

# Vector Approximate Message Passing

**Phil Schniter**



THE OHIO STATE UNIVERSITY

Duke  
UNIVERSITY



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

Supported in part by NSF grant CCF-1527162.

iTWIST @ Aalborg University — Aug 24, 2016

# Standard Linear Regression

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

Examples:

- **Compressive Sensing / Medical Imaging:**

$\mathbf{y}$  = measurements     $\mathbf{x}_o$  = sparse image/signal representation

$\mathbf{w}$  = sensor noise     $\mathbf{A} = \Phi\Psi$ ,  $\Phi$  measurement operator,  $\Psi$  basis

- **Wireless communications:**

$\mathbf{y}$  = received samples     $\mathbf{x}_o$  = finite-alphabet symbols

$\mathbf{w}$  = noise & interference     $\mathbf{A}$  = channel operator

- **Statistics / Machine Learning:**

$\mathbf{y}$  = experimental outcomes     $\mathbf{x}_o$  = prediction coefficients

$\mathbf{w}$  = model error     $\mathbf{A}$  = feature data

# Implicit assumptions used in most of this talk

## Standard linear regression:

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

- $\mathbf{A}$  is a known and high dimensional (e.g.,  $M, N \gtrsim 100$ )
- often  $N \gg M$  (more unknowns than observations)
- $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \tau_w \mathbf{I})$  (additive white Gaussian noise)
- $\mathbf{x}_o$  is “structured” (e.g., sparse, natural image, etc.)
- quantities are real-valued (but can be easily extended to complex-valued)

Later will describe extension to [generalized linear model](#):

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

# Regularized loss minimization

One way to approach this problem is

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

where

- $\frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|^2$  is the **quadratic loss function**
- $f(\mathbf{x})$  is a suitably chosen **regularizer**
  - convex  $f(\cdot)$  leads to a convex optimization problem
  - choosing  $f(\mathbf{x}) = \|\mathbf{x}\|_1$  yields sparse  $\hat{\mathbf{x}}$
- $\lambda > 0$  is a **tuning parameter**

**Bayesian** interpretation:

$$\hat{\mathbf{x}} = \text{MAP estimate of } \mathbf{x} \text{ under } \begin{cases} \text{likelihood } p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mathbf{Ax}, \tau_w \mathbf{I}) \\ \text{prior } p(\mathbf{x}) \propto \exp(-\lambda f(\mathbf{x})/\tau_w) \end{cases}$$

# Iterative thresholding

One approach to regularized loss minimization:

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

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

$$\begin{aligned} \mathbf{v}^t &= \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}^t && \text{compute residual} \\ \hat{\mathbf{x}}^{t+1} &= \mathbf{g}(\hat{\mathbf{x}}^t + \mathbf{A}^\top \mathbf{v}^t) && \text{thresholding} \end{aligned}$$

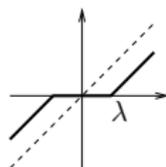
where

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

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

For example,  $f(\mathbf{x}) = \|\mathbf{x}\|_1$  gives “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)

A modification of iterative thresholding:

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

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

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

where

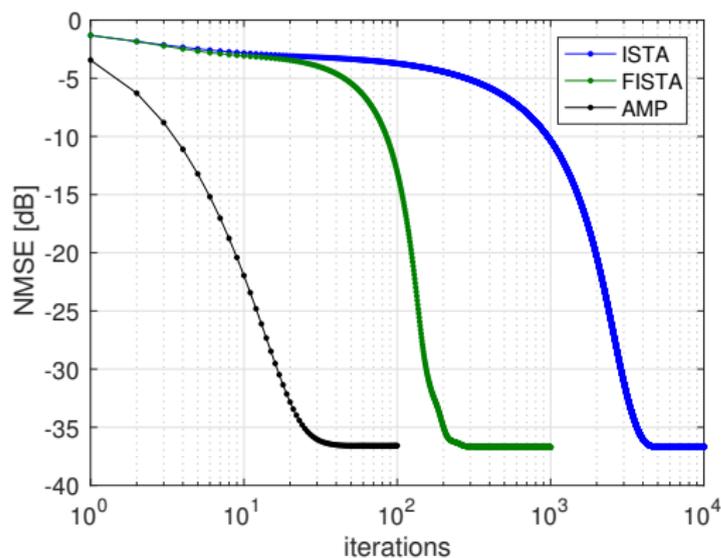
$$\langle \mathbf{g}'(\mathbf{r}) \rangle \triangleq \frac{1}{N} \sum_{j=1}^N \frac{\partial g_j(\mathbf{r})}{\partial r_j} \quad \text{“divergence.”}$$

Note:

- The residual  $\mathbf{v}^t$  now includes an “Onsager correction.”
- The thresholding  $\mathbf{g}^t(\cdot)$  can vary with iteration  $t$ .
- Can be derived using Gaussian & Taylor-series approximations of min-sum belief-propagation / message passing.

# AMP vs ISTA (and FISTA)

Typical convergence behavior with i.i.d. Gaussian  $\mathbf{A}$ :



Experiment:

- $M = 250, N = 500$
- $\Pr\{x_n \neq 0\} = 0.1$
- SNR = 40dB
- ISTA, FISTA<sup>2</sup>, AMP all reach same solution: NMSE = -36.8dB
- Convergence to -35dB:
  - ISTA: 2407 iterations
  - FISTA: 174 iterations
  - AMP: 25 iterations

<sup>2</sup>Beck, Teboulle—JIS'09

# AMP's denoising property

## Assumption 1

- $\mathbf{A} \in \mathbb{R}^{M \times N}$  is i.i.d. Gaussian
- $M, N \rightarrow \infty$  s.t.  $\frac{M}{N} = \delta \in (0, \infty)$
- $f(\mathbf{x}) = \sum_{j=1}^N f(x_j)$  with Lipschitz  $f$

Under Assumption 1, something remarkable happens to the input to the thresholder:<sup>3</sup>

$$\mathbf{r}^t \triangleq \widehat{\mathbf{x}}^t + \mathbf{A}^T \mathbf{v}^t = \mathbf{x}_o + \mathcal{N}(\mathbf{0}, \tau_r^t \mathbf{I})$$

with  $\tau_r^t = \frac{1}{M} \|\mathbf{v}^t\|^2 \triangleq \widehat{\tau}_r^t$

In other words,  $\mathbf{r}^t$  is a noisy version of the true signal  $\mathbf{x}_o$ , where the noise is Gaussian with known variance.

---

<sup>3</sup>Bayati, Montanari–TransIT'11

## AMP's state evolution

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

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

Under Assumption 1, AMP has the following scalar **state evolution** (SE):

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

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

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

The rigorous proof<sup>4</sup> of the SE uses Bolthausen's conditioning trick from the statistical physics literature.

---

<sup>4</sup>Bayati, Montanari–TransIT'11

# Choice of denoiser in AMP

## 1) LASSO/BPDN

- Goal: compute “ $\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \lambda \|\mathbf{x}\|_1$ .”
- Use  $\mathbf{g}^t(\mathbf{r}) = \text{soft}(\mathbf{r}; \alpha \sqrt{\widehat{\tau}_r^t})$ , where  $\alpha$  has a one-to-one map to  $\lambda$ .

## 2) Bayesian MMSE

- Goal: compute/approximate MMSE estimate  $\hat{\mathbf{x}} = \mathbb{E}\{\mathbf{x}|\mathbf{y}\}$ .
- Suppose  $\mathbf{x}_o \sim$  i.i.d.  $p(x_j)$  with known  $p(x_j)$ .
- Use  $[\mathbf{g}^t(\mathbf{r})]_j = \mathbb{E}\{x_j | r_j = x_{o,j} + \mathcal{N}(0, \widehat{\tau}_r^t)\}$  ... scalar denoising!
- MMSE is achieved when the SE has a unique fixed point!

*The choice of denoiser determines the problem solved by AMP.*

## Choice of denoiser in AMP (cont.)

### 3) Non-parametric (or model free) estimation

- Goal: compute MMSE estimate without knowing i.i.d. prior  $p(x_j)$ .
- Assume scalar GMM( $\theta$ ) with unknown parameters  $\theta$ .
- Use MMSE scalar estimator for GMM( $\theta^t$ ) at iteration  $t$ .
- Use EM algorithm to update  $\theta^t$ . Details given later...

### 4) Black-Box Denoisers<sup>5</sup>

- Goal: leverage sophisticated off-the-shelf denoisers like BM3D for natural images or BM4D for image sequences.

- Use  $\mathbf{g}^t(\mathbf{r}) = \text{BM3D}(\mathbf{r}; \tau_r^t)$ .

- Approximate divergence as  $\langle \mathbf{g}^{t'}(\mathbf{r}) \rangle \approx \frac{1}{N} \sum_{j=1}^N \frac{g_j^t(\mathbf{r} + \epsilon \mathbf{s}) - s_j g_j^t(\mathbf{r})}{\epsilon}$   
where  $\{s_j\} \sim$  i.i.d. uniform  $\pm 1$ .

---

<sup>5</sup> Metzler, Maleki, Baraniuk-TIT'16

# The limitations of AMP

The good:

- For large i.i.d. **sub-Gaussian**  $\mathbf{A}$ , AMP performs provably well.<sup>6</sup>
- Finite-sample analysis shows mild degradation with **not-so-large** i.i.d. Gaussian  $\mathbf{A}$ .<sup>7</sup>
- **Empirical evidence** shows good performance in some other cases (e.g., randomly sub-sampled Fourier  $\mathbf{A}$  & i.i.d. sparse  $x$ )

The bad:

- **For general  $\mathbf{A}$ , AMP can perform poorly**

The ugly:

- **For general  $\mathbf{A}$ , AMP may fail to converge!**
  - ill-conditioned  $\mathbf{A}$
  - non-zero mean  $\mathbf{A}$

---

<sup>6</sup> Bayati, Lelarge, Montanari—AAP'15

<sup>7</sup> Rush, Venkataraman—ISIT'16

# This talk: Vector AMP

For SLR  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ , the **vector AMP** algorithm is<sup>8</sup>

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

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

$$\alpha_1^t = \langle \mathbf{g}'(\mathbf{r}_1^t; \gamma_1^t) \rangle \quad \text{divergence}$$

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

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

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

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

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

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

Note similarities with standard AMP.

<sup>8</sup>Rangan, Schniter, Fletcher—arXiv:1610.03082.

## Vector AMP without matrix inverses

Can avoid matrix inverses using an “economy” SVD  $\mathbf{A} = \mathbf{USV}^T$ :

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

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

$$\alpha_1^t = \langle \mathbf{g}'(\mathbf{r}_1^t; \gamma_1^t) \rangle \quad \text{divergence}$$

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

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

---

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

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

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

Note economy SVD computable with  $O(M^3 + M^2N)$  operations.

## 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 can be rigorously characterized by a [state-evolution](#) in the high-dimensional limit of certain random  $\mathbf{A}$ :

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

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

“ $\mathbf{A}$  is [right-rotationally invariant](#)”

# Message-passing derivation of VAMP

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

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

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

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

$\delta(\mathbf{x}_1 - \mathbf{x}_2)$

- Perform<sup>9</sup> message-passing with messages approximated as  $\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ .  
An instance of *expectation-propagation*<sup>10</sup> (EP).

<sup>9</sup>Rangan, Schniter, Fletcher—arXiv:1610.03082.

<sup>10</sup>Minka—Dissertation'01

# Free-energy derivation of VAMP

- Want to compute **posterior density**:

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x})\ell(\mathbf{x})}{Z} \text{ with } \begin{cases} p(\mathbf{x}) = \text{prior} \\ \ell(\mathbf{x}) = N(\mathbf{y}; \mathbf{A}\mathbf{x}, \tau_w \mathbf{I}), \text{ likelihood} \\ Z = \int p(\mathbf{x})\ell(\mathbf{x})d\mathbf{x}, \text{ partition fxn} \end{cases}$$

but difficult due to high-dimensional integral.

- What if we compute the density via

$$\arg \min_{b(\mathbf{x})} D(b(\mathbf{x})||p(\mathbf{x}|\mathbf{y}))$$

where the **KL divergence** can be written as

$$D(b||p) = \underbrace{D(b||p) + D(b||\ell) + H(b)}_{\text{Gibbs free energy}} + \text{const},$$

thus avoiding the partition function  $Z$ . Still difficult...

## Free-energy derivation of VAMP (cont.)

- What about **splitting** the belief  $b(\mathbf{x})$ :

$$\arg \min_{b_1, b_2} \max_q J(b_1, b_2, q) \text{ s.t. } b_1 = b_2 = q$$

$$J(b_1, b_2, q) = D(b_1 \| p) + D(b_2 \| \ell) + H(q)$$

noting that  $D(\cdot \| p)$  is convex and  $H(\cdot)$  is concave?

Still difficult due to the pdf constraint. . .

- So, relax the pdf constraint to moment-matching constraints:

$$b_1 = b_2 = q \longrightarrow \begin{cases} \mathbb{E}\{\mathbf{x}|b_1\} = \mathbb{E}\{\mathbf{x}|b_2\} = \mathbb{E}\{\mathbf{x}|q\} \\ \text{Tr}[\text{Cov}\{\mathbf{x}|b_1\}] = \text{Tr}[\text{Cov}\{\mathbf{x}|b_2\}] = \text{Tr}[\text{Cov}\{\mathbf{x}|q\}] \end{cases}$$

An instance of **expectation-consistent approximation**<sup>11</sup> (EC).

---

<sup>11</sup>Opper, Winther–NIPS'04, Fletcher, Rangan, Schniter–ISIT'16

## Free-energy derivation of VAMP (cont.)

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

$$b_1(\mathbf{x}) \propto p(\mathbf{x})\mathcal{N}(\mathbf{x}; \mathbf{r}_1; \mathbf{I}/\gamma_1)$$

$$b_2(\mathbf{x}) \propto \ell(\mathbf{x})\mathcal{N}(\mathbf{x}; \mathbf{r}_2; \mathbf{I}/\gamma_2)$$

$$q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}; \mathbf{I}/\eta)$$

for parameters  $\mathbf{r}_1, \gamma_1, \mathbf{r}_2, \gamma_2, \hat{\mathbf{x}}, \eta$  that satisfy

$$\hat{\mathbf{x}} = \mathbb{E}\{\mathbf{x}|b_1\} = \mathbb{E}\{\mathbf{x}|b_2\} = \mathbb{E}\{\mathbf{x}|q\}$$

$$1/\eta = \frac{1}{N} \text{Tr}[\text{Cov}\{\mathbf{x}|b_1\}] = \frac{1}{N} \text{Tr}[\text{Cov}\{\mathbf{x}|b_2\}] = \frac{1}{N} \text{Tr}[\text{Cov}\{\mathbf{x}|q\}].$$

- Can then construct algorithms whose **fixed points coincide with these stationary points** (e.g., EC, ADATAP,<sup>12</sup> S-AMP<sup>13</sup>). But **convergence** is not guaranteed.

---

<sup>12</sup>Opper, Winther–NC'00

<sup>13</sup>Cacmak, Winter, Fleury–ITW'14

# Putting things in perspective

- The aforementioned belief-propagation and free-energy derivations are both **well known** and **heuristic** (in general).
  - The resulting algorithms **may not converge** to their fixed points
    - S-AMP diverges with mildly ill-conditioned matrices
  - Even if they do converge, the **accuracy of the fixed points is unclear**:
    - EP generally suboptimal due to approximation of messages
    - EC generally suboptimal due to approximation of constraint
- The important question is whether/when a given heuristic can be **rigorously analyzed** and proven to work well.

*AMP rigorous analyzed under large **i.i.d. Gaussian  $A$**  and Bayes optimal under certain combinations of  $\{p(\mathbf{x}), \ell(\mathbf{x})\}$ .*

# VAMP state evolution

VAMP has a rigorous SE<sup>14</sup>

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 VAMP-SE under  $\hat{\tau}_w = \tau_w$  is as follows:

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

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

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

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

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

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

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

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

<sup>14</sup>Rangan, Schniter, Fletcher—arXiv:1610.03082

## Connections to the replica prediction

- The **replica method** from statistical physics is often used to characterize the average behavior of large disordered systems.
- Although **not fully rigorous**, replica predictions are usually correct.
- For **SLR** under large right-rotationally invariant **A** and matched priors, *The MMSE  $\mathcal{E}_1(\bar{\gamma}_1)$  should satisfy the fixed-point equation<sup>15</sup>*

$$\bar{\gamma}_1 = R_{\mathbf{A}^\top \mathbf{A} / \tau_w}(-\mathcal{E}_1(\bar{\gamma}_1)),$$

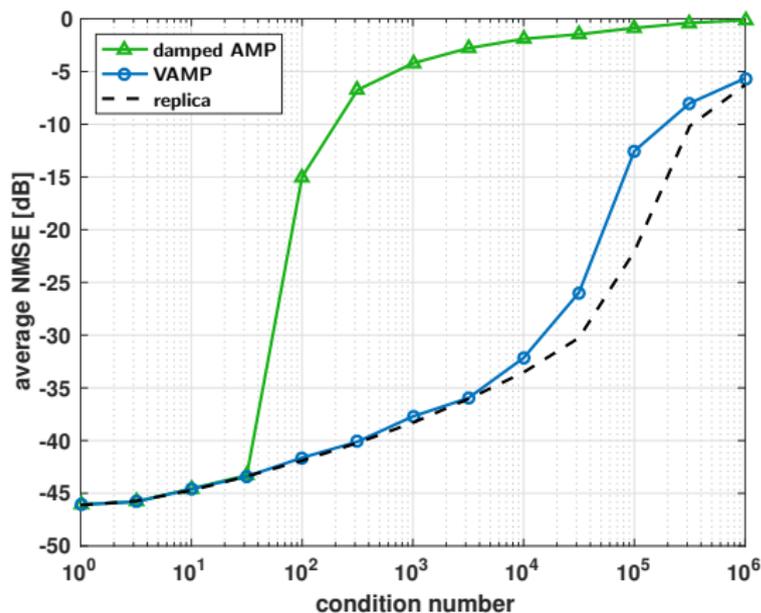
where  $R_{\mathbf{C}}(\cdot)$  denotes the *R-transform* of matrix  $\mathbf{C}$  and  $\mathcal{E}_1(\bar{\gamma}_1) \triangleq \mathbb{E} \left\{ [g_{mmse}(X_o + \mathcal{N}(0, 1/\bar{\gamma}_1); \bar{\gamma}_1) - X_o]^2 \right\}$ .

- It can be shown that **VAMP's matched SE** obeys the above equation.
- Thus, if the replica prediction is correct, then VAMP's estimates **will be MMSE** whenever the replica fixed-point equation has a unique solution.

---

<sup>15</sup>Tulino, Caire, Verdu, Shamai-TIT'13

# Experiment with Matched Priors I



$$N = 1024$$

$$M/N = 0.5$$

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

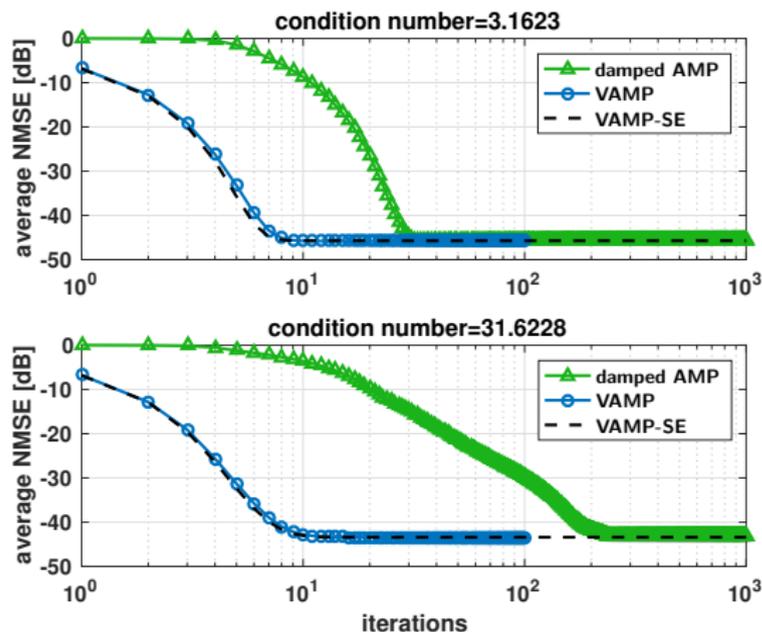
$\mathbf{U}, \mathbf{V}$  drawn uniform  
 $s_n/s_{n-1} = \phi \quad \forall n$   
 $\phi$  determines  $\kappa(\mathbf{A})$

$$X_o \sim \text{Bernoulli-Gaussian}$$
$$\Pr\{X_0 \neq 0\} = 0.1$$

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

Note robustness w.r.t. condition number of  $\mathbf{A}$ .

# Experiment with Matched Priors II



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

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

$X_o \sim \text{Bernoulli-Gaussian}$   
 $\Pr\{X_0 \neq 0\} = 0.1$

SNR = 40dB

Note convergence speed relative to (damped) EM-AMP.

# Non-parametric (model-free) regression

- So far we considered recovering  $\mathbf{x}_o$  from

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

when  $p(\mathbf{x})$  and  $\tau_w$  are **known**.

- Can we **learn**  $\tau_w$ ? Yes, through an EM procedure.<sup>16</sup>

Can we **learn**  $p(\mathbf{x})$ ? Yes if  $p(\mathbf{x}) = \prod_j p(x_j)$ .

- Why is  $p(x_j)$  learnable with VAMP?

- Recall that  $\mathbf{r}_1^t = \mathbf{x}_o + \mathcal{N}(\mathbf{0}, \tau_1^t \mathbf{I})$ .

- Thus  $\mathbf{r}_1^t$  contains i.i.d. samples of  $p(x_j) * \mathcal{N}(x_j; 0, \tau_1^t)$ .

- Should be able to deconvolve  $p(x_j)$  from the empirical distribution of  $\mathbf{r}_1^t$ .

- A practical method: Model  $p(x_j) = \text{GMM}(x_j; \boldsymbol{\theta}_x)$ .  
Learn parameters  $\boldsymbol{\theta}_x$  using EM.

---

<sup>16</sup>Fletcher, Schniter—arXiv:1602.08207

# EM-VAMP

- Recall  $\begin{cases} \text{prior } p(\mathbf{x}; \boldsymbol{\theta}_x) \\ \text{likelihood } \ell(\mathbf{x}; \tau_w) \end{cases} \rightarrow \text{Learn parameters } \boldsymbol{\theta} \triangleq (\boldsymbol{\theta}_x, \tau_w).$

- EM: iterate

$$Q(\boldsymbol{\theta}; \hat{\boldsymbol{\theta}}^k) = \int p(\mathbf{x}|\mathbf{y}; \hat{\boldsymbol{\theta}}^k) \ln p(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) d\mathbf{x} \quad \text{“expectation”}$$

$$\hat{\boldsymbol{\theta}}^{k+1} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \hat{\boldsymbol{\theta}}^k) \quad \text{“maximization”}$$

which uses the posterior  $p(\mathbf{x}|\mathbf{y}; \hat{\boldsymbol{\theta}}^k)$  in the E step.

- With VAMP's posterior approx, EM is an alternating approach to

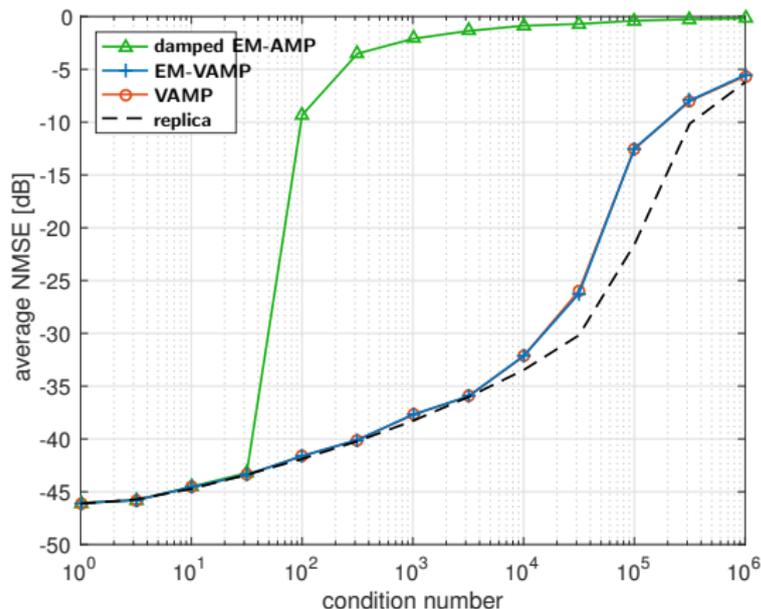
$$\min_{b_1, b_2, \boldsymbol{\theta}} \max_q D(b_1 \| p(\boldsymbol{\theta}_x)) + D(b_2 \| \ell(\tau_w)) + H(q)$$

$$\text{s.t. } \begin{cases} \mathbb{E}\{\mathbf{x}|b_1\} = \mathbb{E}\{\mathbf{x}|b_2\} = \mathbb{E}\{\mathbf{x}|q\} \\ \text{Tr}[\text{Cov}\{\mathbf{x}|b_1\}] = \text{Tr}[\text{Cov}\{\mathbf{x}|b_2\}] = \text{Tr}[\text{Cov}\{\mathbf{x}|q\}] \end{cases}$$

- Can make faster by putting  $\boldsymbol{\theta}$  optimization in the inner loop.

# Experiment with Learned Parameters I

Learning both  $\tau_w$  and  $\theta_x$ :



$N = 1024$

$M/N = 0.5$

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

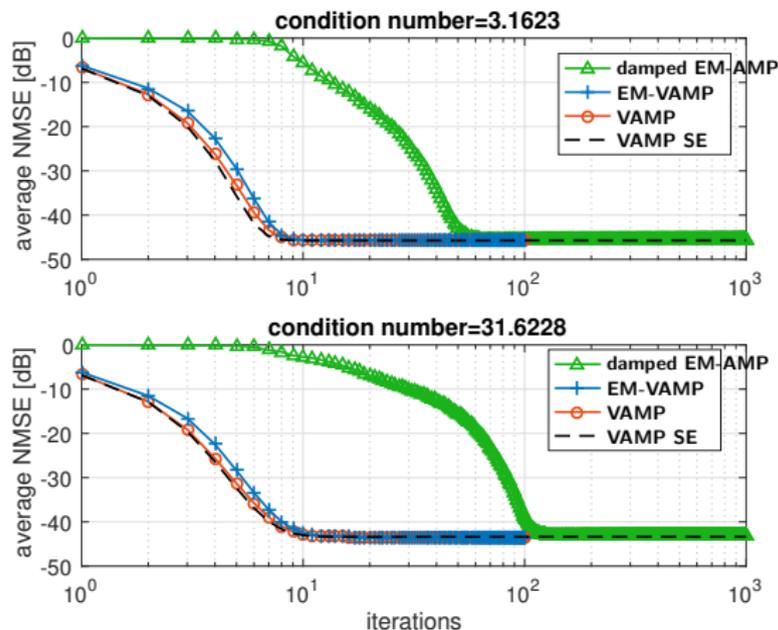
$X_o \sim \text{Bernoulli-Gaussian}$   
 $\Pr\{X_o \neq 0\} = 0.1$

SNR = 40dB

EM-VAMP achieves oracle performance at all condition numbers.

# Experiment with Learned Parameters II

Learning both  $\tau_w$  and  $\theta_x$ :



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

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

$X_o \sim \text{Bernoulli-Gaussian}$   
 $\Pr\{X_0 \neq 0\} = 0.1$

SNR = 40dB

EM-VAMP nearly as fast as VAMP and much faster than EM-AMP.

# Noiseless Image Recovery with BM3D

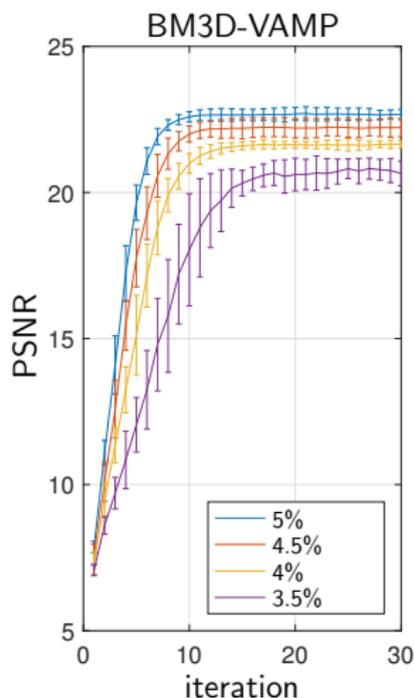
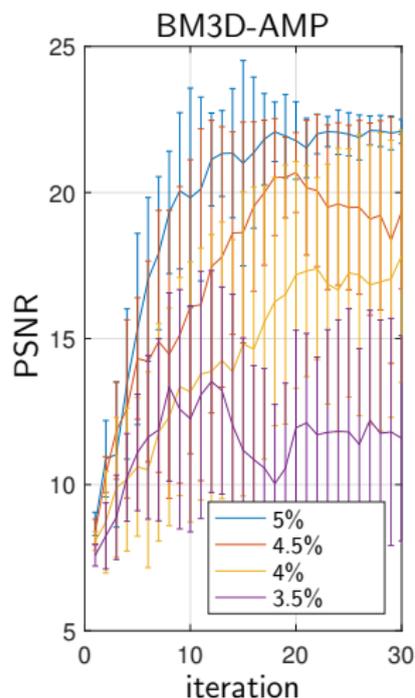
		PSNR	time
10%	L1-AMP	17.7 dB	0.5 s
	L1-VAMP	17.6 dB	0.5 s
	BM3D-AMP	<b>25.2 dB</b>	<b>10.1 s</b>
	BM3D-VAMP	<b>25.2 dB</b>	10.4 s
20%	L1-AMP	20.2 dB	1.0 s
	L1-VAMP	20.2 dB	0.9 s
	BM3D-AMP	<b>30.0 dB</b>	8.8 s
	BM3D-VAMP	<b>30.0 dB</b>	<b>8.5 s</b>
30%	L1-AMP	22.4 dB	1.6 s
	L1-VAMP	22.4 dB	1.4 s
	BM3D-AMP	<b>32.5 dB</b>	8.6 s
	BM3D-VAMP	<b>32.5 dB</b>	<b>8.2 s</b>
40%	L1-AMP	24.6 dB	2.3 s
	L1-VAMP	24.8 dB	1.8 s
	BM3D-AMP	35.1 dB	9.1 s
	BM3D-VAMP	<b>35.2 dB</b>	<b>8.5 s</b>
50%	L1-AMP	27.0 dB	3.1 s
	L1-VAMP	27.2 dB	2.3 s
	BM3D-AMP	37.4 dB	9.8 s
	BM3D-VAMP	<b>37.7 dB</b>	<b>8.8 s</b>

Avg results for recovering 128x128 lena, barbara, boat, fingerprint, house, and peppers from  $\mathbf{y} = \mathbf{A}\mathbf{x}_o$  with i.i.d. Gaussian  $\mathbf{A}$  at various sampling ratios.

All algorithms use 20 iterations and learn the noise variance  $\tau_w$ .

VAMP slightly outperforms AMP in accuracy and runtime.

## Noiseless Image Recovery with BM3D (cont.)



Now look a sampling rates  $\leq 5\%$ .

Goal: recover  $128 \times 128$  *lena* from  $\mathbf{y} = \mathbf{A}\mathbf{x}_o$  with i.i.d. Gaussian  $\mathbf{A}$  and unknown  $\tau_w$ .

**BM3D-VAMP** does much better than **BM3D-AMP**.

# Generalized linear models

- Until now we have considered SLR,  $\mathbf{y} = \mathbf{A}\mathbf{x}_o + \mathbf{w}$ .
- VAMP can also support the **generalized linear model** (GLM)

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

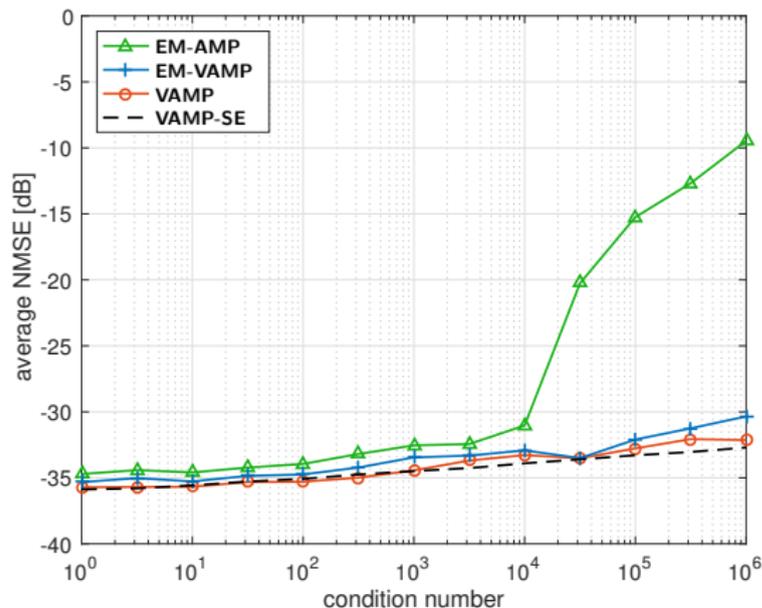
which supports, e.g.,

- $y_i = z_i + w_i$ : additive, possibly **non-Gaussian noise**
- $y_i = \text{sgn}(z_i + w_i)$ : **binary classification / one-bit sensing**
- $y_i = |z_i + w_i|$ : **phase retrieval in noise**
- Poisson  $y_i$ : **photon-limited imaging**

- Trick:  $\mathbf{z} = \mathbf{A}\mathbf{x} \Leftrightarrow \underbrace{\mathbf{0}}_{\tilde{\mathbf{z}}} = \underbrace{[\mathbf{A} - \mathbf{I}]}_{\tilde{\mathbf{A}}} \underbrace{\begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix}}_{\tilde{\mathbf{x}}}$

# One-bit compressed sensing / Probit regression

Learning both  $\tau_w$  and  $\theta_x$ :



$N = 512$

$M/N = 4$

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

$\mathbf{U}, \mathbf{V}$  drawn uniform

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

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

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

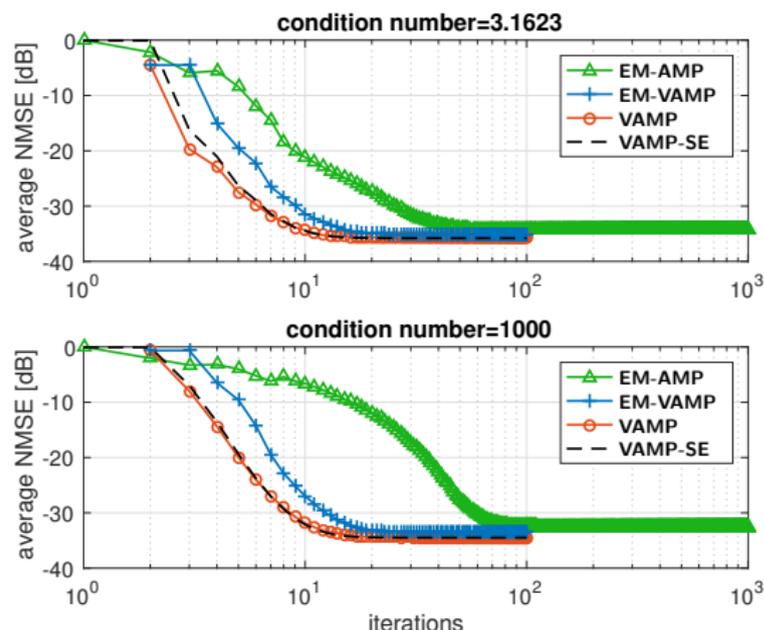
$\Pr\{X_o \neq 0\} = 1/32$

SNR = 40dB

VAMP and EM-VAMP robust to ill-conditioned  $\mathbf{A}$ .

# One-bit compressed sensing / Probit regression

Learning both  $\tau_w$  and  $\theta_x$ :



$N = 512$

$M/N = 4$

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

$\mathbf{U}, \mathbf{V}$  drawn uniform

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

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

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

$$\Pr\{X_0 \neq 0\} = 1/32$$

SNR = 40dB

EM-VAMP mildly slower than VAMP but much faster than damped AMP.

## Conclusions

AMP exhibits some remarkable properties

- low cost-per-iteration and relatively few iterations to convergence,
- intermediate estimates of form  $\mathbf{r}^t = \mathbf{x}_o + \mathcal{N}(\mathbf{0}, \tau_r^t \mathbf{I})$ ,
- rigorous state evolution,
- easy tuning of prior & likelihood,
- compatibility with plug-in denoisers like BM3D,

but those properties are guaranteed only under large i.i.d. Gaussian  $\mathbf{A}$ .

Vector AMP has the same properties, but for a much larger class of  $\mathbf{A}$ .

Ongoing work: analysis of EM procedure, bilinear extensions, connections with deep learning, various applications. . .

*Thanks for listening!*