

# Bilinear Recovery using EM Vector-AMP

Phil Schniter and Subrata Sarkar



THE OHIO STATE UNIVERSITY

Collaborators: **Sundeep Rangan** @ NYU and **Allie Fletcher** @ UCLA

With support from: NSF CCF-1527162 and NSF CCF-1716388

ITA Workshop (San Diego) — Feb 12, 2018

# Bilinear Recovery

Goal: Recover signal  $\mathbf{X}$  and parameters  $\theta_A$  from noisy measurements

$$\mathbf{Y} = \mathbf{A}(\theta_A)\mathbf{X} + \mathbf{W} \quad \text{with affine linear map } \mathbf{A}(\cdot).$$

Applications:

- Self calibration<sup>1</sup>
- Compressed sensing with matrix uncertainty<sup>2</sup>
- Blind deconvolution<sup>3</sup>
- Dictionary learning<sup>4</sup>
- Joint channel estimation and symbol detection

---

<sup>1</sup>Ling,Strohmer'15   <sup>2</sup>Zhu,Leus,Giannakis'11   <sup>3</sup>Ahmed,Recht,Romberg'12   <sup>4</sup>Aharon,Elad,Bruckstein'06

# Statistical Model

Measurements:  $\mathbf{Y} = \mathbf{A}(\boldsymbol{\theta}_A)\mathbf{X} + \mathbf{W}$

Assumptions:

- $\mathbf{A}(\cdot) : \mathbb{R}^Q \rightarrow \mathbb{R}^{M \times N}$  measurement operator
- $\mathbf{X} \in \mathbb{R}^{N \times L}$  with  $x_{nl} \stackrel{\text{i.i.d.}}{\sim} p_X(\cdot; \boldsymbol{\theta}_X)$  random signal
- $\mathbf{W} \in \mathbb{R}^{M \times L}$  with  $w_{ml} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \theta_w^{-1})$  AWGN
- $\boldsymbol{\Theta} \triangleq \{\boldsymbol{\theta}_A, \boldsymbol{\theta}_X, \theta_w\}$  unknown deterministic parameters

Goal: compute...

- $\hat{\boldsymbol{\Theta}}_{\text{ML}} = \arg \max_{\boldsymbol{\Theta}} p(\mathbf{Y}; \boldsymbol{\Theta})$  maximum likelihood
- $\hat{\mathbf{X}}_{\text{MMSE}} = \mathbb{E}[\mathbf{X} | \mathbf{Y}; \hat{\boldsymbol{\Theta}}_{\text{ML}}]$  “empirical Bayes”

## Related Work: $\Theta = \{\theta_A, \theta_X, \theta_w\}$ Known

Consider the case where  $\Theta \triangleq \{\theta_A, \theta_X, \theta_w\}$  (and thus  $\mathbf{A}$ ,  $p_X$ ,  $p_W$ ) are **known**.

- The “**vector AMP**” (VAMP) algorithm<sup>5</sup> can be applied.
- When  $\mathbf{A}$  is a **large right-rotationally invariant** random matrix, the macroscopic behavior of VAMP is rigorously characterized by a scalar **state-evolution**.<sup>5</sup>
- When the state-evolution has a unique fixed point, it is “good” in the sense that **VAMP’s MSE agrees with the replica prediction**<sup>6</sup> of the MMSE.
- VAMP is **more robust than AMP**<sup>7</sup>, which requires large i.i.d sub-Gaussian  $\mathbf{A}$ .

<sup>5</sup>Rangan, Schniter, Fletcher’16

<sup>6</sup>Tulino, Caire, Verdú, Shamai’13

<sup>7</sup>Donoho, Maleki, Montanari’10

## Related Work: $\theta_A$ Known, $\{\theta_X, \theta_w\}$ Unknown

Now consider case where  $\theta_A$  (and thus  $\mathbf{A}$ ) is known, but  $\theta_X$  and  $\theta_w$  are not.

- The EM-VAMP algorithm<sup>8</sup> can be applied.
- When  $\mathbf{A}$  is a large right-rotationally invariant random matrix, the macroscopic behavior of EM-VAMP is rigorously characterized by a state-evolution.<sup>4</sup>
- For certain classes of  $p_X$  (e.g., exponential family), EM-VAMP's parameter estimates  $\{\hat{\theta}_X, \hat{\theta}_w\}$  are asymptotically consistent.<sup>4</sup>

---

<sup>8</sup>Fletcher, Rangan, Schniter'17

## Related Work: $\{\boldsymbol{\theta}_A, \boldsymbol{\theta}_X, \theta_w\}$ Unknown

Finally, consider the case of interest, where  $\boldsymbol{\theta}_A$ ,  $\boldsymbol{\theta}_X$  and  $\theta_w$  are **unknown**.

- Suppose we have  $\mathbf{Y} = \mathbf{A}(\boldsymbol{\theta}_A)\mathbf{X} + \mathbf{W}$  with affine linear map  $\mathbf{A}(\cdot)$ .
- The **Parametric Bilinear Generalized AMP** (P-BiG-AMP) algorithm<sup>910</sup> can be used to recover  $\boldsymbol{\theta}_A$  and  $\mathbf{X} \stackrel{\text{i.i.d.}}{\sim} p_x(\cdot; \boldsymbol{\theta}_X)$ .
- But P-BiG-AMP is **based on an i.i.d. Gaussian model** for  $\mathbf{A}(\cdot)$ , and it may not perform well when  $\mathbf{A}(\cdot)$  deviates from that model!

<sup>9</sup>Parker, Schniter'16

<sup>10</sup>Schulke, Schniter, Zdeborova'16

# This Work: $\{\theta_A, \theta_X, \theta_w\}$ Unknown

In the next few slides, we will outline the EM-VAMP methodology:

- 1 Inference via Expectation Consistent Approximation (EC)
- 2 Algorithmic implementation via VAMP
- 3 Learning  $\Theta$  via Expectation Maximization (EM)
- 4 Joint inference & learning via EM-VAMP

# Variational Inference

For the moment, let's suppose that  $\Theta = \{\theta_A, \theta_X, \theta_w\}$  is **known**.

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

$$p(\mathbf{X}|\mathbf{Y}; \Theta) = \frac{p(\mathbf{X}; \Theta)p(\mathbf{Y}|\mathbf{X}; \Theta)}{Z(\Theta)} \quad \text{for } Z(\Theta) \triangleq \int p(\mathbf{X}; \Theta)p(\mathbf{Y}|\mathbf{X}; \Theta) d\mathbf{X},$$

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

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

$$\begin{aligned} p(\mathbf{X}|\mathbf{Y}; \Theta) &= \arg \min_b D(b(\mathbf{X}) \| p(\mathbf{X}|\mathbf{Y}; \Theta)) \quad \text{where } D(\cdot \| \cdot) \text{ is KL divergence} \\ &= \arg \min_b \underbrace{D(b(\mathbf{X}) \| p(\mathbf{X}; \Theta)) + D(b(\mathbf{X}) \| p(\mathbf{Y}|\mathbf{X}; \Theta)) + H(b(\mathbf{X}))}_{\text{Gibbs free energy}} \\ &= \arg \min_{b, \bar{b}, q} \underbrace{D(b(\mathbf{X}) \| p(\mathbf{X}; \Theta)) + D(\bar{b}(\mathbf{X}) \| p(\mathbf{Y}|\mathbf{X}; \Theta)) + H(q(\mathbf{X}))}_{\triangleq J_{\text{Gibbs}}(b, \bar{b}, q; \Theta)} \\ &\quad \text{s.t. } b = \bar{b} = q, \end{aligned}$$

but the density constraint keeps the problem difficult.



# Expectation Consistent Approximation

- In **expectation-consistent approximate inference (EC)**,<sup>11</sup> the density constraint  $b = \bar{b} = q$  is relaxed to moment-matching constraints:

$$p(\mathbf{X}|\mathbf{Y}; \Theta) \approx \arg \min_{b, \bar{b}, q} J_{\text{Gibbs}}(b, \bar{b}, q; \Theta)$$

$$\text{s.t.} \quad \begin{cases} \mathbb{E}\{\mathbf{x}_l|b\} = \mathbb{E}\{\mathbf{x}_l|\bar{b}\} = \mathbb{E}\{\mathbf{x}_l|q\} \quad \forall l \\ \text{tr}(\text{Cov}\{\mathbf{x}_l|b\}) = \text{tr}(\text{Cov}\{\mathbf{x}_l|\bar{b}\}) = \text{tr}(\text{Cov}\{\mathbf{x}_l|q\}) \quad \forall l. \end{cases}$$

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

$$b(\mathbf{X}) \propto \prod_{l=1}^L p(\mathbf{x}_l; \Theta) \mathcal{N}(\mathbf{x}_l; \mathbf{r}_l, \mathbf{I}/\gamma_l)$$

$$\bar{b}(\mathbf{X}) \propto \prod_{l=1}^L p(\mathbf{y}_l|\mathbf{x}_l; \Theta) \mathcal{N}(\mathbf{x}_l; \bar{\mathbf{r}}_l, \mathbf{I}/\bar{\gamma}_l)$$

$$q(\mathbf{X}) = \prod_{l=1}^L \mathcal{N}(\mathbf{x}_l; \hat{\mathbf{x}}_l, \mathbf{I}/\eta_l)$$

with parameters  $\{\mathbf{R}, \gamma, \bar{\mathbf{R}}, \bar{\gamma}, \hat{\mathbf{X}}, \eta\}$  such that

$$\mathbb{E}\{\mathbf{x}_l|b\} = \mathbb{E}\{\mathbf{x}_l|\bar{b}\} = \hat{\mathbf{x}}_l \quad \forall l$$

$$\frac{1}{N} \text{tr}(\text{Cov}\{\mathbf{x}_l|b\}) = \frac{1}{N} \text{tr}(\text{Cov}\{\mathbf{x}_l|\bar{b}\}) = 1/\eta_l \quad \forall l.$$

<sup>11</sup>Opper, Winther'04

# The VAMP Algorithm

An **iterative approach** to finding  $\{\mathbf{R}, \gamma, \bar{\mathbf{R}}, \bar{\gamma}, \hat{\mathbf{X}}, \eta\}$ :

Initialize  $\{\mathbf{R}, \gamma\}$  and select the estimation functions

$$\mathbf{g}(\mathbf{r}_l; \gamma_l) = \mathbb{E}\{\mathbf{x}_l | b; \mathbf{r}_l, \gamma_l\}$$

$$\bar{\mathbf{g}}(\bar{\mathbf{r}}_l; \bar{\gamma}_l) = \mathbb{E}\{\mathbf{x}_l | \bar{b}; \bar{\mathbf{r}}_l, \bar{\gamma}_l\}.$$

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

$$\hat{\mathbf{x}}_l \leftarrow \mathbf{g}(\mathbf{r}_l; \gamma_l), \quad \forall l \quad \text{MMSE estimation}$$

$$\eta_l \leftarrow \gamma_l N / \text{tr} [\partial \mathbf{g}(\mathbf{r}_l; \gamma_l) / \partial \mathbf{r}_l], \quad \forall l$$

$$\bar{\mathbf{r}}_l \leftarrow (\eta_l \hat{\mathbf{x}}_l - \gamma_l \mathbf{r}_l) / (\eta_l - \gamma_l), \quad \forall l \quad \text{pseudo-measurement}$$

$$\bar{\gamma}_l \leftarrow \eta_l - \gamma_l, \quad \forall l$$

$$\bar{\mathbf{x}}_l \leftarrow \bar{\mathbf{g}}(\bar{\mathbf{r}}_l; \bar{\gamma}_l), \quad \forall l \quad \text{LMMSE estimation}$$

$$\bar{\eta}_l \leftarrow \bar{\gamma}_l N / \text{tr} [\partial \bar{\mathbf{g}}(\bar{\mathbf{r}}_l; \bar{\gamma}_l) / \partial \bar{\mathbf{r}}_l], \quad \forall l$$

$$\mathbf{r}_l \leftarrow (\bar{\eta}_l \bar{\mathbf{x}}_l - \bar{\gamma}_l \bar{\mathbf{r}}_l) / (\bar{\eta}_l - \bar{\gamma}_l), \quad \forall l \quad \text{pseudo-prior}$$

$$\gamma_l \leftarrow \bar{\eta}_l - \bar{\gamma}_l, \quad \forall l$$

Note: this specialization of VAMP is equivalent to **expectation propagation (EP)**.

# Expectation Maximization

We now return to the case where  $\Theta = \{\theta_A, \theta_X, \theta_w\}$  is **unknown**.

- The **EM algorithm** is a well-known iterative approach to maximum-likelihood estimation of  $\Theta$ .
- The EM algorithm can be written in terms of the **Gibbs free energy** as<sup>12</sup>

$$\hat{\Theta}^{(t+1)} = \arg \min_{\Theta} \underbrace{D_{KL}(b^{(t)}(\mathbf{X}) || p(\mathbf{X}; \Theta)) + D_{KL}(b^{(t)}(\mathbf{X}) || p(\mathbf{Y}|\mathbf{X}; \Theta)) + H(b^{(t)})}_{= J_{\text{Gibbs}}(b^{(t)}, b^{(t)}, b^{(t)}; \Theta)}$$

using the belief  $b^{(t)} \triangleq p(\mathbf{X}|\mathbf{Y}; \hat{\Theta}^{(t)})$

- Thus EM and VAMP can be **combined** to solve

$$\min_{\Theta} \min_{b, \bar{b}, q} J_{\text{Gibbs}}(b, \bar{b}, q; \Theta) \text{ s.t. } \begin{cases} \mathbb{E}\{\mathbf{x}_l|b\} = \mathbb{E}\{\mathbf{x}_l|\bar{b}\} = \mathbb{E}\{\mathbf{x}_l|q\} \quad \forall l \\ \text{tr}[\text{Cov}\{\mathbf{x}_l|b\}] = \text{tr}[\text{Cov}\{\mathbf{x}_l|\bar{b}\}] = \text{tr}[\text{Cov}\{\mathbf{x}_l|q\}] \quad \forall l. \end{cases}$$

<sup>12</sup>Neal, Hinton'98

# The EM-VAMP Algorithm

Initialize  $\{\mathbf{R}, \gamma, \hat{\Theta}\}$  and select  $\mathbf{g}(\cdot)$  &  $\bar{\mathbf{g}}(\cdot)$  as before.

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

$$\hat{\mathbf{x}}_l \leftarrow \mathbf{g}(\mathbf{r}_l; \gamma_l, \hat{\Theta}), \quad \forall l \quad \text{MMSE estimation}$$

$$\eta_l \leftarrow \gamma_l N / \text{tr} \left[ \partial \mathbf{g}(\mathbf{r}_l; \gamma_l, \hat{\Theta}) / \partial \mathbf{r}_l \right], \quad \forall l$$

$$\bar{\mathbf{r}}_l \leftarrow (\eta_l \hat{\mathbf{x}}_l - \gamma_l \mathbf{r}_l) / (\eta_l - \gamma_l), \quad \forall l \quad \text{pseudo-measurement}$$

$$\bar{\gamma}_l \leftarrow \eta_l - \gamma_l, \quad \forall l$$

$$\hat{\Theta} \leftarrow \arg \max_{\Theta} \mathbb{E} \{ \ln p(\mathbf{Y} | \mathbf{X}; \Theta) | \bar{\mathbf{R}}; \bar{\gamma}, \hat{\Theta} \} \quad \text{EM update}$$

$$\bar{\mathbf{x}}_l \leftarrow \bar{\mathbf{g}}(\bar{\mathbf{r}}_l; \bar{\gamma}_l, \hat{\Theta}), \quad \forall l \quad \text{LMMSE estimation}$$

$$\bar{\eta}_l \leftarrow \bar{\gamma}_l N / \text{tr} \left[ \partial \bar{\mathbf{g}}(\bar{\mathbf{r}}_l; \bar{\gamma}_l, \hat{\Theta}) / \partial \bar{\mathbf{r}}_l \right], \quad \forall l$$

$$\mathbf{r}_l \leftarrow (\bar{\eta}_l \bar{\mathbf{x}}_l - \bar{\gamma}_l \bar{\mathbf{r}}_l) / (\bar{\eta}_l - \bar{\gamma}_l), \quad \forall l \quad \text{pseudo-prior}$$

$$\gamma_l \leftarrow \bar{\eta}_l - \bar{\gamma}_l, \quad \forall l$$

$$\hat{\Theta} \leftarrow \arg \max_{\Theta} \mathbb{E} \{ \ln p(\mathbf{X}; \Theta) | \mathbf{R}; \gamma, \hat{\Theta} \} \quad \text{EM update}$$

## Variance Auto-Tuning

- Problem: The precisions  $\{\gamma, \bar{\gamma}\}$  of the pseudo- $\{\text{measurement}, \text{prior}\}$  are imperfect when  $\hat{\Theta}$  is imperfect.
- Thus, at each iteration, we **estimate these precisions** jointly with the unknown parameters  $\Theta$ . For example, with the prior parameters  $\theta_X$ :

$$\begin{aligned}
 (\hat{\gamma}, \hat{\theta}_X) &\leftarrow \arg \max_{\gamma, \theta_X} p(\mathbf{R}; \gamma, \theta_X) \\
 &\text{under } \mathbf{r}_l = \mathbf{x}_l + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_l), \quad \mathbf{x}_l \sim p_X(\cdot; \theta_X), \quad \forall l.
 \end{aligned}$$

- In practice, inner iterations of EM can be used to solve the above, e.g.,

$$(\hat{\gamma}, \hat{\theta}_X) = \arg \max_{\gamma, \theta_X} \mathbb{E} \left[ \ln p(\mathbf{X}, \mathbf{R}; \gamma, \theta_X) \mid \mathbf{R}; \hat{\gamma}, \hat{\theta}_X \right].$$

- This “**variance auto-tuning**” procedure<sup>13</sup> leads to asymptot. consistent  $\hat{\theta}_X$ .

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

# State-Evolution and Consistency

Suppose that

- $\mathbf{Y} = \mathbf{A}(\boldsymbol{\theta}_A)\mathbf{X} + \mathcal{N}(\mathbf{0}, \mathbf{I}/\theta_w) \in \mathbb{R}^{M \times L}$ ,  $\mathbf{X} \stackrel{\text{i.i.d.}}{\sim} p_X(\cdot; \boldsymbol{\theta}_X)$
- $\mathbf{A}(\boldsymbol{\theta}) = \mathbf{A}_0 + \sum_{q=1}^Q \theta_q \mathbf{A}_q$  with right-rotationally invariant  $\mathbf{A}_q \in \mathbb{R}^{M \times N}$ .

Conjecture: the behavior of the proposed EM-VAMP algorithm is rigorously characterized by a **state-evolution** with  $M = O(N)$ ,  $N \rightarrow \infty$ , and either

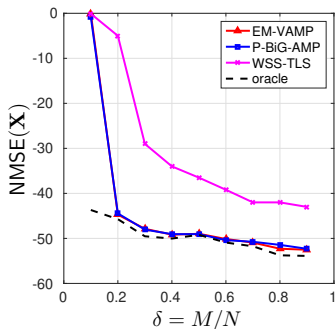
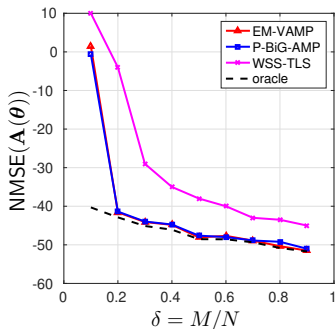
- |                                      |                              |
|--------------------------------------|------------------------------|
| 1 fixed $Q$ and $L = 1$              | (CS with matrix uncertainty) |
| 2 $Q = O(N)$ and $L$ fixed           | (self-calibration)           |
| 3 $Q = O(N^2)$ and $L = O(N \log N)$ | (dictionary learning)        |

Technical conditions include:

- all vectors converge empirically with second-order moments to random variables
- singular values of  $\mathbf{A}_q$  converge empirically with second-order moments to a bounded positive random variable
- Lipschitz  $g(\cdot)$  and  $g'(\cdot)$ .
- exponential-family  $p_X$
- etc...

## Compressed Sensing with Matrix Uncertainty

Problem: Recover 10-sparse  $\mathbf{x} \in \mathbb{R}^N$  from  $\mathbf{y} = (\mathbf{A}_0 + \sum_{q=1}^{10} \theta_q \mathbf{A}_q) \mathbf{x} + \mathbf{w} \in \mathbb{R}^M$ .



EM-VAMP performs similarly to P-BiG-AMP and much better than WSS-TLS.<sup>14</sup>

Details:  $N = 256$ ,  $\mathbf{A}_0 \sim \text{i.i.d. } \mathcal{N}(0, 10)$ ,  $\mathbf{A}_{q>1} \sim \text{i.i.d. } \mathcal{N}(0, 1)$ , SNR=40dB, 10 trials.

<sup>14</sup>Zhu, Leus, Giannakis'11

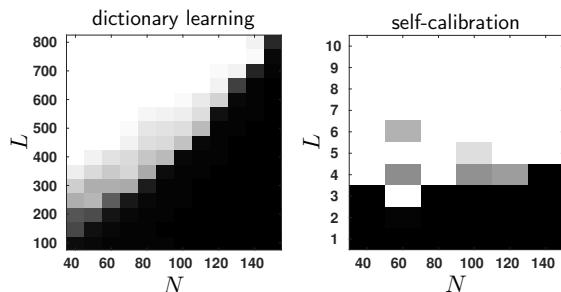
# How many snapshots $L$ are needed?

Problem: Recover  $\mathbf{A} \in \mathbb{R}^{N \times N}$  and sparse  $\mathbf{X} \in \mathbb{R}^{N \times L}$  from  $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{W}$  with

1  $\mathbf{A} \sim$  i.i.d.  $\mathcal{N}(0, 1)$  (dictionary learning)

2  $\mathbf{A} = \sum_{q=1}^N \theta_q \mathbf{A}_q$  with known  $\mathbf{A}_q \sim$  i.i.d.  $\mathcal{N}(0, 1)$  (self-calibration)

NMSE (dB) versus  $N$  and  $L$ :



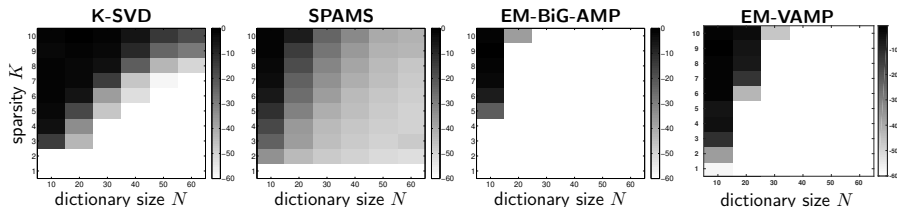
Details: i.i.d. Bernoulli-Gaussian  $\mathbf{X}$ ,  $K = \lceil 0.2N \rceil$ , SNR= 40 dB.



# Dictionary Learning: Sparsity vs Size

Problem: Recover i.i.d. Gaussian  $\mathbf{A} \in \mathbb{R}^{N \times N}$  and Bernoulli-Gaussian  $\mathbf{X} \in \mathbb{R}^{N \times L}$  with  $K$ -sparse columns from  $\mathbf{Y} = \mathbf{A}\mathbf{X}$ .

NMSE (in dB) over 10 realizations for  $L = 5N \log N$ :



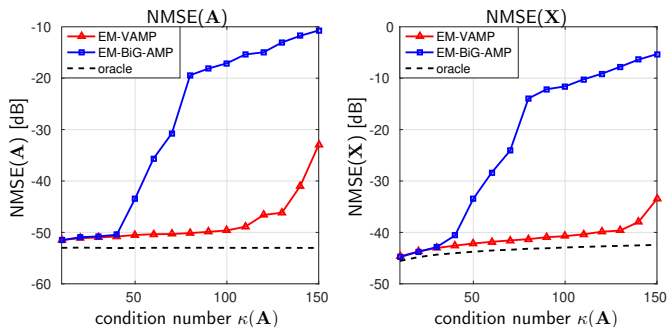
EM-VAMP has slightly worse phase-transition than BiG-AMP, but better than K-SVD<sup>15</sup> and SPAMS<sup>16</sup>.

<sup>15</sup>Aharon, Elad, Bruckstein'06

<sup>16</sup>Mairal, Bach, Ponce, Sapiro'10

# Dictionary Learning: Robustness to Condition Number

Same problem, but now with geometrically spaced singular values in  $\mathbf{A}$ .



EM-VAMP is **more robust than EM-BiG-AMP** to condition number  $\kappa(\mathbf{A})$ .

Details:  $N \times N$  dictionary,  $N = 64$ ,  $K = 13$ -sparse  $\mathbf{x}_l$ ,  $L = 5N \log N$ , SNR=40dB, median NMSE over 100 realizations.

# Conclusions

- We propose a bilinear recovery algorithm, with applications in self-calibration, CS with matrix uncertainty, blind deconvolution, dictionary learning, and joint channel-estimation and symbol detection.
- Broadly speaking, our approach is empirical-Bayesian and uses a combination of EC and EM.
- More specifically, our approach builds on the recently proposed EM-VAMP algorithm, by extending the set of unknown parameters to those that describe the measurement matrix  $\mathbf{A}$ .
- Numerical results suggest performance that is similar to P-BiG-AMP but more robust to non-iid matrices.
- We are currently working on proving the state-evolution conjectures.