

Learning and Free Energies for Vector Approximate Message Passing

Philip Schniter (Ohio State/Duke)
Alyson Fletcher (UCLA)



THE OHIO STATE UNIVERSITY



ICASSP 2017 – New Orleans

Supported by NSF grants 1527162 and 1254204.

Linear Regression with Unknown Prior/Likelihood

Consider the following linear regression problem:

- Observations: $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ with $\begin{cases} \mathbf{x} : \text{unknown signal} \\ \mathbf{A} : \text{known linear operator in } \mathbb{R}^{M \times N} \\ \mathbf{w} : \text{white Gaussian noise.} \end{cases}$
- Prior: $p(\mathbf{x}; \boldsymbol{\theta}_1)$ with deterministic unknown parameters $\boldsymbol{\theta}_1$.
- Likelihood: $\ell(\mathbf{x}; \boldsymbol{\theta}_2) = \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \boldsymbol{\theta}_2\mathbf{I})$ with deterministic unknown variance $\boldsymbol{\theta}_2$.

Goal: jointly infer \mathbf{x} and estimate $\boldsymbol{\theta} \triangleq [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2]$.

Approach: combine variational inference with ML estimation.

Variational Inference

- For now, let's suppose that $\boldsymbol{\theta}$ is known.

- We would like to compute the posterior density

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x}; \boldsymbol{\theta}_1)\ell(\mathbf{x}; \boldsymbol{\theta}_2)}{Z(\boldsymbol{\theta})} \text{ for } Z(\boldsymbol{\theta}) \triangleq \int p(\mathbf{x}; \boldsymbol{\theta}_1)\ell(\mathbf{x}; \boldsymbol{\theta}_2) d\mathbf{x},$$

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

- We can avoid computing $Z(\boldsymbol{\theta})$ through variational optimization:

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \arg \min_b D(b(\mathbf{x})||p(\mathbf{x}|\mathbf{y})) \text{ where } D(\cdot||\cdot) \text{ is KL divergence} \\ &= \arg \min_b \underbrace{D(b(\mathbf{x})||p(\mathbf{x}; \boldsymbol{\theta}_1)) + D(b(\mathbf{x})||\ell(\mathbf{x}; \boldsymbol{\theta}_2)) + H(b(\mathbf{x}))}_{\text{Gibbs free energy}} \\ &= \arg \min_{b_1, b_2} \max_q \underbrace{D(b_1(\mathbf{x})||p(\mathbf{x}; \boldsymbol{\theta}_1)) + D(b_2(\mathbf{x})||\ell(\mathbf{x}; \boldsymbol{\theta}_2)) + H(q(\mathbf{x}))}_{\triangleq J(b_1, b_2, q; \boldsymbol{\theta})} \end{aligned}$$

such that $b_1 = b_2 = q$,

but the density constraint keeps the problem difficult.

- Expectation consistent approximation (EC) [1] relaxes the density constraint to moment-matching constraints:

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &\approx \arg \min_{b_1, b_2} \max_q J(b_1, b_2, q; \boldsymbol{\theta}) \\ \text{such that } &\begin{cases} E\{\mathbf{x}|b_1\} = E\{\mathbf{x}|b_2\} = E\{\mathbf{x}|q\} \\ \text{tr}[\text{Cov}\{\mathbf{x}|b_1\}] = \text{tr}[\text{Cov}\{\mathbf{x}|b_2\}] = \text{tr}[\text{Cov}\{\mathbf{x}|q\}]. \end{cases} \end{aligned}$$

- The stationary points of EC are

$$\begin{aligned} b_1(\mathbf{x}) &\propto p(\mathbf{x}; \boldsymbol{\theta}_1)\mathcal{N}(\mathbf{x}; \mathbf{r}_1, v_1\mathbf{I}) \\ b_2(\mathbf{x}) &\propto \ell(\mathbf{x}; \boldsymbol{\theta}_2)\mathcal{N}(\mathbf{x}; \mathbf{r}_2, v_2\mathbf{I}) \\ q(\mathbf{x}) &= \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \hat{v}\mathbf{I}) \end{aligned} \text{ s.t. } \begin{cases} E\{\mathbf{x}|b_1\} = E\{\mathbf{x}|b_2\} = \hat{\mathbf{x}} \\ \text{tr}[\text{Cov}\{\mathbf{x}|b_1\}] = \text{tr}[\text{Cov}\{\mathbf{x}|b_2\}] = N\hat{v}. \end{cases}$$

Vector AMP (VAMP)

- There exist several algorithms (e.g., EC, ADATAP [2], S-AMP [3]) whose fixed points coincide with the EC stationary points, but often they don't converge.

- An exception is Vector AMP [4], which can be derived using a form of approximation message passing on the vector-valued factor graph

$$p(\mathbf{x}_1; \boldsymbol{\theta}_1) \text{ --- } \mathbf{x}_1 \text{ --- } \mathbf{x}_2 \text{ --- } \ell(\mathbf{x}_2, \boldsymbol{\theta}_2)$$

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

In particular, VAMP is provably convergent under either

- strictly log-concave prior $p(\mathbf{x}; \boldsymbol{\theta}_1)$ and arbitrary \mathbf{A} (after damping),
- iid prior $p(\mathbf{x}; \boldsymbol{\theta}_1)$ and large, right-rotationally invariant \mathbf{A} .

- $\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s})\mathbf{V}^T$ is said to be "right-rotationally invariant" when \mathbf{V} is uniformly distributed on the set of unitary matrices.

- The other SVD quantities, \mathbf{U} and \mathbf{S} , are deterministic and arbitrary.
- This model includes mean-perturbed and ill-conditioned \mathbf{A} , known to break regular AMP.

- With large, right-rotationally invariant \mathbf{A} , VAMP has a rigorous state evolution [4] whose fixed points match the replica prediction of MMSE [5].

Expectation maximization (EM)

- We now return to the problem of estimating $\boldsymbol{\theta}$.

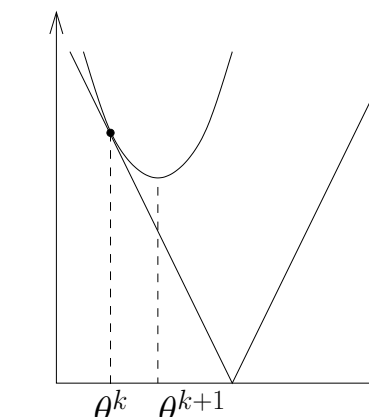
- The maximum-likelihood (ML) estimate is

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} p(\mathbf{y}|\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \{-\ln p(\mathbf{y}|\boldsymbol{\theta})\},$$

which is difficult to compute directly.

- Let's instead consider majorization-minimization: Iteratively minimize a tight upper bound on $-\ln p(\mathbf{y}|\boldsymbol{\theta})$:

$$\begin{aligned} \hat{\boldsymbol{\theta}}^{k+1} &= \arg \min_{\boldsymbol{\theta}} \left\{ -\ln p(\mathbf{y}|\boldsymbol{\theta}) + \underbrace{D(b^k(\mathbf{x})||p(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}))}_{\geq 0} \right\} \\ \text{with } b^k(\mathbf{x}) &= p(\mathbf{x}|\mathbf{y}; \hat{\boldsymbol{\theta}}^k) \end{aligned}$$



- The upper bound " $Q(\boldsymbol{\theta}, b^k)$ " can be rewritten in the form

$$\begin{aligned} Q(\boldsymbol{\theta}, b^k) &\triangleq -\ln p(\mathbf{y}|\boldsymbol{\theta}) + D(b^k(\mathbf{x})||p(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta})) \\ &= -E\{\ln p(\mathbf{y}, \mathbf{x}; \boldsymbol{\theta}) | b^k\} + \text{const}, \end{aligned}$$

which is the usual way of writing the EM algorithm, but it can also be written in terms of the Gibbs free energy

$$\begin{aligned} Q(\boldsymbol{\theta}, b^k) &= D(b^k(\mathbf{x})||p(\mathbf{x}; \boldsymbol{\theta}_1)) + D(b^k(\mathbf{x})||\ell(\mathbf{x}; \boldsymbol{\theta}_2)) + H(b^k(\mathbf{x})) \\ &= J(b^k, b^k, b^k; \boldsymbol{\theta}) \end{aligned}$$

which yields a variational interpretation of EM [6].

The Proposed EM-VAMP Algorithm

- Recall that VAMP iteratively computes a posterior approximation $b^k(\mathbf{x})$ by minimizing $J(b_1, b_2, q; \boldsymbol{\theta})$ (under moment constraints) with known $\boldsymbol{\theta}$.
- Likewise, EM iteratively estimates $\boldsymbol{\theta}$ by minimizing $J(b^k, b^k, b^k; \boldsymbol{\theta})$ assuming the posterior approximation $b^k(\mathbf{x}) = p(\mathbf{x}|\mathbf{y}; \boldsymbol{\theta}^k)$ is available.
- We propose to combine EM and VAMP as follows:

Input \mathbf{g}_1 and \mathbf{g}_2 , and initialize $\mathbf{r}_1 = \mathbf{0}$ and $v_1 = \infty$.

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

$$\hat{\boldsymbol{\theta}}_1 \leftarrow \arg \max_{\boldsymbol{\theta}_1} E\{\ln p(\mathbf{x}; \boldsymbol{\theta}_1) | \mathbf{r}_1, v_1, \hat{\boldsymbol{\theta}}_1\} \quad \text{EM update}$$

$$\hat{v}_1 \leftarrow N^{-1} \text{tr}[\text{Cov}\{\mathbf{x}|\mathbf{r}_1, v_1, \hat{\boldsymbol{\theta}}_1\}] \quad \text{posterior variance}$$

$$\hat{\mathbf{x}}_1 \leftarrow E\{\mathbf{x}|\mathbf{r}_1, v_1, \hat{\boldsymbol{\theta}}_1\} \quad \text{denoising}$$

$$1/v_2 \leftarrow 1/\hat{v}_1 - 1/v_1 \quad \text{variance update}$$

$$\mathbf{r}_2 \leftarrow (\hat{\mathbf{x}}_1/\hat{v}_1 - \mathbf{r}_1/v_1)v_2 \quad \text{Onsager correction}$$

$$\hat{\boldsymbol{\theta}}_2 \leftarrow \arg \max_{\boldsymbol{\theta}_2} E\{\ln \ell(\mathbf{x}; \boldsymbol{\theta}_2) | \mathbf{r}_2, v_2, \hat{\boldsymbol{\theta}}_2\} \quad \text{EM update}$$

$$\hat{v}_2 \leftarrow N^{-1} \text{tr}[\text{Cov}\{\mathbf{x}|\mathbf{r}_2, v_2, \hat{\boldsymbol{\theta}}_2\}] \quad \text{posterior variance}$$

$$\hat{\mathbf{x}}_2 \leftarrow E\{\mathbf{x}|\mathbf{r}_2, v_2, \hat{\boldsymbol{\theta}}_2\} \quad \text{LMMSE estimation}$$

$$1/v_1 \leftarrow 1/\hat{v}_2 - 1/v_2 \quad \text{variance update}$$

$$\mathbf{r}_1 \leftarrow (\hat{\mathbf{x}}_2/\hat{v}_2 - \mathbf{r}_2/v_2)v_1 \quad \text{Onsager correction}$$

where

$$E\{f(\mathbf{x}) | \mathbf{r}_1, v_1, \boldsymbol{\theta}_1\} \triangleq \int f(\mathbf{x}) \frac{p(\mathbf{x}; \boldsymbol{\theta}_1)\mathcal{N}(\mathbf{x}; \mathbf{r}_1, v_1\mathbf{I})}{\int p(\mathbf{x}'; \boldsymbol{\theta}_1)\mathcal{N}(\mathbf{x}'; \mathbf{r}_1, v_1\mathbf{I}) d\mathbf{x}'} d\mathbf{x}$$

$$E\{f(\mathbf{x}) | \mathbf{r}_2, v_2, \boldsymbol{\theta}_2\} \triangleq \int f(\mathbf{x}) \frac{\ell(\mathbf{x}; \boldsymbol{\theta}_2)\mathcal{N}(\mathbf{x}; \mathbf{r}_2, v_2\mathbf{I})}{\int \ell(\mathbf{x}'; \boldsymbol{\theta}_2)\mathcal{N}(\mathbf{x}'; \mathbf{r}_2, v_2\mathbf{I}) d\mathbf{x}'} d\mathbf{x}$$

and similar for the covariances.

- If the SVD $\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s})\mathbf{V}^T$ is precomputed, then

$$\hat{\mathbf{x}}_2 \leftarrow \mathbf{V} (v_2 \text{Diag}(\mathbf{s})^2 + \boldsymbol{\theta}_2\mathbf{I})^{-1} (v_2 \text{Diag}(\mathbf{s})\mathbf{U}^H\mathbf{y} + \boldsymbol{\theta}_2\mathbf{V}^H\mathbf{r}_2)$$

$$\hat{v}_2 \leftarrow \frac{1}{N} \sum_{n=0}^{N-1} \frac{1}{s_n^2/\boldsymbol{\theta}_2 + 1/v_2}, \quad \hat{\boldsymbol{\theta}}_2 \leftarrow \frac{1}{N} \left[\|\mathbf{y} - \mathbf{A}\mathbf{r}_2\|^2 + \sum_n \frac{s_n^2}{s_n^2/\boldsymbol{\theta}_2 + 1/v_2} \right],$$

so EM-VAMP requires only two matrix-vector mults per iteration.

- Other algorithmic variants result when $\boldsymbol{\theta}_1$ and/or $\boldsymbol{\theta}_2$ are updated more or less often than once per VAMP iteration.

Theorem: Fixed Points of EM-VAMP

At any fixed point of EM-VAMP we have

$$\begin{aligned} \hat{v}_1 &= \hat{v}_2 = \frac{v_1 v_2}{v_1 + v_2} \triangleq \hat{v} \\ \hat{\mathbf{x}}_1 &= \hat{\mathbf{x}}_2 = \left(\frac{\mathbf{r}_1}{v_1} + \frac{\mathbf{r}_2}{v_2} \right) \hat{v} \triangleq \hat{\mathbf{x}}. \end{aligned}$$

Also, EM-VAMP's fixed-points are stationary points of the EM-EC optimization

$$\begin{aligned} \min_{\boldsymbol{\theta}} \min_{b_1, b_2} \max_q J(b_1, b_2, q; \boldsymbol{\theta}) \\ \text{such that } &\begin{cases} E\{\mathbf{x}|b_1\} = E\{\mathbf{x}|b_2\} = E\{\mathbf{x}|q\} \\ \text{tr}[\text{Cov}\{\mathbf{x}|b_1\}] = \text{tr}[\text{Cov}\{\mathbf{x}|b_2\}] = \text{tr}[\text{Cov}\{\mathbf{x}|q\}]. \end{cases} \end{aligned}$$

Numerical Experiments

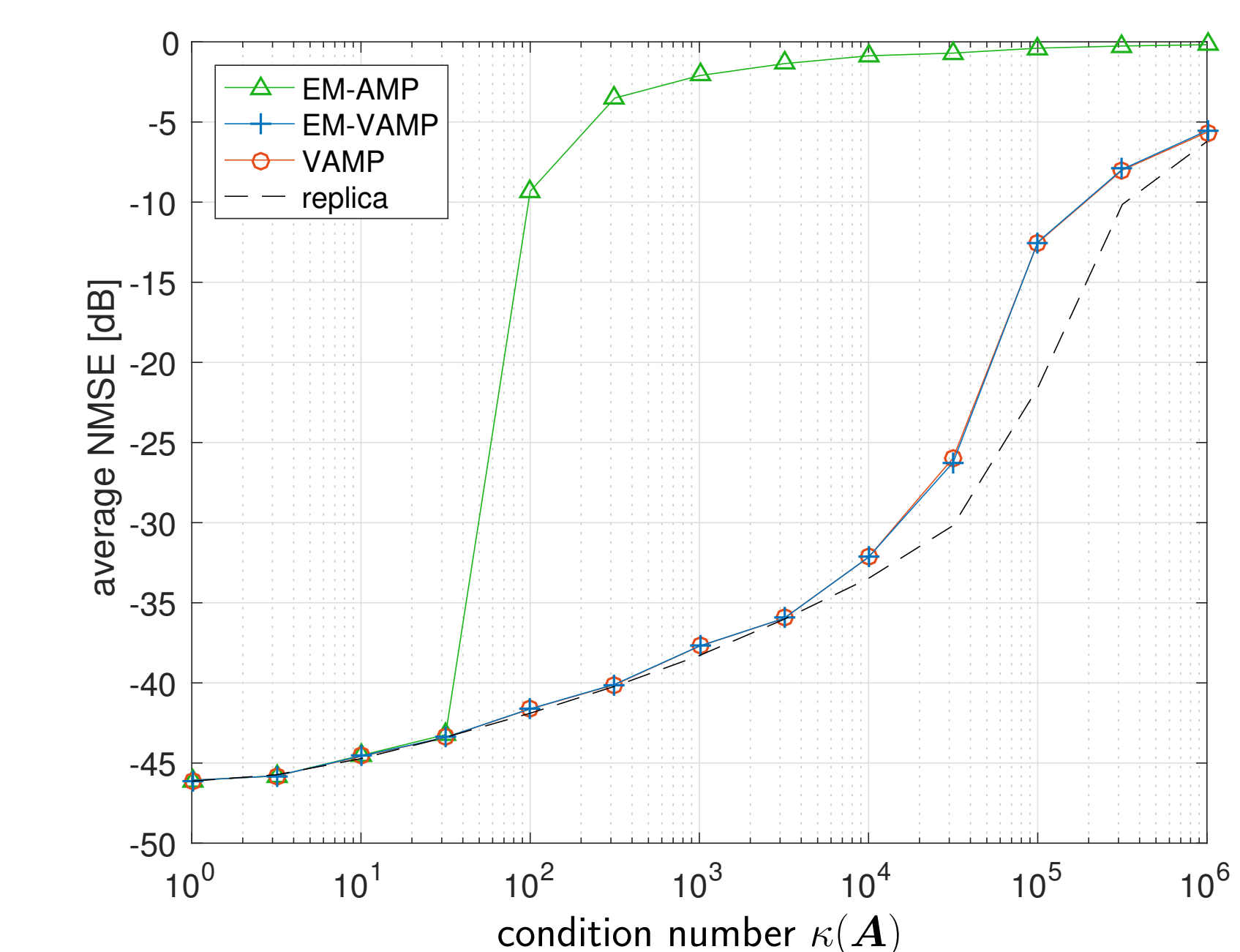
Goal recover $N=1024$ -length i.i.d. Bernoulli-Gaussian \mathbf{x}

$$p(x_n; \boldsymbol{\theta}_1) = (1 - \theta_{11})\delta(x_n) + \theta_{11}\mathcal{N}(x_n; \theta_{12}, \theta_{13}) \text{ with } \boldsymbol{\theta}_1 = [0.1, 0, 1]$$

from $M=512$ measurements

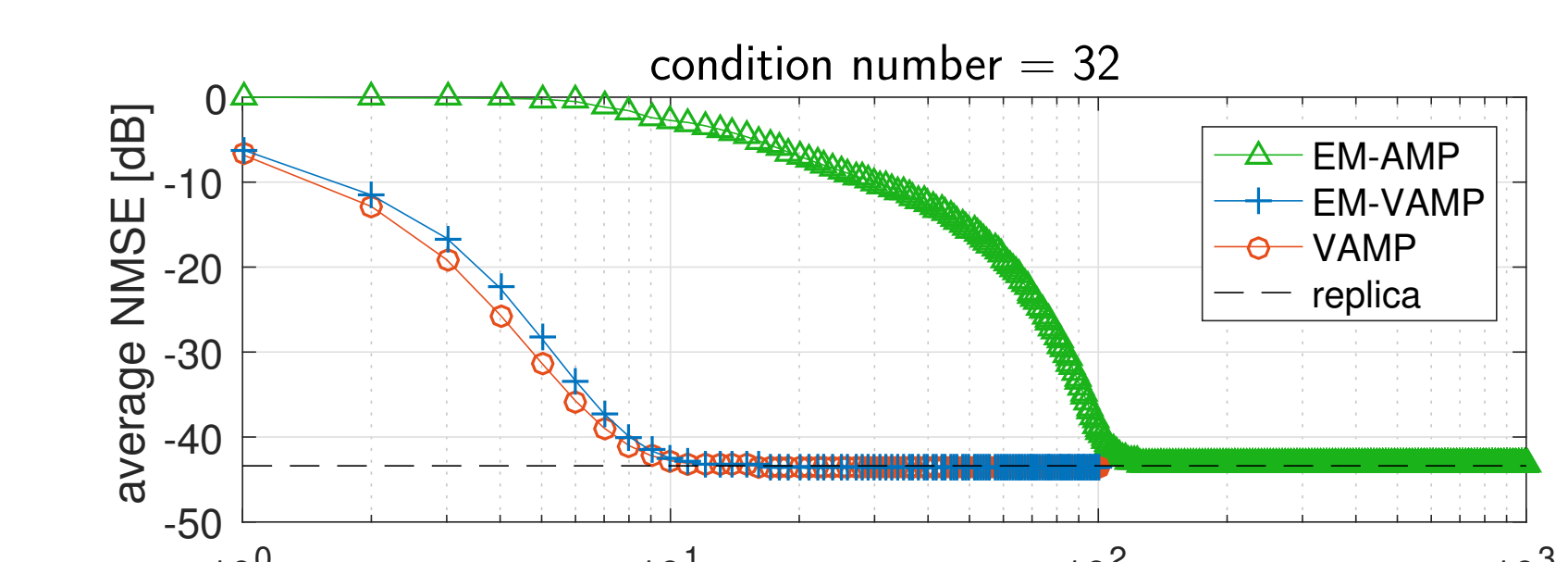
$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathcal{N}(\mathbf{0}, \boldsymbol{\theta}_2\mathbf{I}) \text{ with } \boldsymbol{\theta}_2 \text{ giving SNR}=40 \text{ dB.}$$

Here, $\mathbf{A} = \mathbf{U} \text{Diag}(\mathbf{s})\mathbf{V}^T$ with random orthogonal \mathbf{U}, \mathbf{V} and $s_n/s_{n-1} = \phi \forall n$, where ϕ determines the condition number $\kappa(\mathbf{A})$.

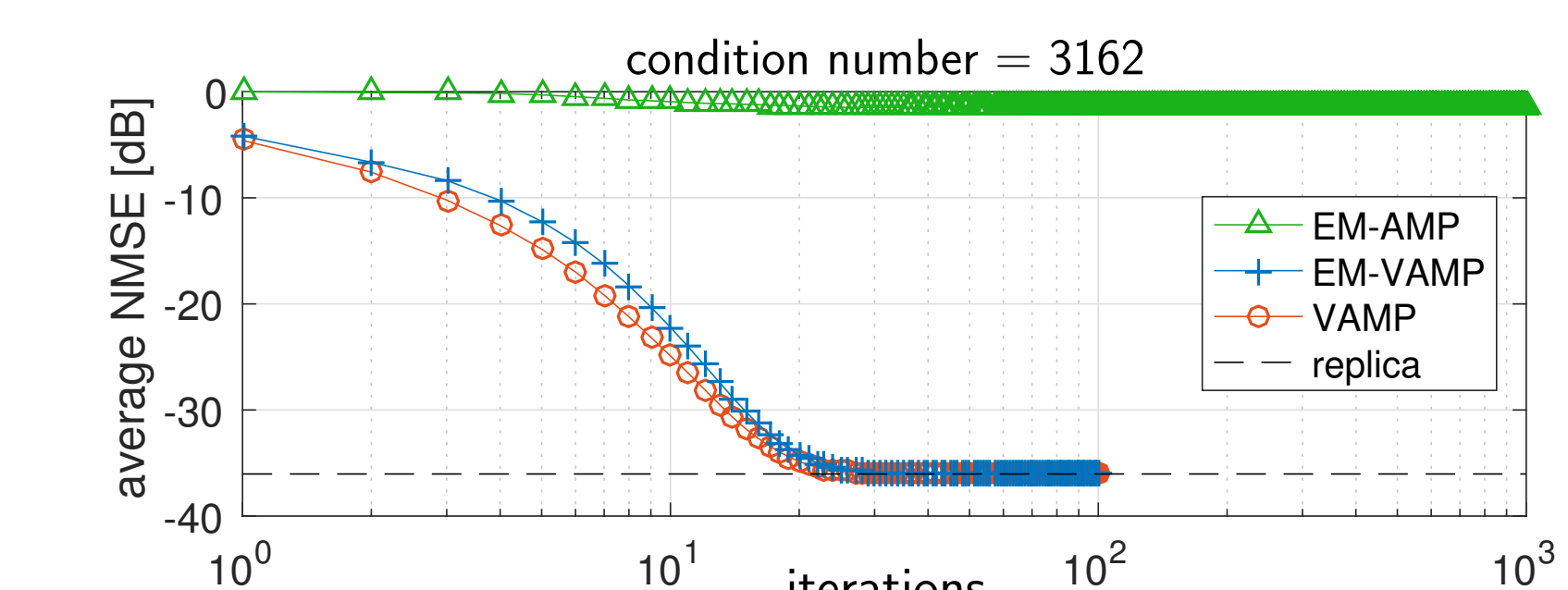


EM-VAMP matches known- $\boldsymbol{\theta}$ -VAMP which matches the replica prediction of MMSE for a wide range of $\kappa(\mathbf{A})$.

EM-AMP [7] only works at small $\kappa(\mathbf{A})$.



EM-VAMP and known- $\boldsymbol{\theta}$ -VAMP take ~ 10 iterations to converge, whereas EM-AMP takes ~ 100 .



EM-VAMP and known- $\boldsymbol{\theta}$ -VAMP take only ~ 25 iterations to converge, while EM-AMP fails.

References

- M. Opper and O. Winther, "Expectation consistent approximate inference," *J. Mach. Learning Res.*, 2005.
- M. Opper and O. Winther, "Adaptive and self-averaging Thouless-Anderson-Palmer mean-field theory for probabilistic modeling," *Phys. Rev. E*, 2001.
- B. Cakmak, O. Winther, and B.H. Fleury, "S-AMP: Approximate Message Passing for General Matrix Ensembles," *ISIT* 2014.
- S. Rangan, P. Schniter, and A. K. Fletcher, "Vector Approximate Message Passing," *arXiv:1610.03082*.
- A. M. Tulino, G. Caire, S. Verdú, and S. Shamai (Shitz), "Support recovery with sparsely sampled free random matrices," *IEEE Trans. Info. Thy.*, 2013.
- R. Neal and G. Hinton, "A view of the EM algorithm that justifies incremental, sparse, and other variants," in *Learning in Graphical Models*, 1998.
- J. P. Vila and P. Schniter, "Expectation-maximization Gaussian-mixture approximate message passing," *IEEE TSP*, 2013.