

Approximate Message Passing for High-Dimensional Inference, I

Ramji Venkataramanan
University of Cambridge

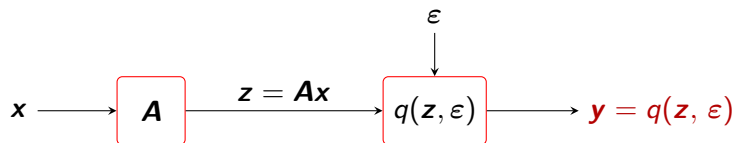
ISIT 2023

Focus of tutorial

Approximate Message Passing (AMP) for

1. Estimation in linear and generalized linear models
2. Low-rank matrix estimation

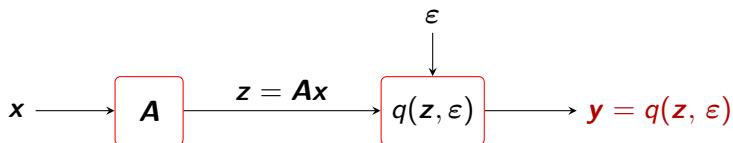
Generalized Linear Models (GLMs)



GOAL:

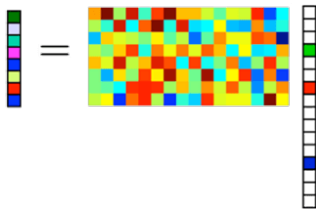
- ▶ Estimate signal $\mathbf{x} \in \mathbb{R}^d$ from observations $\mathbf{y} \equiv (y_1, \dots, y_n)$
- ▶ Known sensing matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ and output function q

Example: Linear model

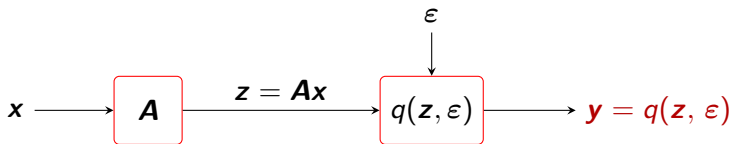


Linear model: $\mathbf{y} = \mathbf{A}\mathbf{x} + \epsilon$

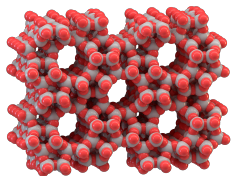
- ▶ Widely used model in signal processing and communications: CDMA, MIMO, sparse regression codes ...
- ▶ Compressed sensing: Signal \mathbf{x} assumed to be sparse



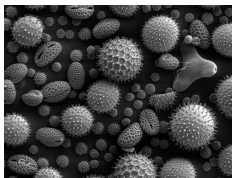
Example: Phase retrieval



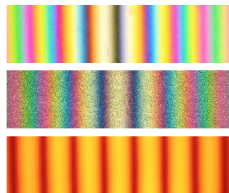
Phase retrieval: $y = |\mathbf{A}x|^2 + \epsilon$



X-ray crystallography

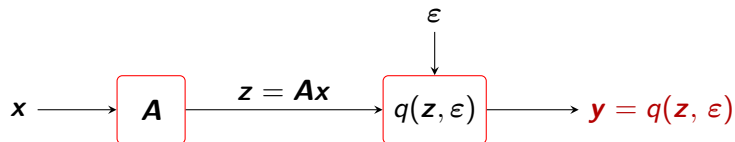


Microscopy



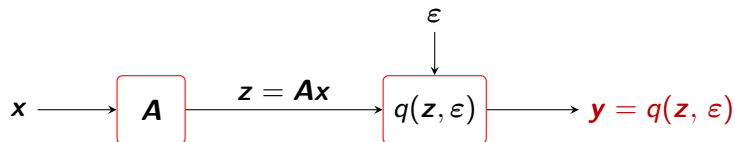
Interferometry

Example: 1-bit compressed sensing



1-bit compressed sensing [Boufounos '08]: $\mathbf{y} = \text{sign}(\mathbf{A}\mathbf{x} + \epsilon)$

Example: 1-bit compressed sensing



1-bit compressed sensing [Boufounos '08]: $\mathbf{y} = \text{sign}(\mathbf{A}\mathbf{x} + \epsilon)$

Many other popular GLMs, e.g.,

- ▶ Logistic, probit regression (Binary classification)
- ▶ Poisson regression (count data)

Low-rank models

$$\begin{bmatrix} & & \\ & \mathbf{A} & \\ & & \end{bmatrix} \approx \begin{bmatrix} & & \\ & \mathbf{U} & \\ & & \end{bmatrix} \begin{bmatrix} & & \\ & \mathbf{V}^T & \\ & & \end{bmatrix}$$

$n \times d$ $n \times k$ $k \times d$

Topic Modelling

- ▶ Each row of \mathbf{A} is a document
- ▶ Each row of \mathbf{V}^T is a topic
- ▶ Each document convex combination of k topics

Hidden clique

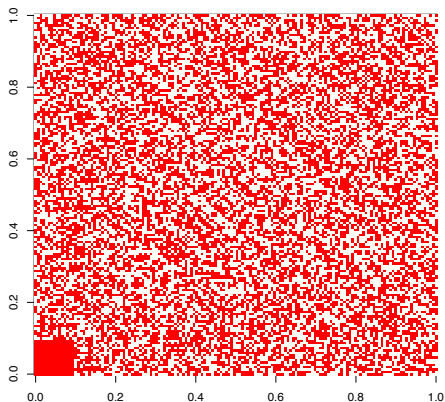


Image from *Statistical Estimation: From Denoising to Sparse Regression and Hidden Cliques* by A. Montanari

Hidden clique

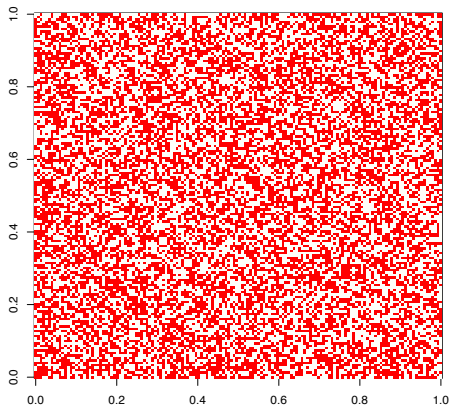


Image from *Statistical Estimation: From Denoising to Sparse Regression and Hidden Cliques* by A. Montanari

Hidden clique

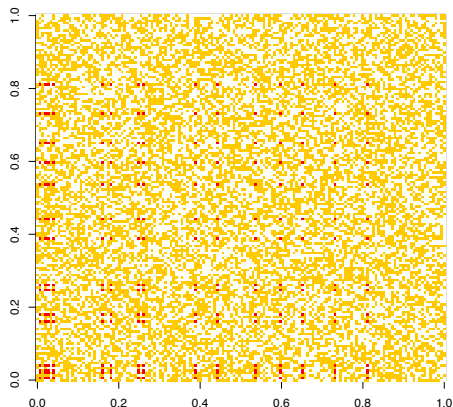
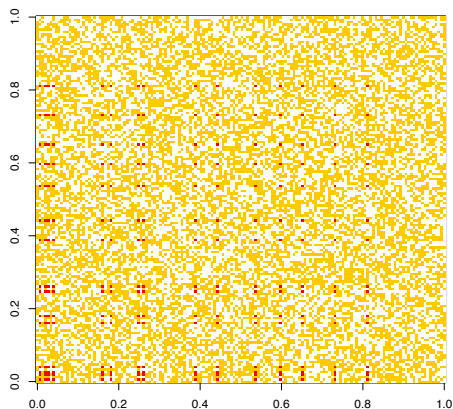


Image from *Statistical Estimation: From Denoising to Sparse Regression and Hidden Cliques* by A. Montanari

Hidden clique



For hidden clique S , adjacency matrix has the form

$$\mathbf{A} = \mathbf{1}_S \mathbf{1}_S^T + \mathbf{W}$$

Structure of Tutorial

1. Introduction to AMP, application to low-rank matrix estimation
2. AMP to derive exact asymptotics in generalized linear models (Cynthia Rush)
3. AMP as a flexible tool in high-dimensional statistics (Marco Mondelli)

Origins of AMP

- ▶ Relaxation of belief propagation for CDMA multiuser detection:
[Kabashima '03], [Caire, Muller, Tanaka '04], [Tanaka, Okada '05]
- ▶ Via systematic approximation of BP iterations:
 1. Compressed sensing (linear models): [Donoho, Maleki, Montanari '09], [Krzakala et. al '11]
 2. Generalized linear models: [Rangan '11]
 3. Low-rank matrix estimation: [Parker, Schniter, Cevher '14], [Fletcher, Rangan '18], [Lesieur et al., '17]

Origins of AMP

- ▶ Relaxation of belief propagation for CDMA multiuser detection:
[Kabashima '03], [Caire, Muller, Tanaka '04], [Tanaka, Okada '05]
- ▶ Via systematic approximation of BP iterations:
 1. Compressed sensing (linear models): [Donoho, Maleki, Montanari '09], [Krzakala et. al '11]
 2. Generalized linear models: [Rangan '11]
 3. Low-rank matrix estimation: [Parker, Schniter, Cevher '14], [Fletcher, Rangan '18], [Lesieur et al., '17]

We'll take a different approach to understanding AMP:
Study it as an iteration defined via a random matrix

Gaussian Orthogonal Ensemble (GOE)

Consider a **symmetric** Gaussian matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$

W_{ij} independent for $1 \leq i \leq j \leq n$

$$W_{ij} \sim N\left(0, \frac{1}{n}\right) \text{ for } i \neq j, \quad W_{ij} \sim N\left(0, \frac{2}{n}\right) \text{ for } i = j.$$

We write $\mathbf{W} \sim \text{GOE}(n)$

Gaussian Orthogonal Ensemble (GOE)

Consider a **symmetric** Gaussian matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$

W_{ij} independent for $1 \leq i \leq j \leq n$

$$W_{ij} \sim N\left(0, \frac{1}{n}\right) \text{ for } i \neq j, \quad W_{ij} \sim N\left(0, \frac{2}{n}\right) \text{ for } i = j.$$

We write $\mathbf{W} \sim \text{GOE}(n)$

Property

If $\mathbf{W} \sim \text{GOE}(n)$ and \mathbf{Q} is any $n \times n$ orthogonal matrix, then:

$$\mathbf{Q}^T \mathbf{W} \mathbf{Q} \sim \text{GOE}(n)$$

An iteration with a GOE matrix

Let \mathbf{W} be a GOE matrix

Starting with an initialization $\mathbf{h}^0 \in \mathbb{R}^n$, define for $t \geq 0$:

$$\mathbf{m}^t = f_t(\mathbf{h}^t), \quad \mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - b_t \mathbf{m}^{t-1}$$

An iteration with a GOE matrix

Let \mathbf{W} be a GOE matrix

Starting with an initialization $\mathbf{h}^0 \in \mathbb{R}^n$, define for $t \geq 0$:

$$\mathbf{m}^t = f_t(\mathbf{h}^t), \quad \mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - b_t \mathbf{m}^{t-1}$$

- ▶ Function f_t is Lipschitz and acts component-wise, for $t \geq 1$
- ▶ Coefficient $b_t = \frac{1}{n} \sum_{i=1}^n f_t'(h_i^t)$
- ▶ First step: $\mathbf{h}^1 = \mathbf{W} f_0(\mathbf{h}^0)$

We call this the **abstract AMP** recursion

State Evolution

$$\mathbf{m}^t = f_t(\mathbf{h}^t), \quad \mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - \mathbf{b}_t \mathbf{m}^{t-1}$$

Key result (informal): If initialization \mathbf{h}^0 is independent of \mathbf{W} , then for $t \geq 1$, as $n \rightarrow \infty$, the empirical distribution of \mathbf{h}^t converges to $N(0, \tau_t^2)$, where

$$\tau_{t+1}^2 = \mathbb{E}\{(f_t(G_t))^2\}, \quad G_t \sim N(0, \tau_t^2)$$

- ▶ The $\tau_t \rightarrow \tau_{t+1}$ recursion is called **state evolution**
- ▶ Initialized with $\tau_1^2 = \lim_{n \rightarrow \infty} \frac{\|f_0(\mathbf{h}^0)\|^2}{n}$

State Evolution

$$\mathbf{m}^t = f_t(\mathbf{h}^t), \quad \mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - \mathbf{b}_t \mathbf{m}^{t-1}$$

Key result (informal): If initialization \mathbf{h}^0 is independent of \mathbf{W} , then for $t \geq 1$, as $n \rightarrow \infty$, the empirical distribution of \mathbf{h}^t converges to $\mathcal{N}(0, \tau_t^2)$, where

$$\tau_{t+1}^2 = \mathbb{E}\{(f_t(G_t))^2\}, \quad G_t \sim \mathcal{N}(0, \tau_t^2)$$

- ▶ The $\tau_t \rightarrow \tau_{t+1}$ recursion is called **state evolution**
- ▶ Initialized with $\tau_1^2 = \lim_{n \rightarrow \infty} \frac{\|f_0(\mathbf{h}^0)\|^2}{n}$

Why is this true? Why is it interesting?

Heuristic for state evolution

First step: $\mathbf{h}^1 = \mathbf{W}\mathbf{m}^0$

- ▶ Let $\nu_n(\mathbf{h}^1)$ denote empirical distribution of \mathbf{h}^1
- ▶ Since \mathbf{m}^0 is independent of \mathbf{W} , we have \mathbf{h}^1 Gaussian with $\nu_n(\mathbf{h}^1) \rightarrow \mathcal{N}(0, \tau_1^2)$ where

$$\tau_1^2 = \lim_{n \rightarrow \infty} \frac{\|\mathbf{m}^0\|^2}{n} = \lim_{n \rightarrow \infty} \frac{\|f_0(\mathbf{h}^0)\|^2}{n}$$

Heuristic for state evolution

First step: $\mathbf{h}^1 = \mathbf{W}\mathbf{m}^0$

- ▶ Let $\nu_n(\mathbf{h}^1)$ denote empirical distribution of \mathbf{h}^1
- ▶ Since \mathbf{m}^0 is independent of \mathbf{W} , we have \mathbf{h}^1 Gaussian with $\nu_n(\mathbf{h}^1) \rightarrow \mathcal{N}(0, \tau_1^2)$ where

$$\tau_1^2 = \lim_{n \rightarrow \infty} \frac{\|\mathbf{m}^0\|^2}{n} = \lim_{n \rightarrow \infty} \frac{\|f_0(\mathbf{h}^0)\|^2}{n}$$

Second step: $\mathbf{h}^2 = \mathbf{W}\mathbf{m}^1 - b_1\mathbf{m}^0$, with $\mathbf{m}^1 = f_1(\mathbf{h}^1)$

- ▶ \mathbf{W} and \mathbf{m}^1 are *dependent*, so $\mathbf{W}\mathbf{m}^1$ is **not** Gaussian

Heuristic for state evolution

First step: $\mathbf{h}^1 = \mathbf{W}\mathbf{m}^0$

- ▶ Let $\nu_n(\mathbf{h}^1)$ denote empirical distribution of \mathbf{h}^1
- ▶ Since \mathbf{m}^0 is independent of \mathbf{W} , we have \mathbf{h}^1 Gaussian with $\nu_n(\mathbf{h}^1) \rightarrow \mathcal{N}(0, \tau_1^2)$ where

$$\tau_1^2 = \lim_{n \rightarrow \infty} \frac{\|\mathbf{m}^0\|^2}{n} = \lim_{n \rightarrow \infty} \frac{\|f_0(\mathbf{h}^0)\|^2}{n}$$

Second step: $\mathbf{h}^2 = \mathbf{W}\mathbf{m}^1 - b_1\mathbf{m}^0$, with $\mathbf{m}^1 = f_1(\mathbf{h}^1)$

- ▶ \mathbf{W} and \mathbf{m}^1 are *dependent*, so $\mathbf{W}\mathbf{m}^1$ is **not** Gaussian
- ▶ For $\tilde{\mathbf{W}} \sim \text{GOE}(n)$ independent of \mathbf{m}^1 , we have $\tilde{\mathbf{W}}\mathbf{m}^1$ Gaussian with $\nu_n(\tilde{\mathbf{W}}\mathbf{m}^1) \rightarrow \mathcal{N}(0, \tau_2^2)$, where

$$\tau_2^2 = \lim_{n \rightarrow \infty} \frac{\|\mathbf{m}^1\|^2}{n} = \lim_{n \rightarrow \infty} \frac{\|f_1(\mathbf{h}^1)\|^2}{n} = \mathbb{E}\{f_1(G_1)^2\}, \quad G_1 \sim \mathcal{N}(0, \tau_1^2)$$

Debiasing term

$$\mathbf{h}^2 = \mathbf{W} \mathbf{m}^1 - b_1 \mathbf{m}^0, \quad b_1 = \frac{1}{n} \sum_{i=1}^n f_1'(h_i^1)$$

- ▶ The 'Onsager' correction $-b_1 \mathbf{m}^0$ is as a **debiasing** term
- ▶ Ensures that \mathbf{h}^2 asymptotically has the same empirical distribution as $\tilde{\mathbf{W}} \mathbf{m}^1$. That is, $\nu_n(\mathbf{h}^2) \rightarrow \mathcal{N}(0, \tau_2^2)$

$$\mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - b_t \mathbf{m}^{t-1}, \quad b_t = \frac{1}{n} \sum_{i=1}^n f_t'(h_i^t)$$

Debiasing term

$$\mathbf{h}^2 = \mathbf{W} \mathbf{m}^1 - b_1 \mathbf{m}^0, \quad b_1 = \frac{1}{n} \sum_{i=1}^n f_1'(h_i^1)$$

- ▶ The 'Onsager' correction $-b_1 \mathbf{m}^0$ is as a **debiasing** term
- ▶ Ensures that \mathbf{h}^2 asymptotically has the same empirical distribution as $\tilde{\mathbf{W}} \mathbf{m}^1$. That is, $\nu_n(\mathbf{h}^2) \rightarrow \mathcal{N}(0, \tau_2^2)$

$$\mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - b_t \mathbf{m}^{t-1}, \quad b_t = \frac{1}{n} \sum_{i=1}^n f_t'(h_i^t)$$

- ▶ Conditional distribution of $\mathbf{W} \mathbf{m}^t$ given $(\mathbf{m}^0, \dots, \mathbf{m}^t)$ can be decomposed into Gaussian component and a non-Gaussian one
- ▶ Non-Gaussian part asymptotically cancelled out by $-b_t \mathbf{m}^{t-1}$

Pseudo-Lipschitz test functions

A function $\phi : \mathbb{R}^m \rightarrow \mathbb{R}$ is called **pseudo-Lipschitz** if for all inputs $\mathbf{x}, \mathbf{y} \in \mathbb{R}^m$,

$$|\phi(\mathbf{x}) - \phi(\mathbf{y})| \leq C \|\mathbf{x} - \mathbf{y}\| (1 + \|\mathbf{x}\| + \|\mathbf{y}\|)$$

for some constant $C > 0$

- ▶ Roughly: Functions with at most quadratic growth
- ▶ Examples: $\phi(x) = x^2$, $\phi(x, y) = xy$

State evolution results for AMP often stated in terms of pseudo-Lipschitz test functions

E.g., mean-squared error (MSE) of estimate $\phi(x, y) = (x - y)^2$

Main result for abstract AMP

$$\mathbf{m}^t = f_t(\mathbf{h}^t), \quad \mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - \mathbf{b}_t \mathbf{m}^{t-1}$$

Assumptions:

- ▶ Functions f_t Lipschitz, for $t \geq 1$
- ▶ Initialization \mathbf{h}^0 is independent of \mathbf{W}

Theorem [Bolthausen '10, Bayati-Montanari '11]

For $t \geq 1$, and any pseudo-Lipschitz function $\phi : \mathbb{R} \rightarrow \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \phi(h_i^t) = \mathbb{E}\{\phi(G_t)\} \text{ almost surely}$$

where $G_t \sim N(0, \tau_t^2)$, with $\tau_{t+1}^2 = \mathbb{E}\{f_t(G_t)^2\}$.

Main result for abstract AMP

$$\mathbf{m}^t = f_t(\mathbf{h}^t), \quad \mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - \mathbf{b}_t \mathbf{m}^{t-1}$$

Assumptions:

- ▶ Functions f_t Lipschitz, for $t \geq 1$
- ▶ Initialization \mathbf{h}^0 is independent of \mathbf{W}

Theorem [Bolthausen '10, Bayati-Montanari '11]

For $t \geq 1$, and any pseudo-Lipschitz function $\phi : \mathbb{R} \rightarrow \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \phi(h_i^t) = \mathbb{E}\{\phi(G_t)\} \text{ almost surely}$$

where $G_t \sim N(0, \tau_t^2)$, with $\tau_{t+1}^2 = \mathbb{E}\{f_t(G_t)^2\}$.

Equivalent to: empirical distribution $\nu_n(\mathbf{h}^t)$ converges to $N(0, \tau_t^2)$ almost surely (in Wasserstein-2 distance)

Stronger statement

$$\mathbf{m}^t = f_t(\mathbf{h}^t), \quad \mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - \mathbf{b}_t \mathbf{m}^{t-1}$$

Theorem [Javanmard-Montanari '13]

For $t \geq 1$, and any pseudo-Lipschitz function $\phi : \mathbb{R}^t \rightarrow \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \phi(h_i^1, h_i^2, \dots, h_i^t) = \mathbb{E}\{\phi(G_1, G_2, \dots, G_t)\} \text{ almost surely}$$

where $(G_1, \dots, G_t) \sim \mathcal{N}(0, \Sigma_t)$, where $\Sigma_t \in \mathbb{R}^{t \times t}$ can be recursively computed via state evolution, for $t \geq 1$.

Empirical distribution of rows of $\nu_n(\mathbf{h}^1, \dots, \mathbf{h}^t)$ converges (in Wasserstein-2 distance) to $\mathcal{N}(0, \Sigma_t)$ almost surely

Rank-1 matrix estimation

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}$$

- ▶ Signal $\mathbf{v} \in \mathbb{R}^n$, entries $v_j \sim \text{iid } P_V$
- ▶ Noise matrix $\mathbf{W} \sim \text{GOE}(n)$

Rank-1 matrix estimation

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}$$


- ▶ Signal $\mathbf{v} \in \mathbb{R}^n$, entries $v_i \sim_{\text{iid}} P_V$
- ▶ Noise matrix $\mathbf{W} \sim \text{GOE}(n)$

Natural estimator: $\hat{\varphi}$ the principal eigenvector of \mathbf{A}

Random matrix theory shows phase transition:

$$\text{Principal eigenvalue } \lambda_1(\mathbf{A}) \rightarrow \begin{cases} \lambda + \lambda^{-1}, & \text{if } \lambda > 1, \\ 2, & \text{if } \lambda \in (0, 1] \end{cases}$$

$$\text{Correlation } \frac{|\langle \hat{\varphi}, \mathbf{v} \rangle|}{\|\hat{\varphi}\| \|\mathbf{v}\|} \rightarrow \begin{cases} \sqrt{1 - \lambda^{-2}}, & \text{if } \lambda > 1, \\ 0, & \text{if } \lambda \in (0, 1] \end{cases}$$

[Baik, Ben Arous, P\'ech\'e '05], [Baik, Silverstein '06], [Capitaine, Donati-Martin, F\'eral '09], [Benaych-Georges and Nadakuditi '11], 

Structural information

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}$$

Spectral estimator $\hat{\varphi}$ doesn't use *structural* information about \mathbf{v}

- ▶ For example, \mathbf{v} may be sparse, bounded, non-negative etc.
- ▶ Relevant in sparse PCA, non-negative PCA, hidden clique, community detection under stochastic block model, ...

[Deshpande, Montanari '14], [Barbier *et al.* '16], [Lesieur *et al.* '17],
[Miolane, Lelarge '16] ...

Structural information

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}$$

Spectral estimator $\hat{\varphi}$ doesn't use *structural* information about \mathbf{v}

- ▶ For example, \mathbf{v} may be sparse, bounded, non-negative etc.
- ▶ Relevant in sparse PCA, non-negative PCA, hidden clique, community detection under stochastic block model, ...

If we know prior P_V on entries of \mathbf{v} , MMSE estimator is

$$\widehat{\mathbf{M}}_{\text{Bayes}} = \mathbb{E} \left[\mathbf{v} \mathbf{v}^T \mid \mathbf{A} \right]$$

$\widehat{\mathbf{M}}_{\text{Bayes}}$ is generally not computable, but computable formula for asymptotic Bayes risk available

[Deshpande, Montanari '14], [Barbier *et al.* '16], [Lesieur *et al.* '17],
[Miolane, Lelarge '16] ...

AMP for rank-1 estimation

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^\top + \mathbf{W}, \quad \mathbf{W} \sim \text{GOE}(n)$$

Let's try same AMP iteration as before, but defined via \mathbf{A}

$$\hat{\mathbf{v}}^t = f_t(\mathbf{v}^t), \quad \mathbf{v}^{t+1} = \mathbf{A} \hat{\mathbf{v}}^t - b_t \hat{\mathbf{v}}^{t-1}, \quad b_t = \frac{1}{n} \sum_{i=1}^n f_t'(v_i^t)$$

AMP for rank-1 estimation

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}, \quad \mathbf{W} \sim \text{GOE}(n)$$

Let's try same AMP iteration as before, but defined via \mathbf{A}

$$\hat{\mathbf{v}}^t = f_t(\mathbf{v}^t), \quad \mathbf{v}^{t+1} = \mathbf{A} \hat{\mathbf{v}}^t - b_t \hat{\mathbf{v}}^{t-1}, \quad b_t = \frac{1}{n} \sum_{i=1}^n f_t'(v_i^t)$$

Using the expression for \mathbf{A} :

$$\mathbf{v}^{t+1} = \lambda \frac{\langle \mathbf{v}, \hat{\mathbf{v}}^t \rangle}{n} \mathbf{v} + \mathbf{W} \hat{\mathbf{v}}^t - b_t \hat{\mathbf{v}}^{t-1}$$

Shift + abstract AMP iterate

First iteration

Suppose

$$\mathbf{v}^0 = \mu_0 \mathbf{v} + \mathbf{g}^0, \quad \text{with } \mathbf{g}_0 \sim \mathcal{N}(0, \sigma_0^2 \mathbf{I}_n)$$

for some constants μ_0, σ_0 . Then

$$\mathbf{v}^1 = \lambda \frac{\langle \mathbf{v}, \hat{\mathbf{v}}^0 \rangle}{n} \mathbf{v} + \mathbf{W} \hat{\mathbf{v}}^0, \quad \hat{\mathbf{v}}^0 = f_0(\mathbf{v}_0)$$

First iteration

Suppose

$$\mathbf{v}^0 = \mu_0 \mathbf{v} + \mathbf{g}^0, \quad \text{with } \mathbf{g}_0 \sim \mathcal{N}(0, \sigma_0^2 \mathbf{I}_n)$$

for some constants μ_0, σ_0 . Then

$$\mathbf{v}^1 = \lambda \frac{\langle \mathbf{v}, \hat{\mathbf{v}}^0 \rangle}{n} \mathbf{v} + \mathbf{W} \hat{\mathbf{v}}^0, \quad \hat{\mathbf{v}}^0 = f_0(\mathbf{v}_0)$$

► Signal term:

$$\frac{\lambda \langle \mathbf{v}, \hat{\mathbf{v}}^0 \rangle}{n} = \frac{\lambda}{n} \sum_{i=1}^n v_i f_0(v_i^0) \rightarrow \mathbb{E}\{V f_0(\mu_0 V + G_0)\} =: \mu_1$$

where $V \sim P_V$ and $G_0 \sim \mathcal{N}(0, \sigma_0^2)$ are independent

First iteration

Suppose

$$\mathbf{v}^0 = \mu_0 \mathbf{v} + \mathbf{g}^0, \quad \text{with } \mathbf{g}_0 \sim \mathcal{N}(0, \sigma_0^2 \mathbf{I}_n)$$

for some constants μ_0, σ_0 . Then

$$\mathbf{v}^1 = \lambda \frac{\langle \mathbf{v}, \hat{\mathbf{v}}^0 \rangle}{n} \mathbf{v} + \mathbf{W} \hat{\mathbf{v}}^0, \quad \hat{\mathbf{v}}^0 = f_0(\mathbf{v}_0)$$

► Signal term:

$$\frac{\lambda \langle \mathbf{v}, \hat{\mathbf{v}}^0 \rangle}{n} = \frac{\lambda}{n} \sum_{i=1}^n v_i f_0(v_i^0) \rightarrow \mathbb{E}\{V f_0(\mu_0 V + G_0)\} =: \mu_1$$

where $V \sim P_V$ and $G_0 \sim \mathcal{N}(0, \sigma_0^2)$ are independent

► Empirical distribution of $\mathbf{W} f_0(\mathbf{v}^0) \rightarrow \mathcal{N}(0, \sigma_1^2)$ where

$$\sigma_1^2 := \lim_{n \rightarrow \infty} \frac{\|\mathbf{v}^0\|^2}{n} = \mathbb{E}\{f_0(\mu_0 V + G_0)^2\}$$

⇒ Empirical dist. $\nu_n(\mathbf{v}^1) \rightarrow \mu_1 V + G_1$, with $G_1 \sim \mathcal{N}(0, \sigma_1^2)$

Subsequent iterations

Recall the AMP iteration:

$$\hat{\mathbf{v}}^t = f_t(\mathbf{v}^t), \quad \mathbf{v}^{t+1} = \mathbf{A} \hat{\mathbf{v}}^t - b_t \hat{\mathbf{v}}^{t-1}, \quad b_t = \frac{1}{n} \sum_{i=1}^n f_t'(v_i^t)$$

Suppose $\nu_n(\mathbf{v}^t) \rightarrow \mu_t V + G_t$, with $G_t \sim N(0, \sigma_t^2)$

$$\mathbf{v}^{t+1} = \underbrace{\lambda \frac{\langle \mathbf{v}, \hat{\mathbf{v}}^t \rangle}{n}}_{\approx \mu_{t+1} \mathbf{v}} + \underbrace{\mathbf{W} \hat{\mathbf{v}}^t - b_t \hat{\mathbf{v}}^{t-1}}_{\approx N(0, \sigma_{t+1}^2 \mathbf{I}_n)}$$

State evolution recursion

$$\mu_{t+1} = \lambda \mathbb{E}[V f_t(\mu_t V + G_t)], \quad \sigma_{t+1}^2 = \mathbb{E}[f_t(\mu_t V + G_t)^2]$$

where $G_t \sim N(0, \sigma_t^2)$ indep. of $V \sim P_V$. Initialize with μ_0, σ_0

Main result for rank-one AMP

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}, \quad \mathbf{W} \sim \text{GOE}(n)$$

Assumptions:

- ▶ Functions f_t Lipschitz, for $t \geq 1$
- ▶ Initialization \mathbf{v}^0 is independent of \mathbf{W}

Theorem

For $t \geq 1$, and any pseudo-Lipschitz function $\phi : \mathbb{R} \rightarrow \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \phi(v_i, v_i^t) = \mathbb{E}\{\phi(V, \mu_t V + G_t)\} \text{ almost surely}$$

where $G_t \sim \text{N}(0, \sigma_t^2)$ independent of $V \sim P_V$

Main result for rank-one AMP

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}, \quad \mathbf{W} \sim \text{GOE}(n)$$

Assumptions:

- ▶ Functions f_t Lipschitz, for $t \geq 1$
- ▶ Initialization \mathbf{v}^0 is independent of \mathbf{W}

Theorem

For $t \geq 1$, and any pseudo-Lipschitz function $\phi : \mathbb{R} \rightarrow \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \phi(v_i, v_i^t) = \mathbb{E}\{\phi(V, \mu_t V + G_t)\} \text{ almost surely}$$

where $G_t \sim \text{N}(0, \sigma_t^2)$ independent of $V \sim P_V$

$$\text{Implies } \lim_{n \rightarrow \infty} \frac{\langle \mathbf{v}, \hat{\mathbf{v}}^t \rangle}{n} = \mathbb{E}\{V f_t(\mu_t V + G_t)\}, \text{ for each } t \geq 1$$

Choosing f_t

AMP result says $\mathbf{v}^t \stackrel{d}{\approx} \mu_t V + G_t$, with $G_t \sim N(0, \sigma_t^2)$

$$\mu_{t+1} = \lambda \mathbb{E}[V f_t(\mu_t V + G_t)], \quad \sigma_{t+1}^2 = \mathbb{E}[f_t(\mu_t V + G_t)^2]$$

- ▶ Given μ_t, σ_t , want to choose f_t to maximize

$$\gamma_{t+1} := \frac{\mu_{t+1}^2}{\sigma_{t+1}^2}$$

Choosing f_t

AMP result says $\mathbf{v}^t \stackrel{d}{\approx} \mu_t V + G_t$, with $G_t \sim N(0, \sigma_t^2)$

$$\mu_{t+1} = \lambda \mathbb{E}[V f_t(\mu_t V + G_t)], \quad \sigma_{t+1}^2 = \mathbb{E}[f_t(\mu_t V + G_t)^2]$$

- ▶ Given μ_t, σ_t , want to choose f_t to maximize

$$\gamma_{t+1} := \frac{\mu_{t+1}^2}{\sigma_{t+1}^2}$$

- ▶ If we know the prior distribution $V \sim P_V$, optimal choice is

$$f_t^*(s) = \mathbb{E}\{V \mid \mu_t V + \sigma_t G_t = s\}$$

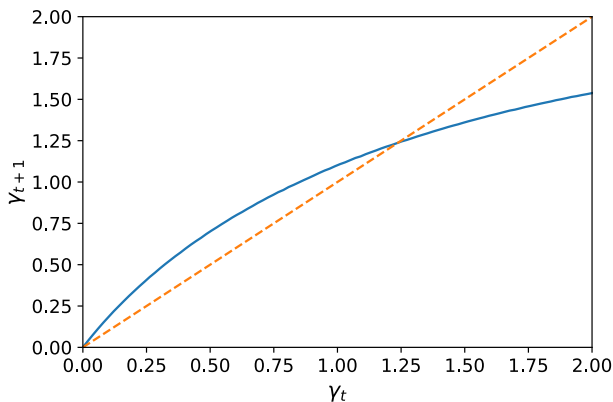
- ▶ State evolution with Bayes-optimal f_t^*

$$\gamma_{t+1} = \lambda^2 \{1 - \text{mmse}(\gamma_t)\}$$

where $\text{mmse}(\gamma) = \mathbb{E}\{(V - \mathbb{E}\{V \mid V + \sqrt{\gamma}G = s\})^2\}$

Fixed point of state evolution

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}, \quad P_V \sim \text{Unif}\{1, -1\}, \quad \lambda = \sqrt{2}$$

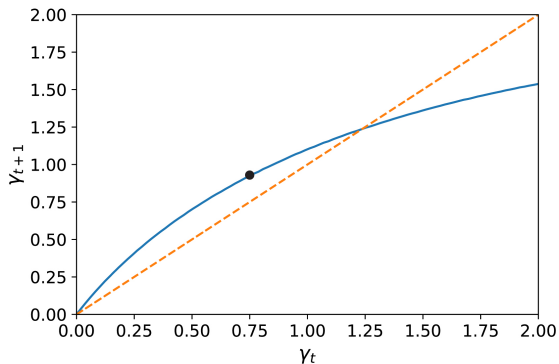


Recall $\mathbf{v}^0 \stackrel{d}{=} \mu_0 V + \sigma_0 G$

- ▶ $\gamma_t = 0$ is an (unstable) fixed point: if $\gamma_0 = \frac{\mu_0^2}{\sigma_0^2} = 0$ then $\gamma_t = 0$ for all t !

Fixed point of state evolution

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}, \quad P_V \sim \text{Unif}\{1, -1\}, \quad \lambda = \sqrt{2}$$

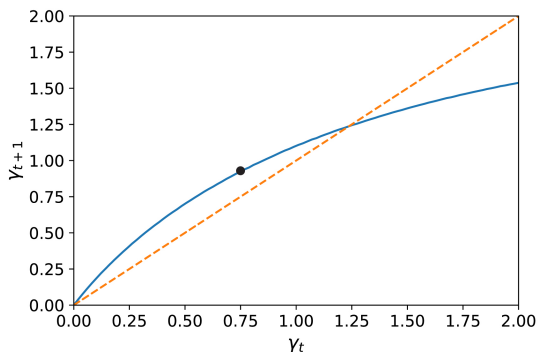


Recall $\mathbf{v}^0 \stackrel{d}{=} \mu_0 \mathbf{V} + \sigma_0 \mathbf{G}$

- ▶ If $\gamma_0 \neq 0$, that is, \mathbf{v}^0 correlated with \mathbf{v} , AMP converges to the 'good' fixed point

Correlated initialization

Assuming correlated initialization often not realistic

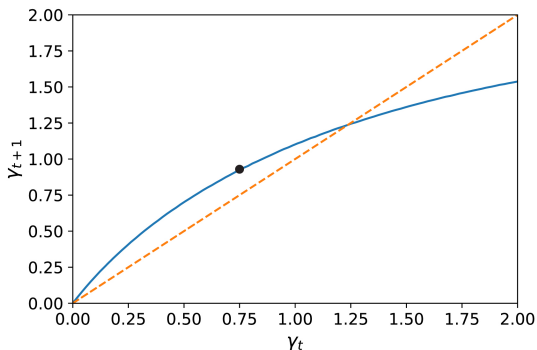


Natural initializer: $\hat{\varphi}$ the principal eigenvector of $\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}$

$$\text{Correlation : } \frac{|\langle \hat{\varphi}, \mathbf{v} \rangle|}{\|\hat{\varphi}\| \|\mathbf{v}\|} \rightarrow \begin{cases} \sqrt{1 - \lambda^{-2}}, & \text{if } \lambda > 1, \\ 0, & \text{if } \lambda \in (0, 1] \end{cases}$$

Spectral initialization

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}$$



- ▶ Standard AMP theory assumes $\hat{\mathbf{v}}^0$ is independent of \mathbf{A}
- ▶ Spectral initialization requires special analysis [Montanari-Venkataramanan '21]
- ▶ With spectral initialization $\gamma_0 = 1 - \lambda^{-2}$, if $\lambda \geq 1$

Example: Two-point mixture

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}$$

$$P_V = \varepsilon \delta_{a_+} + (1 - \varepsilon) \delta_{a_-} \quad a_+ = \sqrt{\frac{1 - \varepsilon}{\varepsilon}} \quad a_- = -\sqrt{\frac{\varepsilon}{1 - \varepsilon}}.$$

Run AMP with spectral initialization

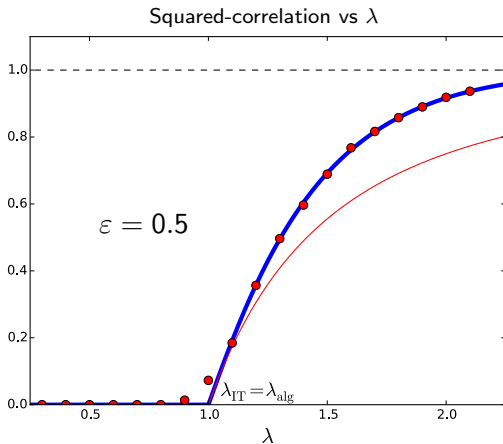
$$\gamma_{t+1} = \lambda^2 \{1 - \text{mmse}(\gamma_t)\}$$

- ▶ Can determine fixed point $\lim_{t \rightarrow \infty} \gamma_t$
- ▶ Initialization $\gamma_0 = 1 - \lambda^{-2}$

Example: Two-point mixture

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}$$

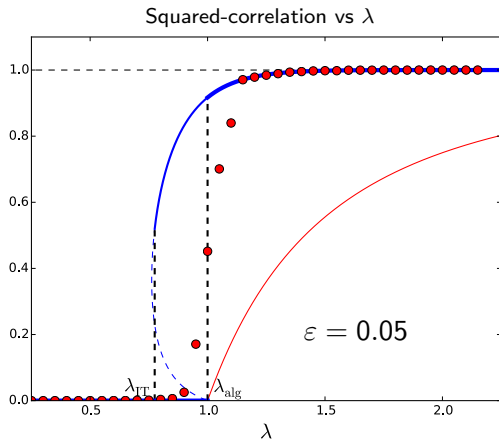
$$P_V = \varepsilon \delta_{a_+} + (1 - \varepsilon) \delta_{a_-} \quad a_+ = \sqrt{\frac{1 - \varepsilon}{\varepsilon}} \quad a_- = -\sqrt{\frac{\varepsilon}{1 - \varepsilon}}$$



Example: Two-point mixture

$$\mathbf{A} = \frac{\lambda}{n} \mathbf{v} \mathbf{v}^T + \mathbf{W}$$

$$P_V = \varepsilon \delta_{a_+} + (1 - \varepsilon) \delta_{a_-} \quad a_+ = \sqrt{\frac{1 - \varepsilon}{\varepsilon}} \quad a_- = -\sqrt{\frac{\varepsilon}{1 - \varepsilon}}$$



Rank- k matrix estimation

Can generalize AMP to estimate rank- k signals

Symmetric:

$$\mathbf{A} = \sum_{i=1}^k \lambda_i \mathbf{v}_i \mathbf{v}_i^T + \mathbf{W} \quad \in \mathbb{R}^{n \times n}$$

GOAL: To estimate the vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ from \mathbf{A}

Non-symmetric:

$$\mathbf{A} = \sum_{i=1}^k \lambda_i \mathbf{u}_i \mathbf{v}_i^T + \mathbf{W} \quad \in \mathbb{R}^{m \times n}$$

GOAL: Estimate the singular vectors $\mathbf{u}_1, \dots, \mathbf{u}_k$ and $\mathbf{v}_1, \dots, \mathbf{v}_k$

Generalizations

Abstract AMP can be generalized to:

1. **Matrix-valued** iterates

$$\mathbf{m}^t = f_t(\mathbf{h}^t), \quad \mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - b_t \mathbf{m}^{t-1}$$

with \mathbf{h}^t , \mathbf{m}^t being $n \times k$ matrices (k is fixed)

Used for analyzing AMP for rank- k matrix estimation

Generalizations

Abstract AMP can be generalized to:

1. **Matrix-valued** iterates

$$\mathbf{m}^t = f_t(\mathbf{h}^t), \quad \mathbf{h}^{t+1} = \mathbf{W} \mathbf{m}^t - \mathbf{b}_t \mathbf{m}^{t-1}$$

with \mathbf{h}^t , \mathbf{m}^t being $n \times k$ matrices (k is fixed)

Used for analyzing AMP for rank- k matrix estimation

2. **Non-symmetric** i.i.d. Gaussian matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$. AMP defined via pairs of functions f_t, g_t for $t \geq 1$:

$$\begin{aligned} \mathbf{e}^t &= \mathbf{A} f_t(\mathbf{h}^t) - \mathbf{b}_t g_{t-1}(\mathbf{e}^{t-1}) \\ \mathbf{h}^{t+1} &= \mathbf{A}^\top g_t(\mathbf{e}^t) - \mathbf{c}_t f_t(\mathbf{h}^t) \end{aligned}$$

- ▶ Empirical distributions of $\mathbf{e}^t \in \mathbb{R}^n$ and $\mathbf{h}^{t+1} \in \mathbb{R}^d$ converge to zero-mean Gaussians with variances given by SE
- ▶ Used for analyzing AMP for linear models

Finite sample analysis of AMP

State evolution (SE) results in the large-but-finite n regime can be established under stronger assumptions

- ▶ [Rush, Venkataramanan '18]: Concentration inequality for AMP performance showing validity of SE for $\sim \frac{\log n}{\log \log n}$ iterations
- ▶ [Li, Wei '22], [Li, Fan, Wei '23]: Refined finite-sample SE for rank-1 AMP showing SE valid for $O\left(\frac{n}{\text{polylog}(n)}\right)$ iterations

Reference:

O. Feng, R. Venkataramanan, C. Rush, R. Samworth,
A unifying tutorial on Approximate Message Passing,
Foundations and Trends in Machine Learning, 2022
<https://arxiv.org/abs/2105.02180>

Free download during ISIT at

<https://nowpublishers.com/conference/ISIT2023>
Access code: ISIT-9502