

Vector Approximate Message Passing and Connections to Deep Learning

Phil Schniter



THE OHIO STATE UNIVERSITY

Duke
UNIVERSITY



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

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

IEEE Information Theory Workshop (ITW) — Sep 13, 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 $t = 0, 1, 2, \dots$

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

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

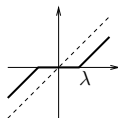
where

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

$\|\mathbf{A}\|_2^2 < 1$ ensures convergence¹ with convex $f(\cdot)$.

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

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



¹Daubechies, Defrise, DeMol-CPAM'04

Approximate Message Passing (AMP)

Donoho, Maleki, and Montanari² 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

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

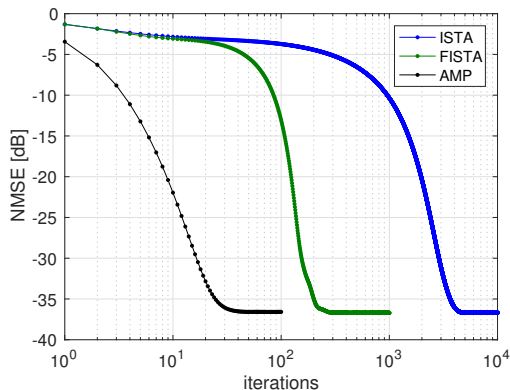
Note:

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

²Donoho, 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 -35dB :
 - ISTA: 2407 iterations
 - FISTA:³174 iterations
 - AMP: 25 iterations

³Beck, Teboulle–JIS'09

AMP's state evolution

Define $\mathcal{E}^t := \frac{1}{N} \mathbb{E} \{ \|\widehat{\mathbf{x}}^t - \mathbf{x}_o\|^2 \}$ as the iteration- t MSE.

For **large i.i.d. sub-Gaussian \mathbf{A}** and separable Lipschitz $\mathbf{g}^t(\cdot)$, AMP has the following scalar **state evolution** (SE):⁴

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

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

$$\mathcal{E}^{t+1} = \frac{1}{N} \mathbb{E} \left\{ \left\| \mathbf{g}^t \left(\underbrace{\mathbf{x}_o + \mathcal{N}(\mathbf{0}, \tau_r^t \mathbf{I})}_{:= \mathbf{r}^t} \right) - \mathbf{x}_o \right\|^2 \right\}$$

But for **generic \mathbf{A}** , AMP is not well justified and **may fail** catastrophically.

⁴Bayati, Montanari–TransIT'11

Vector AMP (VAMP)



The **vector AMP** algorithm for linear regression is

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

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

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

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

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

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

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

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

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

Note similarities with standard AMP.

VAMP without matrix inverses



Can avoid matrix inverses by pre-computing an SVD, $\mathbf{A} = \mathbf{USV}^T$:

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

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

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

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

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

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

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

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

And can tune noise precision $\hat{\tau}_w$ using EM.

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**⁵ under certain large random A :

$$SVD A = USV^T$$

- U is deterministic
- S is deterministic
- V is uniformly distributed on the group of orthogonal matrices

“ A is rotationally invariant”

⁵Rangan, 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**⁶ (EP).
 - Also derivable through **expectation-consistent approximation**⁷ (EC).

⁶Minka–Dissertation'01

⁷Opper, Winther–NIPS'04, Fletcher, Rangan, Schniter–ISIT'16

VAMP state evolution

Assuming empirical convergence of $\{s_j\} \rightarrow S$ and $\{(r_{1,j}^0, x_{o,j})\} \rightarrow (R_1^0, X_o)$ and Lipschitz continuity of g and g' , the SE under $\hat{\tau}_w = \tau_w$ is as follows:

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

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

$$\bar{\alpha}_1^t = \mathbb{E} \left\{ g'(X_o + \mathcal{N}(0, \tau_1^t); \bar{\gamma}_1^t) \right\} \quad \text{divergence}$$

$$\bar{\gamma}_2^t = \bar{\gamma}_1^t \frac{1 - \bar{\alpha}_1^t}{\bar{\alpha}_1^t}, \quad \tau_2^t = \frac{1}{(1 - \bar{\alpha}_1^t)^2} [\mathcal{E}_1^t - (\bar{\alpha}_1^t)^2 \tau_1^t] \quad \text{precision}$$

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

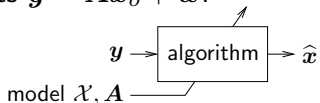
$$\bar{\alpha}_2^t = \bar{\gamma}_2^t \mathbb{E} \left\{ [S^2 / \tau_w + \bar{\gamma}_2^t]^{-1} \right\} \quad \text{divergence}$$

$$\bar{\gamma}_1^{t+1} = \bar{\gamma}_2^t \frac{1 - \bar{\alpha}_2^t}{\bar{\alpha}_2^t}, \quad \tau_1^{t+1} = \frac{1}{(1 - \bar{\alpha}_2^t)^2} [\mathcal{E}_2^t - (\bar{\alpha}_2^t)^2 \tau_2^t] \quad \text{precision}$$

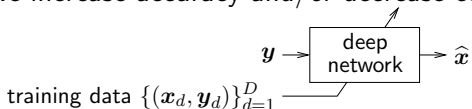
More complicated expressions for \mathcal{E}_2^t and $\bar{\alpha}_2^t$ apply when $\hat{\tau}_w \neq \tau_w$.

Deep learning for sparse reconstruction

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



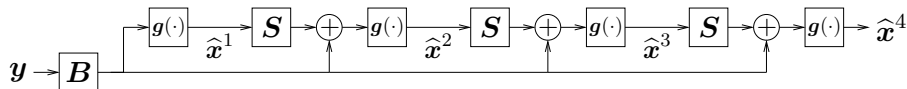
- Are there **connections** between these approaches?

Unrolling ISTA

First, rewrite ISTA as

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

Then “unroll” into a network:

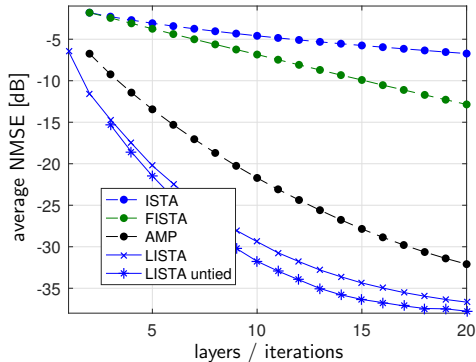


Note cascade of linear “ \mathbf{S} ,” bias “ $\mathbf{B}\mathbf{y}$,” & separable non-linearity “ $\mathbf{g}(\cdot)$.”

ISTA algorithm \Leftrightarrow deep neural network

Learned ISTA (LISTA)

Gregor and LeCun⁸ 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!

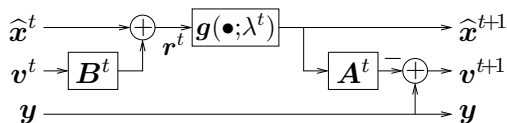
Further improvement when \mathcal{S} is “untied” to $\{\mathcal{S}^t\}_{t=1}^T$.

⁸Gregor, LeCun-ICML'10

Learned AMP (LAMP)

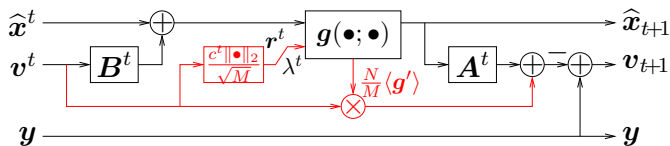


t^{th} LISTA layer:



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

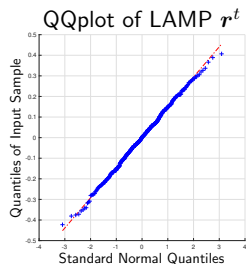
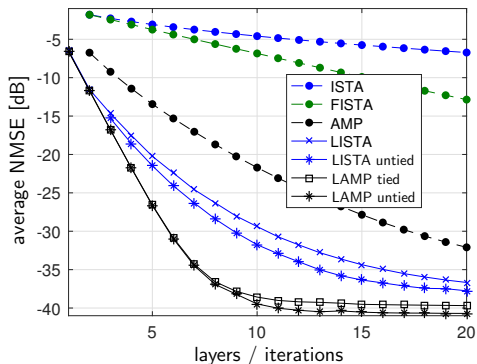
t^{th} LAMP layer:



Onsager correction now aims to decouple errors across layers.

LAMP performance under soft thresholding

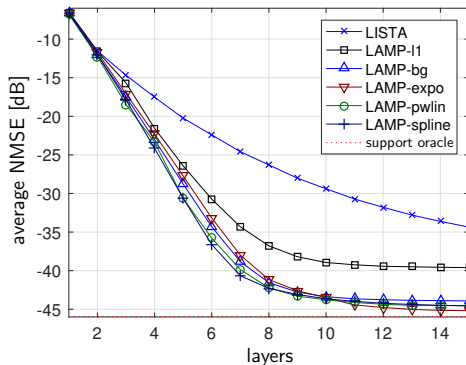
LAMP beats LISTA in both convergence speed and asymptotic MSE.



LAMP with more sophisticated 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
- Bernoulli-Gaussian MMSE
- Exponential kernel⁹
- Piecewise Linear⁹
- Spline¹⁰

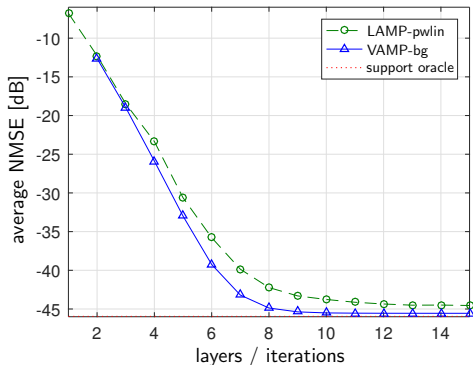
Big improvement!

⁹ Guo, Davies-TSP'15

¹⁰ Kamilov, Mansour-SPL'16



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

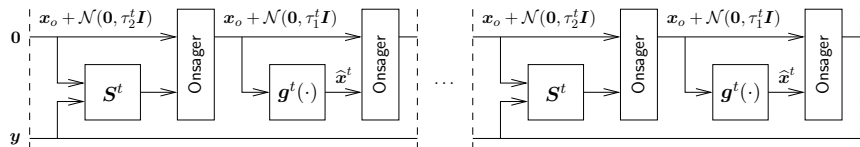


VAMP wins!

So what about “learned VAMP”?

Local optimality of 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.



- Remarkably, backpropagation does not improve matched VAMP!

VAMP is locally optimal

- Essentially, 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

- 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 .
- We proposed a Vector AMP, where the same nice properties hold under large rotationally invariant A .
- “Learned ISTA” results from unrolling ISTA and fitting its parameters to training data. We proposed learned AMP & learned VAMP.
- Remarkably, the original VAMP is locally optimal.