

Memory AMP and Recent Results on Its Implementation

Brian M. Kurkoski

Joint work with Shunqi Huang and Lei Liu

Japan Advanced Institute of Science and Technology and Zhejiang University

March 13, 2024

BIRS Workshop — Algorithmic Structures for Uncoordinated Communications and
Statistical Inference in Exceedingly Large Spaces

Outline

1. Preliminaries

- ▶ Problem Formulation
- ▶ AMP, OAMP/VAMP, and CAMP
- ▶ MAMP and GD-MAMP

2. Overflow-Avoiding GD-MAMP

- ▶ Overflow Problem of GD-MAMP
- ▶ OA-GD-MAMP with eigenvalues of $\mathbf{A}\mathbf{A}^H$
- ▶ OA-GD-MAMP without eigenvalues of $\mathbf{A}\mathbf{A}^H$

3. Complexity-Reduced GD-MAMP

- ▶ Complexity analysis of GD-MAMP
- ▶ Complexity-Reduced GD-MAMP

4. Conclusion

Problem Formulation

- ▶ System model:

$$\Gamma : \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n},$$

$$\Phi : x_i \sim P_X(x), \forall i.$$

where $\mathbf{A} \in \mathbb{C}^{M \times N}$, $\mathbf{n} \sim \mathcal{CN}(0, \sigma^2 \mathbf{I}_M)$, and $\mathbf{A}, \mathbf{y}, \sigma^2, P_X(\cdot)$ are known.

- ▶ Assumptions:

(1) $M, N \rightarrow \infty$ with fixed $\delta = M/N$.

(2) \mathbf{A} is right-unitarily-invariant.

(3) \mathbf{x} is IID. For convenience, $\mathbb{E}\{\mathbf{x}\} = \mathbf{0}$ and $\frac{1}{N}\mathbb{E}\{\|\mathbf{x}\|^2\} = 1$.

- ▶ Goal: Given $\{\mathbf{y}, \mathbf{A}, \Gamma, \Phi\}$, find an MMSE estimate of \mathbf{x} :

$$MSE \rightarrow \text{mmse}\{\mathbf{x}|\mathbf{y}, \mathbf{A}, \Gamma, \Phi\}$$

For non-Gaussian \mathbf{x} , without the assumptions of $M, N \rightarrow \infty$ and \mathbf{A} , finding the optimal solution is generally NP-hard.

Approximate Message Passing (AMP)

- ▶ AMP-type algorithms:

linear estimator (LE) : $\mathbf{r}_t = \gamma_t(\mathbf{x}_t)$,

non-linear estimator (NLE) : $\mathbf{x}_{t+1} = \phi_t(\mathbf{r}_t)$.

- ▶ AMP:

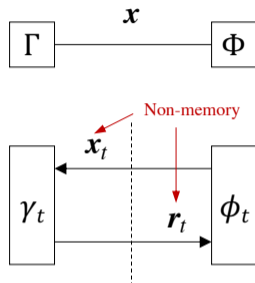
LE : $\mathbf{r}_t = \mathbf{x}_t + \mathbf{A}^H(\mathbf{y} - \mathbf{A}\mathbf{x}_t) + \mathbf{r}_t^{\text{Onsager}}$,

NLE : $\mathbf{x}_{t+1} = \phi(\mathbf{r}_t) = \mathbb{E}\{\mathbf{x}|\mathbf{r}_t\}$,

where $\mathbf{r}_t^{\text{Onsager}} = \beta \langle \phi'(\mathbf{r}_{t-1}) \rangle (\mathbf{r}_{t-1} - \mathbf{x}_{t-1})$.

- ✓ Bayes optimal
- ✓ Low-complexity
- ✗ IID \mathbf{A} is required

- D. L. Donoho, A. Maleki, and A. Montanari, "Message-passing algorithms for compressed sensing," in *Proc. Nat. Acad. Sci.*, 2009.



Orthogonal/Vector AMP (OAMP/VAMP)

► OAMP/VAMP:

$$\text{LE : } \mathbf{r}_t = \mathbf{x}_t + \frac{1}{\epsilon_t^\gamma} \mathbf{A}^H (\rho_t \mathbf{I} + \mathbf{A} \mathbf{A}^H)^{-1} (\mathbf{y} - \mathbf{A} \mathbf{x}_t),$$

$$\text{NLE : } \mathbf{x}_{t+1} = \frac{1}{\epsilon_{t+1}^\phi} [\hat{\phi}_t(\mathbf{r}_t) + (1 - \epsilon_{t+1}^\phi) \mathbf{r}_t],$$

where ϵ_t^γ and ϵ_t^ϕ are orthogonal parameters.

✓ Bayes optimal (replica)

✓ Unitarily-invariant \mathbf{A}

✗ High-complexity

□ J. Ma and L. Ping, "Orthogonal AMP," *IEEE Access*, 2017.

□ S. Rangan, P. Schniter, and A. Fletcher, "Vector approximate message passing," *IEEE Trans. Inf. Theory*, 2019.

Convolutional AMP (CAMP)

► CAMP:

$$\begin{aligned} \text{LE :} \quad & \mathbf{r}_t = \mathbf{x}_t + \mathbf{A}^H(\mathbf{y} - \mathbf{A}\mathbf{x}_t) + \mathbf{r}_t^{\text{Onsager}}, \\ \text{NLE :} \quad & \mathbf{x}_{t+1} = \phi(\mathbf{r}_t) = \text{E}\{\mathbf{x}|\mathbf{r}_t\}, \end{aligned}$$

where $\mathbf{r}_t^{\text{Onsager}} = \sum_{\tau=0}^{t-1} [\prod_{t'=\tau}^{t-1} \langle \phi'(\mathbf{r}_{t'}) \rangle] (\theta_{t-\tau} \mathbf{A}^T \mathbf{A} - g_{t-\tau})(\mathbf{r}_\tau - \mathbf{x}_\tau)$.

- ✓ Bayes optimal (replica), if converges
- ✓ Unitarily-invariant \mathbf{A}
- ✓ Low-complexity
- ✗ Fails to converge for \mathbf{A} with high condition numbers

□ K. Takeuchi, "Bayes-optimal convolutional AMP," *IEEE Trans. Inf. Theory*, 2021.

Memory AMP

- ▶ Memory AMP (MAMP):

$$\text{LE : } \quad \mathbf{r}_t = \gamma_t(\mathbf{X}_t) = \mathbf{Q}_t \mathbf{y} + \sum_{i=1}^t \mathbf{P}_{t,i} \mathbf{x}_i,$$

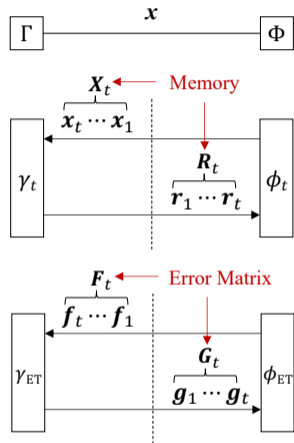
$$\text{NLE : } \quad \mathbf{x}_{t+1} = \phi_t(\mathbf{R}_t),$$

under **orthogonality**: $\forall t \geq 1$,

$$\frac{1}{N} \mathbf{g}_t^H \mathbf{x} \stackrel{\text{a.s.}}{=} 0, \quad \frac{1}{N} \mathbf{F}_t^H \mathbf{g}_t \stackrel{\text{a.s.}}{=} \mathbf{0}, \quad \frac{1}{N} \mathbf{G}_t^H \mathbf{f}_{t+1} \stackrel{\text{a.s.}}{=} \mathbf{0},$$

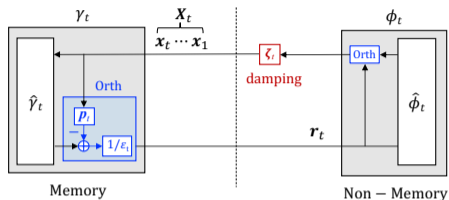
where \mathbf{G}_t and \mathbf{F}_t are error matrices.

- ▶ AMP, OAMP/VAMP, CAMP can be unified under MAMP.



□ L. Liu, S. Huang, and B. M. Kurkoski, "Memory AMP," *IEEE Trans. Inf. Theory*, 2022.

Bayes-Optimal MAMP (BO-MAMP) — Principle



- ▶ The LE has memory and consists of a local estimator $\hat{\gamma}_t$ and orthogonalization.
 1. $\hat{\gamma}_t$ approaches $(\rho_t \mathbf{I} + \mathbf{A} \mathbf{A}^H)^{-1} (\mathbf{y} - \mathbf{A} \mathbf{x}_t)$ in OAMP/VAMP.
 2. The error of r_t is orthogonal to \mathbf{x} and the errors of $\mathbf{x}_1, \dots, \mathbf{x}_t$.
- ▶ Damping ζ is added. ζ is analytically optimized to guarantee convergence. Also improves convergence speed.
- ▶ NLE: Same as OAMP/VAMP and has no memory.

MAMP with Gradient Descent

GD-MAMP:

$$\begin{aligned}\text{LE : } \quad \mathbf{u}_t &= \theta_t \mathbf{B} \mathbf{u}_{t-1} + \xi_t (\mathbf{y} - \mathbf{A} \mathbf{x}_t), \\ \mathbf{r}_t &= \gamma_t (\mathbf{X}_t) = \frac{1}{\varepsilon_t^\gamma} (\mathbf{A}^H \mathbf{u}_t + \sum_{i=1}^t p_{t,i} \mathbf{x}_i),\end{aligned}$$

$$\text{NLE : } \mathbf{x}_{t+1} = [\mathbf{x}_1 \cdots \mathbf{x}_t \phi_t(\mathbf{r}_t)] \cdot \zeta_{t+1},$$

- ▶ \mathbf{u}_t is an estimate of $(\rho_t \mathbf{I} + \mathbf{A} \mathbf{A}^H)^{-1} (\mathbf{y} - \mathbf{A} \mathbf{x}_t)$.
- ▶ The parameters θ_t and ξ_t are optimized, $p_{t,i}$ and ε_t^γ are chosen to ensure orthogonality. Computation not shown, but we'll more talk about these later.
- ▶ ζ_{t+1} is the optimized damping vector
- ▶ $\phi_t(\cdot)$ is the same as that in OAMP/VAMP,

Gradient descent (GD) is used to approximate $\frac{\xi_t}{\theta_t} (\rho_t \mathbf{I} + \mathbf{A} \mathbf{A}^H)^{-1} (\mathbf{y} - \mathbf{A} \mathbf{x}_t)$ by \mathbf{u}_t .

Why Memory? Intuition 1: Gradient Descent Avoids Matrix Inverse

Want to eliminate matrix inverse in OAMP: $\frac{\xi_t}{\theta_t}(\rho_t \mathbf{I} + \mathbf{A}\mathbf{A}^T)^{-1}(\mathbf{y} - \mathbf{A}\mathbf{x}_t)$.

Solve $\mathbf{W}\mathbf{u} = \mathbf{b}$ without finding \mathbf{W}^{-1} .

Intuition 1 Solving $\mathbf{u} = \mathbf{W}^{-1}\mathbf{b}$ is equivalent to

$$\arg \min f(\mathbf{u}) = \frac{1}{2}\mathbf{u}^T \mathbf{W}\mathbf{u} - \mathbf{b}^T \mathbf{u}$$

when \mathbf{W} is positive definite. Find solution using gradient descent with step:

$$\begin{aligned}\mathbf{u}_i &= \mathbf{u}_{i-1} - \alpha \nabla f(\mathbf{u}_{i-1}) \\ &= \mathbf{u}_{i-1} + \alpha(\mathbf{b} - \mathbf{W}\mathbf{u}_{i-1}).\end{aligned}$$

Then \mathbf{u}_i approaches the correct value \mathbf{u} .

Choosing $\alpha = \frac{2}{\lambda_{\max} + \lambda_{\min}}$ is close to optimal. Shown for real-valued case.

ChatGPT: Is it possible to find a matrix inverse using gradient descent?



ChatGPT

While there might be unconventional methods or iterative approaches that can approximate a matrix inverse, gradient descent is not the typical choice for this particular problem. If you need to find the inverse of a matrix, it's recommended to use established linear algebra techniques for accuracy and efficiency.

Why Memory? Intuition 2: Neumann Series for Matrix Inverse

Let $\rho(\mathbf{C})$ denote the spectral radius of \mathbf{C} . If $\rho(\mathbf{C}) < 1$, then

$$(\mathbf{I} - \mathbf{C})^{-1} = \sum_{i=0}^{\infty} \mathbf{C}^i.$$

Choose $\mathbf{C} = \mathbf{I} - \mathbf{W}$. When $\rho(\mathbf{C}) \geq 1$, let $\mathbf{C}' = \mathbf{I} - \theta(\mathbf{I} - \mathbf{C})$, where θ ensures $\rho(\mathbf{C}') < 1$. With $\mathbf{u}_0 = \mathbf{0}$:

$$\mathbf{u}_i = \mathbf{C}'\mathbf{u}_{i-1} + \theta\mathbf{b}.$$

Then \mathbf{u}_i approaches $\mathbf{u} = \mathbf{W}^{-1}\mathbf{b}$. This iteration is identical to gradient descent.

To accelerate convergence, we can minimize $\rho(\mathbf{C}')$ by:

$$\theta = \frac{2}{\lambda_1 + \lambda_2},$$

where λ_1 and λ_2 denote the maximum and minimum eigenvalues of $\mathbf{I} - \mathbf{C}$.

Overview of AMP-Type Algorithms

Table 1: Overview of AMP-Type Algorithms

Algorithm	Matrix \mathbf{A}	Convergence	Time complexity	Optimality
AMP	IID	Converges	Low: $\mathcal{O}(MN)$	Bayes-optimal
OAMP/VAMP	Right unitarily invariant	Converges	High: $\mathcal{O}(M^2N)$	Bayes-optimal
CAMP	Right unitarily invariant	Diverges in high condition numbers	Low: $\mathcal{O}(MN)$	Bayes-optimal
GD-MAMP	Right unitarily invariant	Converges	Low: $\mathcal{O}(MN)$	Bayes-optimal

- D. L. Donoho, A. Maleki, and A. Montanari, "Message-passing algorithms for compressed sensing," in *Proc. Nat. Acad. Sci.*, 2009.
- J. Ma and L. Ping, "Orthogonal AMP," *IEEE Access*, 2017.
- S. Rangan, P. Schniter, and A. Fletcher, "Vector approximate message passing," *IEEE Trans. Inf. Theory*, 2019.
- K. Takeuchi, "Bayes-optimal convolutional AMP," *IEEE Trans. Inf. Theory*, 2021.
- L. Liu, S. Huang, and B. M. Kurkoski, "Memory AMP," *IEEE Trans. Inf. Theory*, 2022.

Comparison of AMP-Type Algorithms

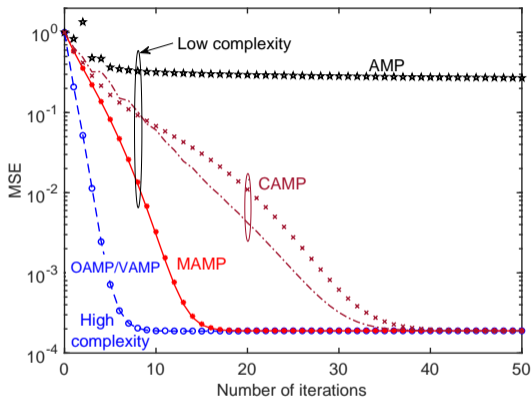


Figure 1: $M = 2^{13}$, $N = 2^{14}$, $\kappa(\mathbf{A}) = 10$, SNR = 30dB

- ▶ AMP:
low complexity
poor MSE
- ▶ OAMP/VAMP:
fastest convergence
high complexity
- ▶ CAMP:
low complexity
slow convergence
incorrect state evolution
- ▶ GD-MAMP:
low complexity
fast convergence
correct state evolution

Outline

1. Preliminaries
 - ▶ Problem Formulation
 - ▶ AMP, OAMP/VAMP, and CAMP
 - ▶ MAMP and GD-MAMP
2. Overflow-Avoiding GD-MAMP
 - ▶ Overflow Problem in GD-MAMP
 - ▶ OA-GD-MAMP with eigenvalues of $\mathbf{A}\mathbf{A}^H$
 - ▶ OA-GD-MAMP without eigenvalues of $\mathbf{A}\mathbf{A}^H$
3. Complexity-Reduced GD-MAMP
 - ▶ Complexity analysis of GD-MAMP
 - ▶ Complexity-Reduced GD-MAMP
4. Conclusion

Overflow Problem in GD-MAMP

- ▶ Representation of floating point numbers from the IEEE 754 technical standard:
 - ▶ binary32 (single precision): $\pm 1.18 \times 10^{-38} \sim \pm 3.4 \times 10^{38}$
 - ▶ binary64 (double precision): $\pm 2.23 \times 10^{-308} \sim \pm 1.8 \times 10^{308}$

Double-precision is widely used, including in Matlab and Python

- ▶ The dreaded NaN will appear for values that are too large.
- ▶ Overflow problem: In GD-MAMP, some intermediate variables may increase exponentially, and exceed the maximum value of double precision.

Which intermediate variables cause overflow in GD-MAMP?

- ▶ In iteration t , the parameters (1) ξ_t (2) $p_{t,i}$ and (3) $v_{t,t}^\gamma$ (variance of \mathbf{r}_t) require w_t .
- ▶ $\lambda^\dagger = (\lambda_{\max} + \lambda_{\min})/2$ and $\mathbf{B} = \lambda^\dagger \mathbf{I} - \mathbf{A}\mathbf{A}^H$, where $\lambda_{\max}, \lambda_{\min}$ denotes the maximum and minimum eigenvalues of $\mathbf{A}\mathbf{A}^H$.

$$b_k \equiv \frac{1}{N} \text{tr}\{\mathbf{B}^k\},$$

$$w_k \equiv \frac{1}{N} \text{tr}\{\mathbf{A}^H \mathbf{B}^k \mathbf{A}\} = \lambda^\dagger b_k - b_{k+1}.$$

- ▶ b_k can be computed if the eigenvalues of $\mathbf{A}\mathbf{A}^H$ are known. Otherwise, there are simple methods to approximate $\lambda_{\max}, \lambda_{\min}$ and b_k .
- ▶ If $\lambda_{\max} > \lambda_{\min} + 2$, the spectral radius $\rho(\mathbf{B}) > 1$, b_{2k} increases exponentially.
- ▶ Even w_k may increase exponentially.

Overflow of w_k

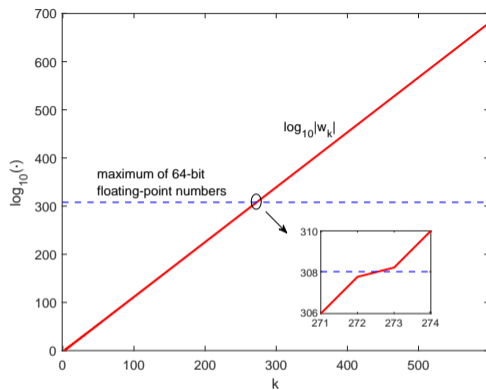


Figure 2: $\log_{10} |w_k|$ versus k

- ▶ w_1, \dots, w_{2T} are required, where T is the maximum number of iterations.
- ▶ As shown in Figure 2, $|w_k|$ increases exponentially as k increases.

Overflow Problem in GD-MAMP

- ▶ To compute ξ_t , $p_{t,i}$ and $v_{t,t}^\gamma$, we need w_k :

$$b_k \equiv \frac{1}{N} \text{tr}\{\mathbf{B}^k\},$$
$$w_k \equiv \frac{1}{N} \text{tr}\{\mathbf{A}^H \mathbf{B}^k \mathbf{A}\} = \lambda^\dagger b_k - b_{k+1}.$$

- ▶ While w_k increases exponentially, it always appears in the product ϑw_k , which is bounded (i.e. ϑ is small).
- ▶ For any $\vartheta \in \mathbb{R} \setminus \{0\}$, the following holds:

$$\vartheta w_k = \frac{\text{sgn}(\vartheta)}{N} \mathbf{1}^T [(\lambda^\dagger \mathbf{1} - \boldsymbol{\lambda}_B) \circ \mathbf{s}_\lambda^{\circ k} \circ e^{\circ \log |\vartheta| \mathbf{1} + k \boldsymbol{\lambda}_B^{\log}}],$$

where $\boldsymbol{\lambda}_B$ denotes the eigenvalues of \mathbf{B} , $\mathbf{s}_\lambda \equiv \text{sgn}(\boldsymbol{\lambda}_B)$, $\boldsymbol{\lambda}_B^{\log} \equiv \log^\circ |\boldsymbol{\lambda}_B|$ and \circ is component-wise operation.

Overflow-Avoiding GD-MAMP with Eigenvalues of AA^H

- ▶ ϑw_k requires $\mathcal{O}(M)$ computations, GD-MAMP requires to compute $\mathcal{O}(T^3)$ terms involving w_k , where T is the number of iterations. The overall complexity is $\mathcal{O}(MT^3)$. How to reduce the complexity?
- ▶ Define

$$\chi_k \equiv \theta_0^k w_k,$$

where $\theta_0 = (\lambda^\dagger + \sigma^2)^{-1} > 0$. We pre-compute $\chi_1, \dots, \chi_{2T-1}$ before the iterations. Computing ϑw_k can be reduced to a scalar operation:

$$\vartheta w_k = \text{sgn}(\vartheta) e^{\log|\alpha| - k \log \vartheta} \chi_k.$$

- ▶ Pre-computing $\chi_1, \dots, \chi_{2T-1}$ costs $\mathcal{O}(MT)$, and computing terms involving w_k in iterations costs $\mathcal{O}(T^3)$. The overall complexity is reduced to $\mathcal{O}(MT + T^3)$.

Overflow-Avoiding GD-MAMP with Eigenvalues of $\mathbf{A}\mathbf{A}^H$

Theorem (2)

For any $k \geq 0$,

$$|\chi_k| \leq \delta(\lambda^\dagger + \theta_0^{-1}),$$

where $\delta = M/N$.

- ▶ Theorem 2 shows that χ_k is bounded. In other words, computing χ_k has no risks of overflow.

Overflow-Avoiding GD-MAMP without Eigenvalues of $\mathbf{A}\mathbf{A}^H$

In large-scale systems, computation of eigenvalues of $\mathbf{A}\mathbf{A}^H$ may be impractical.

- ▶ A method to estimate the maximum and minimum eigenvalue λ_{\max} and λ_{\min} was given in [LHK22].
- ▶ For $k \geq 0$, χ_k can be estimated by

$$\chi_k = \bar{\mathbf{h}}_i^H \bar{\mathbf{h}}_{k-i}$$

where $i = \lceil k/2 \rceil$ and $\bar{\mathbf{h}}_i$ is given by a recursion

$$\bar{\mathbf{h}}_i = \theta_0 (\lambda^\dagger \mathbf{I} - \mathbf{A}\mathbf{A}^H) \bar{\mathbf{h}}_{i-1}$$

with $\bar{\mathbf{h}}_0 = \mathbf{A}\mathbf{h}_0$, $\mathbf{h}_0 \sim \mathcal{N}(\mathbf{0}, \frac{1}{N}\mathbf{I}_N)$.

Simulation Results

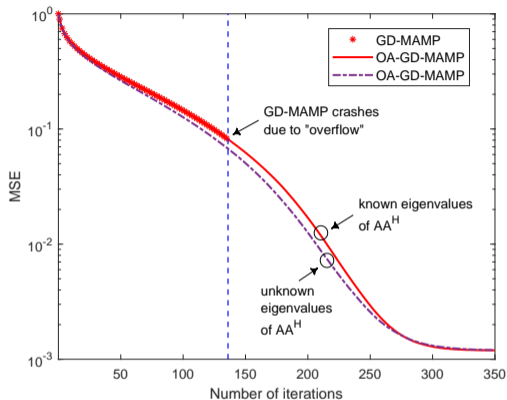


Figure 3: $M = 2^{13}$, $N = 2^{14}$, $\kappa(\mathbf{A}) = 1000$, SNR = 35dB

- ▶ When $t > 136$, the unmodified GD-MAMP does not reach the fixed point since b_{137} and w_{137} overflows.
- ▶ Both OA-GD-MAMP with and without eigenvalues of $\mathbf{A}\mathbf{A}^H$ work properly.

Outline

1. Preliminaries

- ▶ Problem Formulation
- ▶ AMP, OAMP/VAMP, and CAMP
- ▶ MAMP and GD-MAMP

2. Overflow-Avoiding GD-MAMP

- ▶ Overflow Problem in GD-MAMP
- ▶ OA-GD-MAMP with eigenvalues of $\mathbf{A}\mathbf{A}^H$
- ▶ OA-GD-MAMP without eigenvalues of $\mathbf{A}\mathbf{A}^H$

3. Complexity-Reduced GD-MAMP

- ▶ Complexity analysis of GD-MAMP
- ▶ Complexity-Reduced GD-MAMP

4. Conclusion

Complexity Analysis of GD-MAMP

Let T be the number of iterations. The main complexity of GD-MAMP is $\mathcal{O}(MNT)$, dominated by the number of matrix-vector products, each with $\mathcal{O}(MN)$

$$\begin{aligned}\text{LE : } \quad \mathbf{u}_t &= \theta_t \mathbf{B} \mathbf{u}_{t-1} + \xi_t (\mathbf{y} - \mathbf{A} \mathbf{x}_t), \\ \mathbf{r}_t &= \frac{1}{\varepsilon_t^\gamma} (\mathbf{A}^H \mathbf{u}_t + \sum_{i=1}^t p_{t,i} \mathbf{x}_i),\end{aligned}$$

$$\text{NLE : } \mathbf{x}_{t+1} = [\mathbf{x}_1 \cdots \mathbf{x}_t \phi_t(\mathbf{r}_t)] \cdot \zeta_{t+1}.$$

GD-MAMP requires 4 matrix-vector products per iteration?

- ▶ Computing $\mathbf{A} \phi_{t-1}(\mathbf{r}_{t-1})$ to estimate $v_{t,1}^\phi, \dots, v_{t,t}^\phi$ requires one (hidden in ζ_t).
- ▶ Computing $\mathbf{B} \mathbf{u}_{t-1} = (\lambda^\dagger \mathbf{I} - \mathbf{A} \mathbf{A}^H) \mathbf{u}_{t-1}$ requires two.
- ▶ Computing $\mathbf{A}^H \mathbf{u}_t$ requires one.

Easily eliminate one matrix-vector product:

- ▶ Two of the products are $\mathbf{A}^H \mathbf{u}_t$, these only need to be computed once.

GD-MAMP requires 3 matrix-vector products per iteration!

Complexity-Reduced GD-MAMP Using Approximate ξ_t

Finding the damping vector ζ_t nominally requires one matrix-vector product:

1. Compute $\mathbf{z}_t = \mathbf{y} - \mathbf{A}\phi_{t-1}(\mathbf{r}_{t-1})$.
2. Estimate $v_{t,1}^\phi, \dots, v_{t,t}^\phi$ by using \mathbf{z}_t , where $v_{t,i}^\phi$ denotes the covariance of $\phi_{t-1}(\mathbf{r}_{t-1})$ and \mathbf{x}_i for $i < t$.
3. Compute ζ_t from the covariance matrix \mathbf{V}_t^ϕ of $\mathbf{x}_1, \dots, \mathbf{x}_{t-1}$ and $\phi_{t-1}(\mathbf{r}_{t-1})$.

To remove the above matrix-vector product:

- (1) We move the damping from the NLE to the LE (details omitted).
- (2) ξ_t nominally depends on \mathbf{z}_t and \mathbf{V}_t^ϕ . But we found that approximating ξ_t gave little to no performance loss:

$$\tilde{\xi}_t = 1/(v_{t,t}^\phi + \sigma^2).$$

where $v_{t,t}^\phi$ is the variance of \mathbf{x}_t , given as that in OAMP/VAMP.

- The resulting complexity-reduced GD-MAMP requires only 2 matrix-vector products per iteration.

Simulation Results

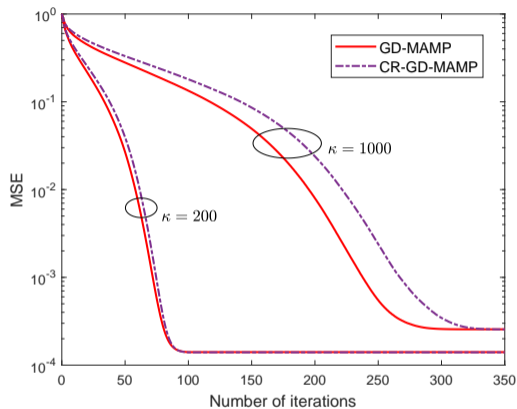


Figure 4: MSE versus number of iterations, $M = 2^{13}$, $N = 2^{14}$, SNR = 35dB

- ▶ CR-GD-MAMP requires a few more iterations to converge.

Simulation Results

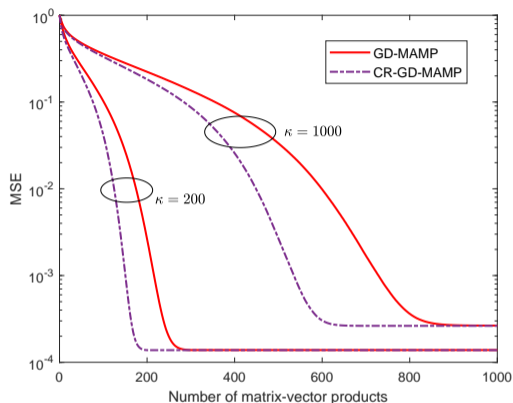


Figure 5: MSE versus number of matrix-vector products, $M = 2^{13}$, $N = 2^{14}$, SNR = 35dB

- ▶ CR-GD-MAMP achieves the same MSE as GD-MAMP while requiring about 2/3 matrix-vector products.

Conclusion

GD-MAMP is a memory AMP algorithm that:

- ▶ Converges for unitarily-invariant matrices
- ▶ is low complexity, avoiding matrix inverse
- ▶ converges for \mathbf{A} of high condition number

overcoming the weakness of AMP, OAMP/VAMP and CAMP.

(1) To solve the overflow problem, we propose OA-GD-MAMP:

- ▶ With known eigenvalues of $\mathbf{A}\mathbf{A}^H$, OA-GD-MAMP is equivalent to GD-MAMP.
- ▶ Otherwise, OA-GD-MAMP can achieve nearly the same performance.

(2) GD-MAMP requires three matrix-vector products per iteration. To reduce it:

- ▶ We propose CR-GD-MAMP as a variant of GD-MAMP. It requires only two matrix-vector products per iteration with almost the same convergence speed.