *and* asymptotic MSE!

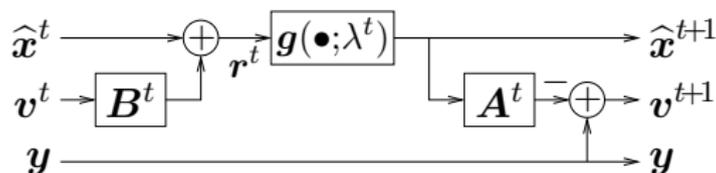
Further improvement when  $\mathcal{S}$  is “untied” to  $\{\mathcal{S}^t\}_{t=1}^T$ .

<sup>9</sup>Gregor, LeCun-ICML'10

# Learned AMP (LAMP)

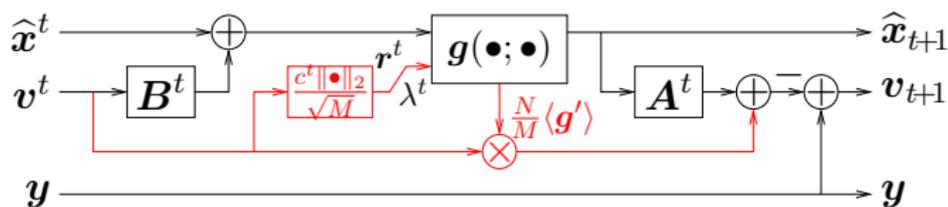


$t^{\text{th}}$  LISTA layer:



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

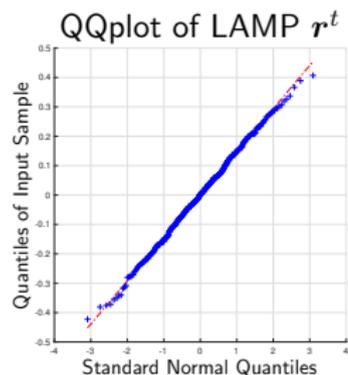
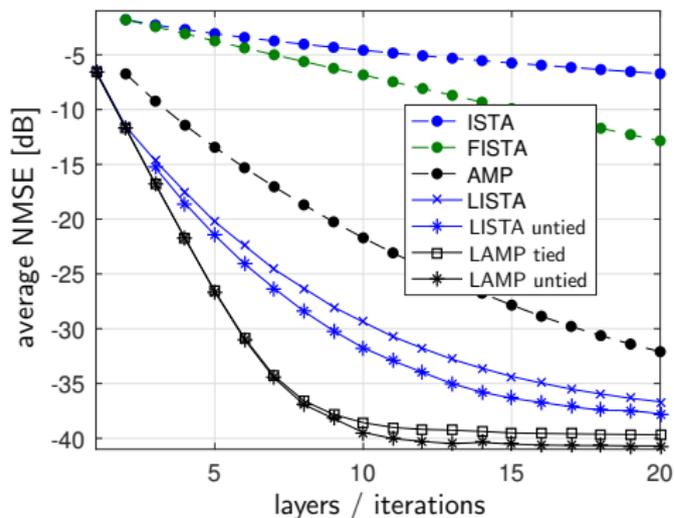
$t^{\text{th}}$  LAMP layer:



Onsager correction now aims to decouple errors across layers.

# LAMP performance under soft thresholding

LAMP beats LISTA in both convergence speed and asymptotic MSE.



## Simulation/Training Details

- Bernoulli-Gaussian  $\mathbf{x} \in \mathbb{R}^{500}$ , known Gaussian  $\mathbf{A} \in \mathbb{R}^{250 \times 500}$ , 40dB SNR
- TensorFlow implementation on GPU
- Adam<sup>10</sup>SGD in mini-batches of 1000 vectors
- Add layers by greedy extension, then whole-network fine-tuning
- Reported results are for an untrained validation set

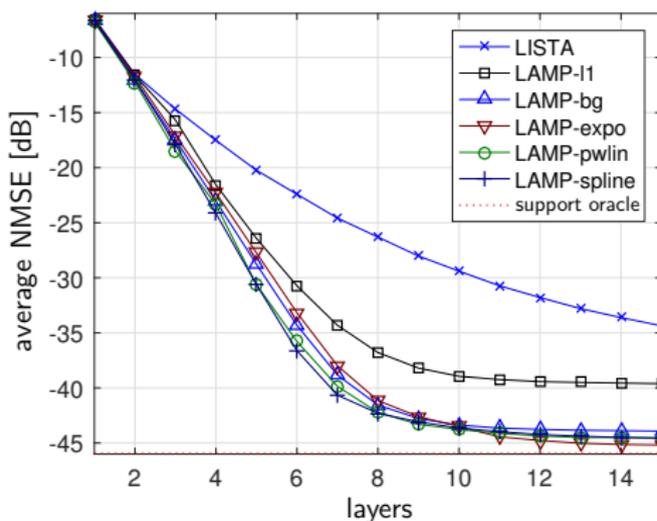
---

<sup>10</sup>Kingma, BA'15

# LAMP with more sophisticated denoisers

So far, we used soft-thresholding, so as to compare directly to LISTA.

What happens when we learn **other denoisers**?



Here we learned the parameters of these denoiser families:

- scaled soft-threshold ( $LAMP-\ell_1$ )
- Bernoulli-Gaussian MMSE
- Exponential kernel<sup>11</sup>
- Piecewise Linear<sup>11</sup>
- Spline<sup>12</sup>

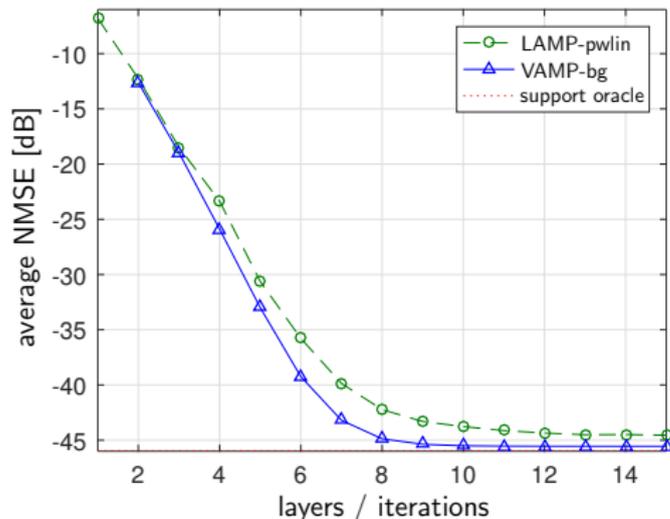
Big improvement!

<sup>11</sup> Guo, Davies-TSP'15

<sup>12</sup> Kamilov, Mansour-SPL'16



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

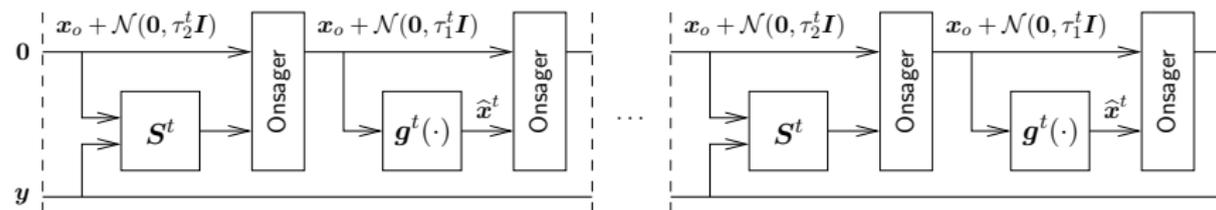


*VAMP wins!*

So what about “learned VAMP”?

# Learned VAMP

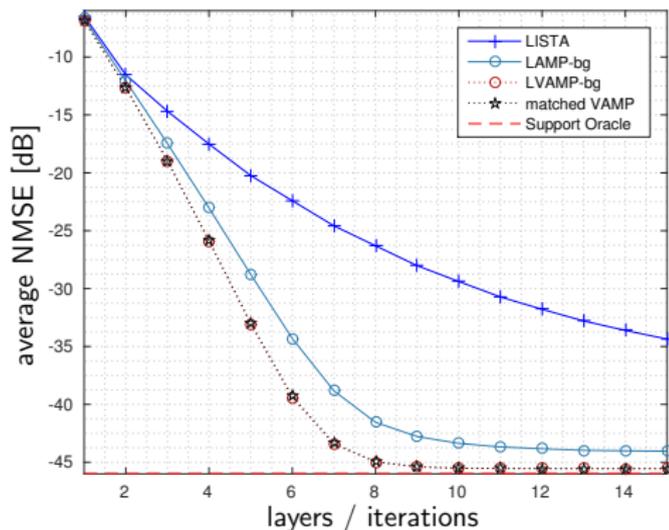
- Suppose we unroll VAMP and learn (via backprop) the parameters  $\{\mathbf{S}^t, \mathbf{g}^t\}_{t=1}^T$  that minimize the training MSE.



Can we improve VAMP with learning?

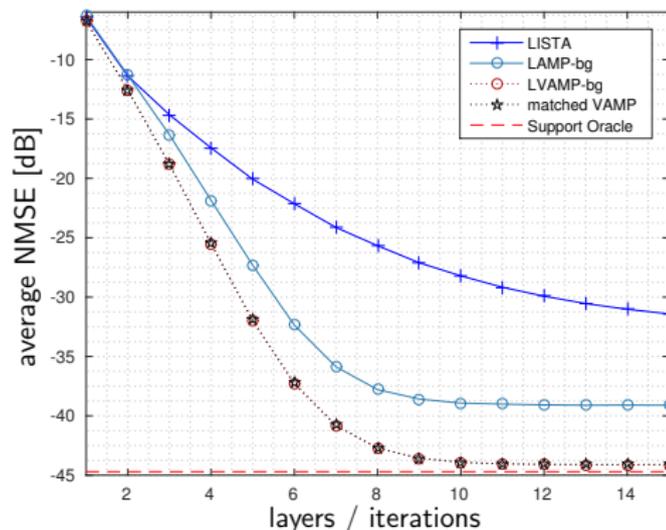
# Learning VAMP

- Remarkably, backpropagation does **not** improve **matched** VAMP!
- Matched VAMP specifies optimal network parameters in closed form, without training.



# Non iid sub-Gaussian $\mathbf{A}$

- What happens if  $\mathbf{A}$  is not iid sub-Gaussian?
- We replaced the singular values from our previous matrix with a log-spaced sequence so  $\kappa(\mathbf{A}) = 15$ .
- Basic AMP fails to converge!
- VAMP's advantages become even clearer!
- In large-system limit, VAMP can handle arbitrarily large  $\kappa(\mathbf{A})$



# Conclusions

- For sparse reconstruction, AMP has some nice properties:
  - low cost-per-iteration
  - fast convergence,
  - rigorous state evolution,but only under large i.i.d. Gaussian  $A$ .
- Vector AMP has the same nice properties under large rotationally invariant  $A$ .
- “Learned ISTA” results from unrolling ISTA and fitting its parameters to training data. We proposed learned AMP & learned VAMP.
- MSE-optimal parameters of VAMP can be specified in closed form, without training, when signal/noise statistics are known.

For details, see Borgerding, Schniter arXiv 1612.01183