

# Local convexity of the TAP free energy and AMP convergence for Z2-synchronization

Song Mei

UC Berkeley

September 15, 2021

Joint work with Michael Celentano and Zhou Fan

# Motivation

- ▶ Variational Bayesian inference:
  - ▶ Optimizing **variational free energy** (VB, Bethe, TAP).
  - ▶ **Iterative algorithms** (BP, EP, MP, AMP).
- ▶ These methods have been established for **more than 20 years**.  
Have been written into softwares and **works well in practice**.
- ▶ Theoretical challenges:
  - ▶ Minimizer of free energy **approximates** Bayesian posterior mean?
  - ▶ **Landscape** of the free energy?
  - ▶ **Convergence** of iterative algorithms?

# Today

- ▶ High dimensional statistical model:  
 $\mathbb{Z}_2$  synchronization in the weak signal regime.
- ▶ Landscape of the TAP free energy.
- ▶ Convergence of iterative algorithms (NGD, AMP).

## $\mathbb{Z}_2$ synchronization

- ▶ Signal:

$$\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{Z}_2^n, \quad x_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

- ▶ Observation: for  $1 \leq i < j \leq n$

$$Y_{ij} = \frac{\lambda}{n} x_i x_j + W_{ij}.$$

- ▶ Noise  $W_{ij} \sim \mathcal{N}(0, 1/n)$ .
- ▶ SNR  $\lambda \in [0, \infty)$  fixed, dimension  $n \rightarrow \infty$ .
- ▶ In matrix notation:

$$\mathbf{Y} = \frac{\lambda}{n} \mathbf{x} \mathbf{x}^\top + \mathbf{W}.$$

- ▶ Task: given  $\mathbf{Y} = (Y_{ij})$ , estimate  $\mathbf{x}$  (or say  $\mathbf{X} = \mathbf{x} \mathbf{x}^\top$ ).

## $\mathbb{Z}_2$ synchronization

- ▶ Signal:

$$\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{Z}_2^n, \quad x_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

- ▶ Observation: for  $1 \leq i < j \leq n$

$$Y_{ij} = \frac{\lambda}{n} x_i x_j + W_{ij}.$$

- ▶ Noise  $W_{ij} \sim \mathcal{N}(0, 1/n)$ .
- ▶ SNR  $\lambda \in [0, \infty)$  fixed, dimension  $n \rightarrow \infty$ .
- ▶ In matrix notation:

$$\mathbf{Y} = \frac{\lambda}{n} \mathbf{x} \mathbf{x}^\top + \mathbf{W}.$$

- ▶ Task: given  $\mathbf{Y} = (Y_{ij})$ , estimate  $\mathbf{x}$  (or say  $\mathbf{X} = \mathbf{x} \mathbf{x}^\top$ ).

## $\mathbb{Z}_2$ synchronization

- ▶ Signal:

$$\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{Z}_2^n, \quad x_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

- ▶ Observation: for  $1 \leq i < j \leq n$

$$Y_{ij} = \frac{\lambda}{n} x_i x_j + W_{ij}.$$

- ▶ Noise  $W_{ij} \sim \mathcal{N}(0, 1/n)$ .
- ▶ SNR  $\lambda \in [0, \infty)$  fixed, dimension  $n \rightarrow \infty$ .
- ▶ In matrix notation:

$$\mathbf{Y} = \frac{\lambda}{n} \mathbf{x} \mathbf{x}^\top + \mathbf{W}.$$

- ▶ Task: given  $\mathbf{Y} = (Y_{ij})$ , estimate  $\mathbf{x}$  (or say  $\mathbf{X} = \mathbf{x} \mathbf{x}^\top$ ).

## $\mathbb{Z}_2$ synchronization

- ▶ Signal:

$$\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{Z}_2^n, \quad x_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

- ▶ Observation: for  $1 \leq i < j \leq n$

$$Y_{ij} = \frac{\lambda}{n} x_i x_j + W_{ij}.$$

- ▶ Noise  $W_{ij} \sim \mathcal{N}(0, 1/n)$ .
- ▶ SNR  $\lambda \in [0, \infty)$  fixed, dimension  $n \rightarrow \infty$ .
- ▶ In matrix notation:

$$\mathbf{Y} = \frac{\lambda}{n} \mathbf{x} \mathbf{x}^\top + \mathbf{W}.$$

- ▶ Task: given  $\mathbf{Y} = (Y_{ij})$ , estimate  $\mathbf{x}$  (or say  $\mathbf{X} = \mathbf{x} \mathbf{x}^\top$ ).

## $\mathbb{Z}_2$ synchronization

- ▶ Signal:

$$\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{Z}_2^n, \quad x_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

- ▶ Observation: for  $1 \leq i < j \leq n$

$$Y_{ij} = \frac{\lambda}{n} x_i x_j + W_{ij}.$$

- ▶ Noise  $W_{ij} \sim \mathcal{N}(0, 1/n)$ .
- ▶ SNR  $\lambda \in [0, \infty)$  fixed, dimension  $n \rightarrow \infty$ .
- ▶ In matrix notation:

$$\mathbf{Y} = \frac{\lambda}{n} \mathbf{x} \mathbf{x}^\top + \mathbf{W}.$$

- ▶ Task: given  $\mathbf{Y} = (Y_{ij})$ , estimate  $\mathbf{x}$  (or say  $\mathbf{X} = \mathbf{x} \mathbf{x}^\top$ ).



## $\mathbb{Z}_2$ synchronization

- ▶ Signal:

$$\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{Z}_2^n, \quad x_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

- ▶ Observation: for  $1 \leq i < j \leq n$

$$Y_{ij} = \frac{\lambda}{n} x_i x_j + W_{ij}.$$

- ▶ Noise  $W_{ij} \sim \mathcal{N}(0, 1/n)$ .
- ▶ SNR  $\lambda \in [0, \infty)$  fixed, dimension  $n \rightarrow \infty$ .
- ▶ In matrix notation:

$$\mathbf{Y} = \frac{\lambda}{n} \mathbf{x} \mathbf{x}^\top + \mathbf{W}.$$

- ▶ Task: given  $\mathbf{Y} = (Y_{ij})$ , estimate  $\mathbf{x}$  (or say  $\mathbf{X} = \mathbf{x} \mathbf{x}^\top$ ).

# Bayes estimation in $\mathbb{Z}_2$ synchronization

- ▶ Settings:

$$\mathbf{x} \sim \text{Unif}(\mathbb{Z}_2^n), \quad \mathbf{Y} = (\lambda/n)\mathbf{x}\mathbf{x}^\top + \mathbf{W}.$$

- ▶ Estimate  $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$  with loss:

$$\ell(\mathbf{X}, \widehat{\mathbf{X}}) = (1/n^2)\|\mathbf{X} - \widehat{\mathbf{X}}\|_F^2.$$

- ▶ For  $\lambda < 1$ , impossible.
- ▶ For  $\lambda > 1$ , possible and efficient, e.g., spectral estimator (BBAP phase transition).
- ▶ Optimal estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{Y}].$$

# Bayes estimation in $\mathbb{Z}_2$ synchronization

- ▶ Settings:

$$\mathbf{x} \sim \text{Unif}(\mathbb{Z}_2^n), \quad \mathbf{Y} = (\lambda/n)\mathbf{x}\mathbf{x}^\top + \mathbf{W}.$$

- ▶ Estimate  $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$  with loss:

$$\ell(\mathbf{X}, \widehat{\mathbf{X}}) = (1/n^2)\|\mathbf{X} - \widehat{\mathbf{X}}\|_F^2.$$

- ▶ For  $\lambda < 1$ , impossible.
- ▶ For  $\lambda > 1$ , possible and efficient, e.g., spectral estimator (BBAP phase transition).
- ▶ Optimal estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{Y}].$$

# Bayes estimation in $\mathbb{Z}_2$ synchronization

- ▶ Settings:

$$\mathbf{x} \sim \text{Unif}(\mathbb{Z}_2^n), \quad \mathbf{Y} = (\lambda/n)\mathbf{x}\mathbf{x}^\top + \mathbf{W}.$$

- ▶ Estimate  $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$  with loss:

$$\ell(\mathbf{X}, \widehat{\mathbf{X}}) = (1/n^2)\|\mathbf{X} - \widehat{\mathbf{X}}\|_F^2.$$

- ▶ For  $\lambda < 1$ , impossible.
- ▶ For  $\lambda > 1$ , possible and efficient, e.g., spectral estimator (BBAP phase transition).
- ▶ Optimal estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{Y}].$$

# Bayes estimation in $\mathbb{Z}_2$ synchronization

- ▶ Settings:

$$\mathbf{x} \sim \text{Unif}(\mathbb{Z}_2^n), \quad \mathbf{Y} = (\lambda/n)\mathbf{x}\mathbf{x}^\top + \mathbf{W}.$$

- ▶ Estimate  $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$  with loss:

$$\ell(\mathbf{X}, \widehat{\mathbf{X}}) = (1/n^2)\|\mathbf{X} - \widehat{\mathbf{X}}\|_F^2.$$

- ▶ For  $\lambda < 1$ , impossible.
- ▶ For  $\lambda > 1$ , possible and efficient, e.g., spectral estimator (BBAP phase transition).
- ▶ Optimal estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{Y}].$$

# Bayes estimation in $\mathbb{Z}_2$ synchronization

- ▶ Settings:

$$\mathbf{x} \sim \text{Unif}(\mathbb{Z}_2^n), \quad \mathbf{Y} = (\lambda/n)\mathbf{x}\mathbf{x}^\top + \mathbf{W}.$$

- ▶ Estimate  $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$  with loss:

$$\ell(\mathbf{X}, \widehat{\mathbf{X}}) = (1/n^2)\|\mathbf{X} - \widehat{\mathbf{X}}\|_F^2.$$

- ▶ For  $\lambda < 1$ , impossible.
- ▶ For  $\lambda > 1$ , possible and efficient, e.g., spectral estimator (BBAP phase transition).
- ▶ Optimal estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{Y}].$$

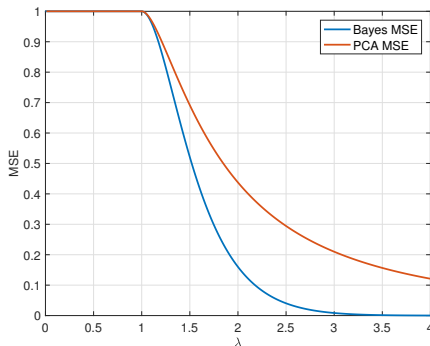
# Bayes estimation in $\mathbb{Z}_2$ synchronization

► Settings:

$$\mathbf{x} \sim \text{Unif}(\mathbb{Z}_2^n), \quad \mathbf{Y} = (\lambda/n)\mathbf{x}\mathbf{x}^\top + \mathbf{W}.$$

► Risk:

$$\text{MSE}_\lambda(\widehat{\mathbf{X}}) = (1/n^2)\mathbb{E}[\|\mathbf{x}\mathbf{x}^\top - \widehat{\mathbf{X}}\|_F^2].$$



# Compute the Bayesian estimator

- ▶ The Bayesian estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top} | \mathbf{Y}] = \sum_{\sigma \in \mathbb{Z}_2^n} \sigma\sigma^{\top} p(\sigma | \mathbf{Y}).$$

- ▶ The posterior distribution: (relationship with SK measure)

$$p(\sigma | \mathbf{Y}) = \frac{1}{Z} \exp\{\lambda \langle \sigma, \mathbf{Y} \sigma \rangle / 2\}.$$

- ▶ Two viewpoints in variational inference:

- ▶ Variational free energies (TAP free energy).
- ▶ Iterative algorithms (NGD, AMP).



# Compute the Bayesian estimator

- ▶ The Bayesian estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top} | \mathbf{Y}] = \sum_{\boldsymbol{\sigma} \in \mathbb{Z}_2^n} \boldsymbol{\sigma}\boldsymbol{\sigma}^{\top} p(\boldsymbol{\sigma} | \mathbf{Y}).$$

- ▶ The posterior distribution: (relationship with SK measure)

$$p(\boldsymbol{\sigma} | \mathbf{Y}) = \frac{1}{Z} \exp\{\lambda \langle \boldsymbol{\sigma}, \mathbf{Y}\boldsymbol{\sigma} \rangle / 2\}.$$

- ▶ Two viewpoints in variational inference:

• Variational free energies (TAP free energy).

• Iterative algorithms (NGD, AMP).

# Compute the Bayesian estimator

- ▶ The Bayesian estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top} | \mathbf{Y}] = \sum_{\boldsymbol{\sigma} \in \mathbb{Z}_2^n} \boldsymbol{\sigma}\boldsymbol{\sigma}^{\top} p(\boldsymbol{\sigma} | \mathbf{Y}).$$

- ▶ The posterior distribution: (relationship with SK measure)

$$p(\boldsymbol{\sigma} | \mathbf{Y}) = \frac{1}{Z} \exp\{\lambda \langle \boldsymbol{\sigma}, \mathbf{Y}\boldsymbol{\sigma} \rangle / 2\}.$$

- ▶ Two viewpoints in variational inference:
  - ▶ Variational free energies (TAP free energy).
  - ▶ Iterative algorithms (NGD, AMP).

# Compute the Bayesian estimator

- ▶ The Bayesian estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top} | \mathbf{Y}] = \sum_{\boldsymbol{\sigma} \in \mathbb{Z}_2^n} \boldsymbol{\sigma}\boldsymbol{\sigma}^{\top} p(\boldsymbol{\sigma} | \mathbf{Y}).$$

- ▶ The posterior distribution: (relationship with SK measure)

$$p(\boldsymbol{\sigma} | \mathbf{Y}) = \frac{1}{Z} \exp\{\lambda \langle \boldsymbol{\sigma}, \mathbf{Y}\boldsymbol{\sigma} \rangle / 2\}.$$

- ▶ Two viewpoints in variational inference:
  - ▶ Variational free energies (TAP free energy).
  - ▶ Iterative algorithms (NGD, AMP).

# Compute the Bayesian estimator

- ▶ The Bayesian estimator:

$$\widehat{\mathbf{X}}_{\text{Bayes}} = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top} | \mathbf{Y}] = \sum_{\boldsymbol{\sigma} \in \mathbb{Z}_2^n} \boldsymbol{\sigma}\boldsymbol{\sigma}^{\top} p(\boldsymbol{\sigma} | \mathbf{Y}).$$

- ▶ The posterior distribution: (relationship with SK measure)

$$p(\boldsymbol{\sigma} | \mathbf{Y}) = \frac{1}{Z} \exp\{\lambda \langle \boldsymbol{\sigma}, \mathbf{Y}\boldsymbol{\sigma} \rangle / 2\}.$$

- ▶ Two viewpoints in variational inference:
  - ▶ Variational free energies (TAP free energy).
  - ▶ Iterative algorithms (NGD, AMP).

# Variational free energies

- ▶ Find a function  $\mathcal{F} : [-1, 1]^n \rightarrow \mathbb{R}$ , such that

$$\hat{\mathbf{m}} \equiv \arg \min_{\mathbf{m}} \mathcal{F}(\mathbf{m}) \longleftrightarrow \widehat{\mathbf{X}}_{\text{Bayes}}.$$

- ▶ The **mean field variational Bayes** (the KL minimization)

$$\mathcal{F}_{\text{MF}}(\mathbf{m}) \equiv - \sum_{i=1}^n h(m_i) - \lambda \langle \mathbf{m}, \mathbf{Y} \mathbf{m} \rangle / 2,$$

where  $h(m) = -\frac{1-m}{2} \log(\frac{1-m}{2}) - \frac{1+m}{2} \log(\frac{1+m}{2})$ .

- ▶ The **TAP free energy** (Thouless, Anderson, and Palmer (1977))

$$\mathcal{F}_{\text{TAP}}(\mathbf{m}) \equiv \underbrace{- \sum_{i=1}^n h(m_i) - \frac{\lambda}{2} \langle \mathbf{m}, \mathbf{Y} \mathbf{m} \rangle}_{\mathcal{F}_{\text{MF}}} - \underbrace{\frac{n\lambda^2}{4} \left[ 1 - \frac{\|\mathbf{m}\|_2^2}{n} \right]^2}_{\text{Onsager's correction term}}.$$

# Variational free energies

- ▶ Find a function  $\mathcal{F} : [-1, 1]^n \rightarrow \mathbb{R}$ , such that

$$\hat{\mathbf{m}} \equiv \arg \min_{\mathbf{m}} \mathcal{F}(\mathbf{m}) \longleftrightarrow \widehat{\mathbf{X}}_{\text{Bayes}}.$$

- ▶ The **mean field variational Bayes** (the KL minimization)

$$\mathcal{F}_{\text{MF}}(\mathbf{m}) \equiv - \sum_{i=1}^n h(m_i) - \lambda \langle \mathbf{m}, \mathbf{Y} \mathbf{m} \rangle / 2,$$

where  $h(m) = -\frac{1-m}{2} \log(\frac{1-m}{2}) - \frac{1+m}{2} \log(\frac{1+m}{2})$ .

- ▶ The **TAP free energy** (Thouless, Anderson, and Palmer (1977))

$$\mathcal{F}_{\text{TAP}}(\mathbf{m}) \equiv \underbrace{- \sum_{i=1}^n h(m_i) - \frac{\lambda}{2} \langle \mathbf{m}, \mathbf{Y} \mathbf{m} \rangle}_{\mathcal{F}_{\text{MF}}} - \underbrace{\frac{n\lambda^2}{4} \left[ 1 - \frac{\|\mathbf{m}\|_2^2}{n} \right]^2}_{\text{Onsager's correction term}}.$$

# Variational free energies

- ▶ Find a function  $\mathcal{F} : [-1, 1]^n \rightarrow \mathbb{R}$ , such that

$$\hat{\mathbf{m}} \equiv \arg \min_{\mathbf{m}} \mathcal{F}(\mathbf{m}) \longleftrightarrow \widehat{\mathbf{X}}_{\text{Bayes}}.$$

- ▶ The **mean field variational Bayes** (the KL minimization)

$$\mathcal{F}_{\text{MF}}(\mathbf{m}) \equiv - \sum_{i=1}^n h(m_i) - \lambda \langle \mathbf{m}, \mathbf{Y} \mathbf{m} \rangle / 2,$$

where  $h(m) = -\frac{1-m}{2} \log(\frac{1-m}{2}) - \frac{1+m}{2} \log(\frac{1+m}{2})$ .

- ▶ The **TAP free energy** (Thouless, Anderson, and Palmer (1977))

$$\mathcal{F}_{\text{TAP}}(\mathbf{m}) \equiv \underbrace{- \sum_{i=1}^n h(m_i) - \frac{\lambda}{2} \langle \mathbf{m}, \mathbf{Y} \mathbf{m} \rangle}_{\mathcal{F}_{\text{MF}}} - \underbrace{\frac{n\lambda^2}{4} \left[ 1 - \frac{\|\mathbf{m}\|_2^2}{n} \right]^2}_{\text{Onsager's correction term}}.$$

# Variational free energies

- ▶ The TAP free energy (Thouless, Anderson, and Palmer (1977))

$$\mathcal{F}_{\text{TAP}}(\mathbf{m}) \equiv \underbrace{-\sum_{i=1}^n h(m_i) - \frac{\lambda}{2} \langle \mathbf{m}, \mathbf{Y} \mathbf{m} \rangle}_{\mathcal{F}_{\text{MF}}} - \underbrace{\frac{n\lambda^2}{4} \left[ 1 - \frac{\|\mathbf{m}\|_2^2}{n} \right]^2}_{\text{Onsager's correction term}}.$$

- ▶ This is **non-convex**.

Q: Nice landscape? Global minimizer corresponds to  $\widehat{\mathbf{X}}_{\text{Bayes}}$ ?



# Variational free energies

- ▶ The TAP free energy (Thouless, Anderson, and Palmer (1977))

$$\mathcal{F}_{\text{TAP}}(\mathbf{m}) \equiv \underbrace{-\sum_{i=1}^n h(m_i) - \frac{\lambda}{2} \langle \mathbf{m}, \mathbf{Y} \mathbf{m} \rangle}_{\mathcal{F}_{\text{MF}}} - \underbrace{\frac{n\lambda^2}{4} \left[ 1 - \frac{\|\mathbf{m}\|_2^2}{n} \right]^2}_{\text{Onsager's correction term}}.$$

- ▶ This is **non-convex**.

Q: Nice landscape? Global minimizer corresponds to  $\widehat{\mathbf{X}}_{\text{Bayes}}$ ?

# Variational free energies

- ▶ The TAP free energy (Thouless, Anderson, and Palmer (1977))

$$\mathcal{F}_{\text{TAP}}(\mathbf{m}) \equiv \underbrace{-\sum_{i=1}^n h(m_i) - \frac{\lambda}{2} \langle \mathbf{m}, \mathbf{Y} \mathbf{m} \rangle}_{\mathcal{F}_{\text{MF}}} - \underbrace{\frac{n\lambda^2}{4} \left[ 1 - \frac{\|\mathbf{m}\|_2^2}{n} \right]^2}_{\text{Onsager's correction term}}.$$

- ▶ This is **non-convex**.

Q: Nice landscape? Global minimizer corresponds to  $\widehat{\mathbf{X}}_{\text{Bayes}}$ ?

# Approximate message passing

- ▶ AMP iteration [Bolthausen, 2012] [Donoho, Maleki, Montanari, 2009]

$$\mathbf{m}^{k+1} = \tanh(\lambda \mathbf{Y} \mathbf{m}^k - \lambda^2 (1 - Q(\mathbf{m}^k)) \mathbf{m}^{k-1}).$$

- ▶ Fixed point is a stationary point of the TAP free energy

$$\mathbf{m}_* = \text{AMP}(\mathbf{m}_*, \mathbf{m}_*) \quad \longleftrightarrow \quad \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}_*) = \mathbf{0}.$$

- ▶ Not a descent algorithm on  $\mathcal{F}_{\text{TAP}}$ .

# Approximate message passing

- ▶ AMP iteration [Bolthausen, 2012] [Donoho, Maleki, Montanari, 2009]

$$\mathbf{m}^{k+1} = \tanh(\lambda \mathbf{Y} \mathbf{m}^k - \lambda^2 (1 - Q(\mathbf{m}^k)) \mathbf{m}^{k-1}).$$

- ▶ Fixed point is a stationary point of the TAP free energy

$$\mathbf{m}_* = \text{AMP}(\mathbf{m}_*, \mathbf{m}_*) \quad \longleftrightarrow \quad \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}_*) = \mathbf{0}.$$

- ▶ Not a descent algorithm on  $\mathcal{F}_{\text{TAP}}$ .

# Approximate message passing

- ▶ AMP iteration [Bolthausen, 2012] [Donoho, Maleki, Montanari, 2009]

$$\mathbf{m}^{k+1} = \tanh(\lambda \mathbf{Y} \mathbf{m}^k - \lambda^2 (1 - Q(\mathbf{m}^k)) \mathbf{m}^{k-1}).$$

- ▶ Fixed point is a stationary point of the TAP free energy

$$\mathbf{m}_\star = \text{AMP}(\mathbf{m}_\star, \mathbf{m}_\star) \quad \longleftrightarrow \quad \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}_\star) = \mathbf{0}.$$

- ▶ Not a descent algorithm on  $\mathcal{F}_{\text{TAP}}$ .

# Approximate message passing

- ▶ Convergence starting from spectral initialization [Montanari, Venkataramanan, 2017]

$$\lim_{n \rightarrow \infty} \|\mathbf{m}^k \mathbf{m}^k - \widehat{\mathbf{X}}_{\text{Bayes}}\|_F^2 / n^2 = \varepsilon(k) \rightarrow 0, \quad \text{as } k \rightarrow \infty.$$

- ▶ Analysis based on state evolution. Best known is  $k = o(\log(n) / \log \log n)$  [Rush, Venkataramanan, 2016].

Q: Convergence of AMP for fixed  $n$  as  $k \rightarrow \infty$ ?

# Approximate message passing

- ▶ Convergence starting from spectral initialization [Montanari, Venkataramanan, 2017]

$$\lim_{n \rightarrow \infty} \|\mathbf{m}^k \mathbf{m}^k - \widehat{\mathbf{X}}_{\text{Bayes}}\|_F^2 / n^2 = \varepsilon(k) \rightarrow 0, \quad \text{as } k \rightarrow \infty.$$

- ▶ Analysis based on state evolution. Best known is  $k = o(\log(n) / \log \log n)$  [Rush, Venkataramanan, 2016].

Q: Convergence of AMP for fixed  $n$  as  $k \rightarrow \infty$ ?

# Approximate message passing

- ▶ Convergence starting from spectral initialization [Montanari, Venkataramanan, 2017]

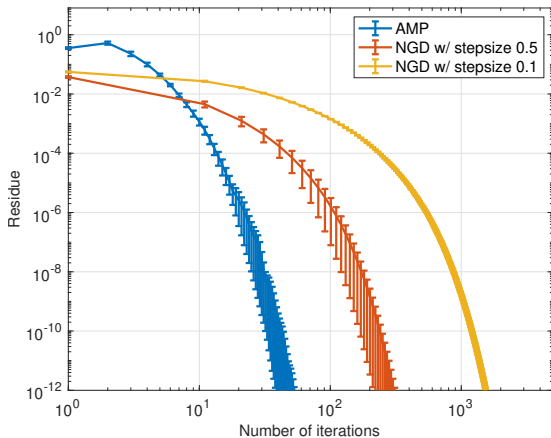
$$\lim_{n \rightarrow \infty} \|\mathbf{m}^k \mathbf{m}^k - \widehat{\mathbf{X}}_{\text{Bayes}}\|_F^2 / n^2 = \varepsilon(k) \rightarrow 0, \quad \text{as } k \rightarrow \infty.$$

- ▶ Analysis based on state evolution. Best known is  $k = o(\log(n) / \log \log n)$  [Rush, Venkataramanan, 2016].

Q: Convergence of AMP for fixed  $n$  as  $k \rightarrow \infty$ ?



# Numerical simulations



## Main results

# Landscape of the TAP free energy

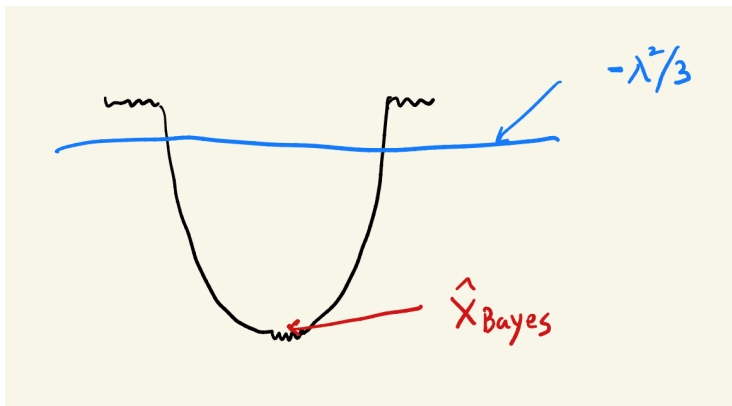
Theorem (Fan, Mei, Montanari, 2018)

Denote  $\mathcal{C}_{\lambda,n} = \{\mathbf{m} \in [-1, 1]^n : \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}) = \mathbf{0}, \mathcal{F}_{\text{TAP}}(\mathbf{m}) \leq -\lambda^2/3\}$ .  
There exists  $\lambda_0 > 0$ , such that for any  $\lambda > \lambda_0$ , we have

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[ \sup_{\mathbf{m} \in \mathcal{C}_{\lambda,n}} \frac{1}{n^2} \|\mathbf{m}\mathbf{m}^\top - \widehat{\mathbf{X}}_{\text{Bayes}}\|_F^2 \wedge 1 \right] = 0. \quad (1)$$

# Landscape of the TAP free energy

All the critical points below a threshold are close to the Bayesian estimator. [Fan, Mei, Montanari, 2018]



# Landscape of the TAP free energy

## Theorem (Celentano, Fan, Mei, 2021)

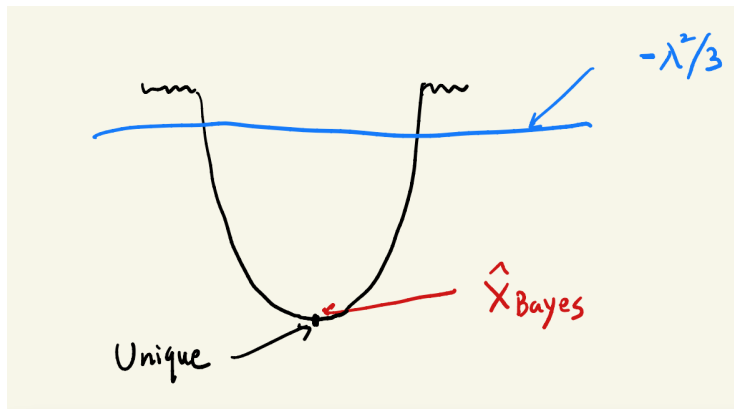
Fix any  $\lambda > 1$  (1 is the IT thresholds). For any sufficiently small  $\iota > 0$ , with probability approaching 1 as  $n \rightarrow \infty$ , there exists a unique critical point and unique local minimizer  $\mathbf{m}_*$  of  $\mathcal{F}_{\text{TAP}}(\mathbf{m})$  (up to  $\pm$  sign) such that

$$\frac{1}{n^2} \|\mathbf{m}_* \mathbf{m}_*^\top - \widehat{\mathbf{X}}_{\text{Bayes}}\|_{\text{F}}^2 < \iota. \quad (2)$$

Moreover,  $\mathcal{F}_{\text{TAP}}$  is strongly convex over  $(-1, 1)^n \cap \mathcal{B}_{\sqrt{\varepsilon n}}(\mathbf{m}_*)$  for some  $n$  independent constant  $\varepsilon$ .

# Landscape of the TAP free energy

There is a local strongly convex region near the Bayes estimator, which contains a unique local minimizer. [Celentano, Fan, Mei, 2021]



# Iterative algorithms

- ▶ Approximate message passing (AMP).

$$\begin{aligned} \mathbf{m}^k &= \tanh(\mathbf{h}^k), \\ \mathbf{h}^{k+1} &= \left( \lambda \mathbf{Y} \mathbf{m}^k - \lambda^2 [1 - Q(\mathbf{m}^k)] \mathbf{m}^{k-1} \right). \end{aligned}$$

- ▶ Natural gradient descent (NGD) or Bregman gradient descent.

$$\begin{aligned} \mathbf{m}^k &= \tanh(\mathbf{h}^k) \\ \mathbf{h}^{k+1} &= \mathbf{h}^k - \eta \mathbf{n} \cdot \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}^k) \\ &= (1 - \eta) \mathbf{h}^k + \eta \left( \lambda \mathbf{Y} \mathbf{m}^k - \lambda^2 [1 - Q(\mathbf{m}^k)] \mathbf{m}^k \right). \end{aligned}$$

# A hybrid algorithm converges

## A hybrid algorithm

Spectral initialization. Run AMP for  $T$  steps, and then continue to run NGD with stepsize  $\eta$ .

## Theorem (Celentano, Fan, Mei, 2021)

Fix any  $\lambda > 1$ . There exist  $\lambda$ -dependent constants  $C, \mu, \eta_0 > 0$  and  $T \geq 1$  such that with probability approaching 1 as  $n \rightarrow \infty$ , we have

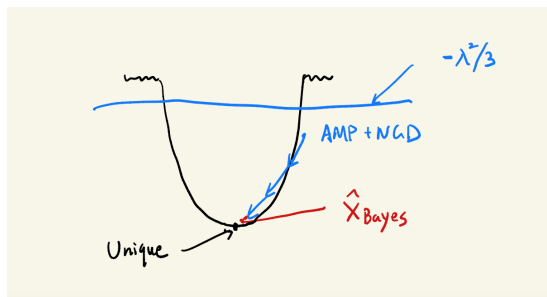
$$\begin{aligned}\mathcal{F}_{\text{TAP}}(\mathbf{m}^{T+k}) - \mathcal{F}_{\text{TAP}}(\pm \mathbf{m}_\star) &< C(1 - \mu\eta)^k, \\ \|\mathbf{m}^{T+k} - (\pm \mathbf{m}_\star)\|_2 &< C(1 - \mu\eta)^k \sqrt{n}.\end{aligned}$$

In particular,  $\lim_{k \rightarrow \infty} \mathbf{m}^{T+k} \in \{+\mathbf{m}_\star, -\mathbf{m}_\star\}$ .



# A hybrid algorithm converges

A hybrid algorithm converges. [Celentano, Fan, Mei, 2021]



- ▶ AMP drives the iterates into the local region.
- ▶ Local convergence of NGD is due to the relative smoothness and relative strong convexity with respect to the entropy function.

## How about AMP or NGD individually?

### Theorem (Celentano, Fan, Mei, 2021)

*There exists  $\lambda_0$ , such that for any  $\lambda > \lambda_0$ , with probability converging to 1 as  $n \rightarrow \infty$ , either **spectral initialized AMP** or **spectral initialized NGD** converges to the global minimizer  $m_*$ .*

The proof for the case  $\lambda \geq \lambda_0$  is much more easier than  $\lambda > 1$ .

## How about AMP or NGD individually?

Theorem (Celentano, Fan, Mei, 2021)

*There exists  $\lambda_0$ , such that for any  $\lambda > \lambda_0$ , with probability converging to 1 as  $n \rightarrow \infty$ , either **spectral initialized AMP** or **spectral initialized NGD** converges to the global minimizer  $m_*$ .*

The proof for the case  $\lambda \geq \lambda_0$  is much more easier than  $\lambda > 1$ .

# What do we know about AMP when $\lambda > 1$ ?

## Theorem (Celentano, Fan, Mei, 2021)

*For any  $\lambda > 1$ , with probability converging to 1 as  $n \rightarrow \infty$ ,  $m_*$  is an **asymptotically stable** fixed point of AMP. That means, there exists a neighborhood  $B(m_*, \delta)$  (the size of  $\delta$  may depend on  $n$ ), such that if AMP is initialized in the neighborhood, it converges to  $m_*$ .*

## Open problem

Show that spectral initialized AMP converges as long as  $\lambda > 1$ .

# What do we know about AMP when $\lambda > 1$ ?

## Theorem (Celentano, Fan, Mei, 2021)

*For any  $\lambda > 1$ , with probability converging to 1 as  $n \rightarrow \infty$ ,  $m_*$  is an **asymptotically stable** fixed point of AMP. That means, there exists a neighborhood  $B(m_*, \delta)$  (the size of  $\delta$  may depend on  $n$ ), such that if AMP is initialized in the neighborhood, it converges to  $m_*$ .*

## Open problem

Show that spectral initialized AMP converges as long as  $\lambda > 1$ .

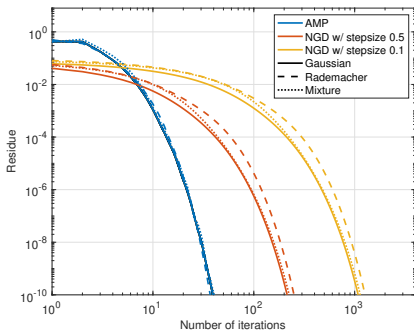
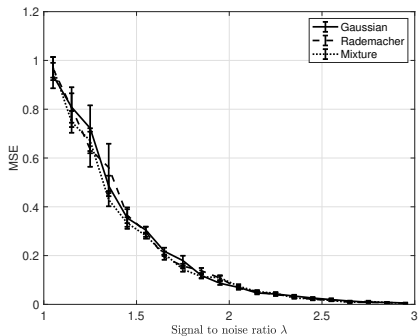
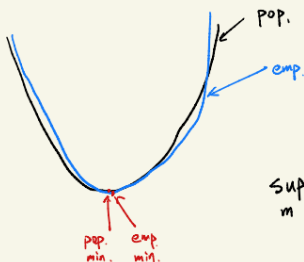


Figure: Universality with respect to the noise distribution.

# Technical challenge

When SNR is super large ( $\lambda = O(\log n)$ )

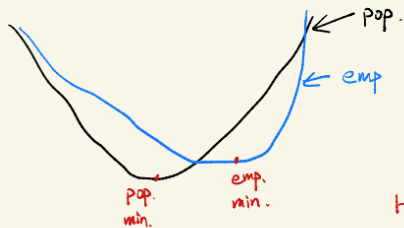


$$\sup_m \left| \nabla^k \text{emp}(m) - \nabla^k \text{pop}(m) \right| \rightarrow 0$$

$k=0, 1, 2.$

# Technical challenge

When SNR is  $O(1)$



No uniform conv.

emp. min. at a  
random region.

Hard to characterize.



## Technical tool 1: Kac-Rice formula

- ▶  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  a “sufficiently regular” random morse function.

$$\text{Crit}(T) = \#\{\mathbf{m} \in T : \nabla f(\mathbf{m}) = \mathbf{0}\}.$$

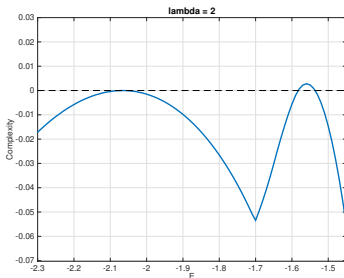
Kac-Rice formula:

$$\mathbb{E}[\text{Crit}(T)] = \int_T \mathbb{E}\left[|\det \nabla^2 f(\mathbf{m})| \mid \nabla f(\mathbf{m}) = \mathbf{0}\right] p_{\mathbf{m}}(\mathbf{0}) d\mathbf{m}.$$

- ▶ Technicality: determinant of Hessian  $\det \nabla^2 f(\mathbf{m})$ .
- ▶ Result: [Fan, Mei, Montanari, 2018]

$$\frac{1}{n} \log \mathbb{E}[\text{Crit}_n(U)] \leq \sup_{(q, \varphi, e) \in U} S_{\star}(q, \varphi, e) + o(1).$$

▶  $S_*(e) = \sup_{q, \varphi} S_*(q, \varphi, e).$



- ▶ All local min below some threshold are close to the global min.
- ▶ There could potentially be many local min near global.

## Technical tool 2: Conditional Gaussian comparison

- Analyze the following Gaussian process ...

$$\min_{\mathbf{m} \in \mathcal{B}(\mathbf{m}_*, \varepsilon \sqrt{n})} \min_{\mathbf{u} \in \mathbb{S}^{n-1}} \langle \mathbf{u}, \nabla^2 \mathcal{F}_{\text{TAP}}(\mathbf{m}) \mathbf{u} \rangle$$

... conditional on  $\nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}_*) = 0$  (a linear constraint on Gaussian noise matrix).

- $[\nabla^2 \mathcal{F}_{\text{TAP}}(\mathbf{m}) | \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}) = 0]$  is a finite rank spiked GOE matrix.

## Technical tool 2: Conditional Gaussian comparison

- ▶ Gaussian comparison lower bound (with many non-trivial idea) [Celentano, Fan, Mei, 2021]

$$\text{rescaled Hessian} \geq \inf_{u,p} \sup_{\alpha,\kappa,\gamma} H_\lambda(p, u; \alpha, \kappa, \gamma) + o(1) > 0.$$

- ▶ So,  $\mathcal{F}_{\text{TAP}}$  is locally strongly convex near  $m_*$ .

Some technical challenges showing this lower bound:

- ▶ Need control on the empirical distribution of coordinates of  $m_*$ . Use Slepian + Kac-Rice to give it a tight control.
- ▶ Need to analyze the variational formula  $\inf \sup H_\lambda$ . First handle the bulk, then show the spikes do not affect the bulk.

## Technical tool 2: Conditional Gaussian comparison

- ▶ Gaussian comparison lower bound (with many non-trivial idea) [Celentano, Fan, Mei, 2021]

$$\text{rescaled Hessian} \geq \inf_{u, p} \sup_{\alpha, \kappa, \gamma} H_\lambda(p, u; \alpha, \kappa, \gamma) + o(1) > 0.$$

- ▶ So,  $\mathcal{F}_{\text{TAP}}$  is locally strongly convex near  $m_*$ .

Some technical challenges showing this lower bound:

- ▶ Need control on the empirical distribution of coordinates of  $m_*$ .  
Use Slepian + Kac-Rice to give it a tight control.
- ▶ Need to analyze the variational formula  $\inf \sup H_\lambda$ .  
First handle the bulk, then show the spikes do not affect the bulk.

## Technical tool 3: Union bound using Kac-Rice

Need to translate

$$\sup_{\mathbf{m}_* \in \text{some region}} \mathbb{P} \left( \min_{\mathbf{m} \in \mathcal{B}(\mathbf{m}_*, \varepsilon \sqrt{n})} \lambda_{\min}(\nabla^2 \mathcal{F}_{\text{TAP}}(\mathbf{m})) < -\varepsilon \right. \\ \left. \mid \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}_*) = 0 \right) < e^{-n\delta}$$

to

$$\mathbb{P} \left( \forall \mathbf{m}_* \in \text{some region}, \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}_*) = 0, \right. \\ \left. \min_{\mathbf{m} \in \mathcal{B}(\mathbf{m}_*, \varepsilon \sqrt{n})} \lambda_{\min}(\nabla^2 \mathcal{F}_{\text{TAP}}(\mathbf{m})) < -\varepsilon \right) < e^{-n\delta}$$

Method: using Kac-Rice to upper bound

$$\mathbb{E}[\#\{\mathbf{m}_* \in \text{some region}, \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}_*) = 0, \\ \min_{\mathbf{m} \in \mathcal{B}(\mathbf{m}_*, \varepsilon \sqrt{n})} \lambda_{\min}(\nabla^2 \mathcal{F}_{\text{TAP}}(\mathbf{m})) < -\varepsilon\}].$$

# Convergence of NGD

- ▶ Local geometry: relatively smooth and relatively strongly convex

$$c \cdot \nabla^2 \text{Ent}(\mathbf{m}) \preceq \nabla^2 \mathcal{F}_{\text{TAP}}(\mathbf{m}) \preceq C \cdot \nabla^2 \text{Ent}(\mathbf{m}).$$

- ▶ Local geometry implies convergence of NGD.

# Convergence of AMP

- ▶ AMP iteration

$$\begin{bmatrix} \mathbf{m}^{k+1} \\ \mathbf{m}^k \end{bmatrix} = \text{AMP} \left( \begin{bmatrix} \mathbf{m}^k \\ \mathbf{m}^{k-1} \end{bmatrix} \right)$$

- ▶ Analyze the linearized AMP operator  $\nabla \text{AMP}(\mathbf{m}_*, \mathbf{m}_*)$ , conditional on  $\nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}_*, \mathbf{m}_*) = \mathbf{0}$ .



# Convergence of AMP

- ▶ Analyze the linearized AMP operator  $\nabla \text{AMP}(\mathbf{m}_*, \mathbf{m}_*)$ , conditional on  $\nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}_*) = \mathbf{0}$ .
- ▶ When  $\lambda > 1$ , w.h.p

$$\sup_{i \in [2n]} |\lambda_i(\nabla \text{AMP}(\mathbf{m}_*, \mathbf{m}_*))| < 1 - \varepsilon,$$

implies asymptotic stability at  $\mathbf{m}_*$ . (Quite non-trivial proof)

- ▶ When  $\lambda \geq \lambda_0$ , w.h.p

$$\sup_{\mathbf{m} \in \mathcal{B}(\mathbf{m}_*, \varepsilon)} \|\nabla \text{AMP}'\|_{\text{op}} < 1 - \varepsilon,$$

implies global convergence.

# Summary

- ▶ TAP has no spurious local min below a threshold ( $\lambda \geq \lambda_0$ )
- ▶ TAP is locally strongly convex near Bayes estimator ( $\lambda > 1$ ).
- ▶ A hybrid spec. + AMP + NGD algorithm converges ( $\lambda > 1$ ).
- ▶ AMP, NGD converges ( $\lambda \geq \lambda_0$ ).
- ▶ AMP is asymptotically stable at  $m_*$  ( $\lambda > 1$ ).
- ▶ The proof strategy can be extended to other problems.