

# Variational inference, spin glasses, and TAP free energy

Song Mei

April 20, 2021

Joint work with Zhou Fan and Andrea Montanari

# General motivation

- ▶ **Bayesian inference**: topic modeling, Bayesian GLM, Bayesian NN... High dim. integration is hard!
- ▶ **Variational inference**: integration  $\rightarrow$  optimization.
- ▶ Popular objective function: “**mean field free energy**”.
- ▶ ... but not optimal for moderate SNR.
- ▶ **Today**: optimal objective “**TAP free energy**”.

# $\mathbb{Z}_2$ 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  $\mathbf{Y} \in \mathbb{R}^{n \times n}$ :

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

- ▶ Noise:  $\mathbf{W} \sim \text{GOE}(n)$ .
- ▶ SNR  $\lambda \in [0, \infty)$  fixed, dimension  $n \rightarrow \infty$ .
- ▶ 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

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

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

- ▶  $\text{MSE}(\widehat{\mathbf{X}}) = \mathbb{E}[(1/n^2) \|\mathbf{X} - \widehat{\mathbf{X}}\|_F^2]$ .
- ▶ For  $\lambda < 1$ , impossible.
- ▶ For  $\lambda > 1$ , possible and efficient. BBP.
- ▶ The optimal estimator is the Bayes estimator:

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

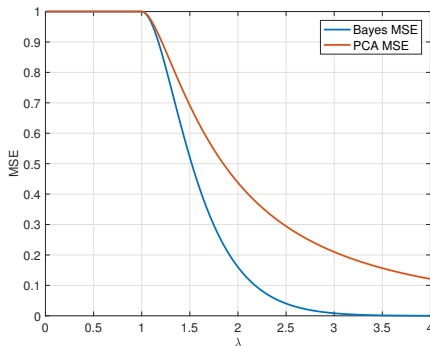
# Bayes estimation in $\mathbb{Z}_2$ synchronization

► 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:

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

- ▶ The posterior distribution:

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

# Mean field variational inference

- ▶ The posterior distribution:

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

- ▶ Approximate  $p(\boldsymbol{\sigma}|\mathbf{Y})$  by  $q \in \mathcal{P}_{\text{MF}}$ :

$$\mathcal{P}_{\text{MF}} = \left\{ q(\boldsymbol{\sigma}) = \prod_{i=1}^n q_i(\sigma_i) : q_i \in \mathcal{P}(\mathbb{Z}_2) \right\} \cong [-1, 1]^n.$$

- ▶ Minimize the relative entropy between  $q$  and  $p(\boldsymbol{\sigma}|\mathbf{Y})$ :

$$\min_{q \in \mathcal{P}_{\text{MF}}} D_{\text{kl}}(q \| p(\boldsymbol{\sigma}|\mathbf{Y})).$$

- ▶ Equivalently minimizing  $\min_{\mathbf{m} \in [-1, 1]^n} \mathcal{F}_{\text{MF}}(\mathbf{m})$

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

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

# Mean field variational inference

- ▶ The posterior distribution:

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

- ▶ Approximate  $p(\boldsymbol{\sigma}|\mathbf{Y})$  by  $q \in \mathcal{P}_{\text{MF}}$ :

$$\mathcal{P}_{\text{MF}} = \left\{ q(\boldsymbol{\sigma}) = \prod_{i=1}^n q_i(\sigma_i) : q_i \in \mathcal{P}(\mathbb{Z}_2) \right\} \cong [-1, 1]^n.$$

- ▶ Minimize the relative entropy between  $q$  and  $p(\boldsymbol{\sigma}|\mathbf{Y})$ :

$$\min_{q \in \mathcal{P}_{\text{MF}}} D_{\text{kl}}(q \| p(\boldsymbol{\sigma}|\mathbf{Y})).$$

- ▶ Equivalently minimizing  $\min_{\mathbf{m} \in [-1, 1]^n} \mathcal{F}_{\text{MF}}(\mathbf{m})$

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

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



# Mean field variational inference

- ▶ The posterior distribution:

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

- ▶ Approximate  $p(\boldsymbol{\sigma}|\mathbf{Y})$  by  $q \in \mathcal{P}_{\text{MF}}$ :

$$\mathcal{P}_{\text{MF}} = \left\{ q(\boldsymbol{\sigma}) = \prod_{i=1}^n q_i(\sigma_i) : q_i \in \mathcal{P}(\mathbb{Z}_2) \right\} \cong [-1, 1]^n.$$

- ▶ Minimize the relative entropy between  $q$  and  $p(\boldsymbol{\sigma}|\mathbf{Y})$ :

$$\min_{q \in \mathcal{P}_{\text{MF}}} D_{\text{kl}}(q \| p(\boldsymbol{\sigma}|\mathbf{Y})).$$

- ▶ Equivalently minimizing  $\min_{\mathbf{m} \in [-1, 1]^n} \mathcal{F}_{\text{MF}}(\mathbf{m})$

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

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

# Mean field variational inference

- ▶ The posterior distribution:

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

- ▶ Approximate  $p(\boldsymbol{\sigma}|\mathbf{Y})$  by  $q \in \mathcal{P}_{\text{MF}}$ :

$$\mathcal{P}_{\text{MF}} = \left\{ q(\boldsymbol{\sigma}) = \prod_{i=1}^n q_i(\sigma_i) : q_i \in \mathcal{P}(\mathbb{Z}_2) \right\} \cong [-1, 1]^n.$$

- ▶ Minimize the relative entropy between  $q$  and  $p(\boldsymbol{\sigma}|\mathbf{Y})$ :

$$\min_{q \in \mathcal{P}_{\text{MF}}} D_{\text{kl}}(q \| p(\boldsymbol{\sigma}|\mathbf{Y})).$$

- ▶ Equivalently minimizing  $\min_{\mathbf{m} \in [-1, 1]^n} \mathcal{F}_{\text{MF}}(\mathbf{m})$

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

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

# Mean field variational inference

- ▶ Mean field free energy:

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

- ▶ For  $\mathbf{m}_\star = \arg \min_{\mathbf{m}} \mathcal{F}_{\text{MF}}(\mathbf{m})$ , we hope

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

- ▶ It was shown that  $\mathbf{m}_\star \mathbf{m}_\star^\top \not\approx \mathbb{E}[\mathbf{x} \mathbf{x}^\top | \mathbf{Y}]$  [Ghorbani, Javadi, and Montanari, 2017].
- ▶ The assumption that posterior distribution can be approximately factorized into the product of marginals is wrong!

# Mean field variational inference

- ▶ Mean field free energy:

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

- ▶ For  $\mathbf{m}_\star = \arg \min_{\mathbf{m}} \mathcal{F}_{\text{MF}}(\mathbf{m})$ , we hope

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

- ▶ It was shown that  $\mathbf{m}_\star \mathbf{m}_\star^\top \not\approx \mathbb{E}[\mathbf{x} \mathbf{x}^\top | \mathbf{Y}]$  [Ghorbani, Javadi, and Montanari, 2017].
- ▶ The assumption that posterior distribution can be approximately factorized into the product of marginals is wrong!

# Mean field variational inference

- ▶ Mean field free energy:

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

- ▶ For  $\mathbf{m}_\star = \arg \min_{\mathbf{m}} \mathcal{F}_{\text{MF}}(\mathbf{m})$ , we hope

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

- ▶ It was shown that  $\mathbf{m}_\star \mathbf{m}_\star^\top \not\approx \mathbb{E}[\mathbf{x} \mathbf{x}^\top | \mathbf{Y}]$  [Ghorbani, Javadi, and Montanari, 2017].
- ▶ The assumption that posterior distribution can be approximately factorized into the product of marginals is wrong!

# Mean field variational inference

- ▶ Mean field free energy:

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

- ▶ For  $\mathbf{m}_\star = \arg \min_{\mathbf{m}} \mathcal{F}_{\text{MF}}(\mathbf{m})$ , we hope

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

- ▶ It was shown that  $\mathbf{m}_\star \mathbf{m}_\star^\top \not\approx \mathbb{E}[\mathbf{x} \mathbf{x}^\top | \mathbf{Y}]$  [Ghorbani, Javadi, and Montanari, 2017].
- ▶ The assumption that posterior distribution can be approximately factorized into the product of marginals is wrong!

# The TAP free energy

- ▶ Thouless, Anderson, and Palmer (1977) proposed the TAP free energy when they study the Sherrington-Kirkpatrick model, whose Gibbs measure gives

$$G_{\beta, \lambda}(\sigma) = \frac{1}{Z_{\beta, \lambda}} \exp\{\beta \langle \sigma, Y \sigma \rangle\}.$$

where  $Y_{ij} \sim \mathcal{N}(\lambda/n, 1/n)$ .

- ▶ When  $\beta = \lambda$ , the Gibbs measure of SK model is the same as the posterior of  $\mathbb{Z}_2$  synchronization
- ▶ The TAP free energy (when  $\beta = \lambda$ ) gives

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

# The TAP free energy

- ▶ Thouless, Anderson, and Palmer (1977) proposed the TAP free energy when they study the Sherrington-Kirkpatrick model, whose Gibbs measure gives

$$G_{\beta, \lambda}(\sigma) = \frac{1}{Z_{\beta, \lambda}} \exp\{\beta \langle \sigma, Y \sigma \rangle\}.$$

where  $Y_{ij} \sim \mathcal{N}(\lambda/n, 1/n)$ .

- ▶ When  $\beta = \lambda$ , the Gibbs measure of SK model is the same as the posterior of  $\mathbb{Z}_2$  synchronization
- ▶ The TAP free energy (when  $\beta = \lambda$ ) gives

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



# The TAP free energy

- ▶ Thouless, Anderson, and Palmer (1977) proposed the TAP free energy when they study the Sherrington-Kirkpatrick model, whose Gibbs measure gives

$$G_{\beta, \lambda}(\sigma) = \frac{1}{Z_{\beta, \lambda}} \exp\{\beta \langle \sigma, Y \sigma \rangle\}.$$

where  $Y_{ij} \sim \mathcal{N}(\lambda/n, 1/n)$ .

- ▶ When  $\beta = \lambda$ , the Gibbs measure of SK model is the same as the posterior of  $\mathbb{Z}_2$  synchronization
- ▶ The TAP free energy (when  $\beta = \lambda$ ) gives

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

# The TAP free energy

- ▶ The TAP free energy

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

- ▶ For  $\mathbf{m}_\star = \arg \min_{\mathbf{m}} \mathcal{F}_{\text{TAP}}(\mathbf{m})$ , we hope

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

- ▶ Our main theorem shows that this is correct.

# The TAP free energy

- ▶ The TAP free energy

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

- ▶ For  $\mathbf{m}_\star = \arg \min_{\mathbf{m}} \mathcal{F}_{\text{TAP}}(\mathbf{m})$ , we hope

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

- ▶ Our main theorem shows that this is correct.

# The TAP free energy

- ▶ The TAP free energy

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

- ▶ For  $\mathbf{m}_\star = \arg \min_{\mathbf{m}} \mathcal{F}_{\text{TAP}}(\mathbf{m})$ , we hope

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

- ▶ Our main theorem shows that this is correct.

# Remarks

- ▶ Derivation of TAP can be obtained by expectation consistency/Plefka's expansion.
- ▶ The stationery equation of TAP  $\nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}) = 0$  is the fixed point equation for AMP algorithm

$$\mathbf{m}^{k+1} = \text{AMP}(\mathbf{m}^k, \mathbf{m}^{k-1}), \quad \mathbf{m}^{k+1} = \mathbf{m}^k = \mathbf{m}^{k-1}.$$

- ▶ AMP is not a descent algorithm on TAP free energy.

# Main theorem

## Theorem (Fan, M., 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

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

All the critical points (below a threshold) are close to the Bayesian estimator.

# Main theorem

## Theorem (Fan, M., 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

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

All the critical points (below a threshold) are close to the Bayesian estimator.

# Related literatures in spin glass theory

TAP free energy in unbiased SK.

- ▶ TAP equations: [Talagrand, 2004], [Chatterjee, 2009], [Chen, 2011], [Auffinger and Jagannath, 2016], Posterior means/Pure states satisfy TAP equations.
- ▶ TAP free energy: [Chen and Panchenko, 2017], constrained TAP minimum are exact.

Calculating the complexity.

- ▶ [Auffinger, Ben Arous, and Cerny, 2010], [Subag, 2016].

And a few more after 2018



# Main theorem

## Theorem (Fan, M., 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

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

All the critical points (below a threshold) are close to the Bayesian estimator.

# Proof idea - Count the number of critical points

- Recall

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

- Define some important statistics of  $\mathbf{m}$ :

$$E(\mathbf{m}) = \mathcal{F}_{\text{TAP}}(\mathbf{m})/n, \quad Q(\mathbf{m}) = \|\mathbf{m}\|_2^2/n, \quad M(\mathbf{m}) = \langle \mathbf{m}, \mathbf{x} \rangle/n.$$

- For any  $U \subseteq \mathbb{R}^3$ , define

$$\text{Crit}_n(U) \equiv \#\{\mathbf{m} : \nabla E(\mathbf{m}) = \mathbf{0}, (Q(\mathbf{m}), M(\mathbf{m}), E(\mathbf{m})) \in U\}. \quad (2)$$

# Calculating the Crit: Kac-Rice formula

Lemma (Kac-Rice formula, c.f. [Adler and Taylor, 2007])

*] Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a “sufficiently regular” random morse function. Let  $p_m(z)$  be the density of  $\nabla f(m)$  at  $z$ . For any Borel measurable set  $T \subseteq \mathbb{R}^d$ , denote*

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

*Then*

$$\begin{aligned}\mathbb{E}[\text{Crit}(T)] &= \mathbb{E}\left[\int_T |\det \nabla^2 f(m)| \cdot \delta(\nabla f(m)) \cdot dm\right] \\ &= \int_T \mathbb{E}\left[|\det \nabla^2 f(m)| \mid \nabla f(m) = 0\right] p_m(0) dm.\end{aligned}$$

- $|\det \nabla^2 f(m)|$  is the correct weight function so that each critical point count exactly once.

# Proof idea - Count the number of critical points

## Proposition

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

$$S_\star(q, \varphi, e) = \sup_{a \in \mathbb{R}} \inf_{(\mu, \nu, \tau, \gamma) \in \mathbb{R}^4} S(q, \varphi, a, e; \mu, \nu, \tau, \gamma),$$

where

$$\begin{aligned} S(q, \varphi, a, e; \mu, \nu, \tau, \gamma) = & \frac{1}{4\beta^2} \left[ \frac{a}{q} - \frac{\beta\lambda\varphi^2}{q} - \beta^2(1-q) \right]^2 \\ & - q\mu - \varphi\nu - a\tau - \left[ -\frac{\beta^2}{4}(1-q^2) + \frac{a}{2} - e \right] \gamma + \log I, \end{aligned}$$

and

$$\begin{aligned} I = & \int_{-\infty}^{\infty} \frac{1}{(2\pi\beta^2q)^{1/2}} \exp \left\{ -\frac{(x - \beta\lambda\varphi)^2}{2\beta^2q} \right. \\ & \left. + \mu \tanh^2(x) + \nu \tanh(x) + \tau x \tanh(x) + \gamma \log[2 \cosh(x)] \right\} dx. \end{aligned}$$

# Proof idea - Count the number of critical points

- ▶ Key proposition: for  $U \subseteq \mathbb{R}^3$ ,

$$\mathbb{E}[\text{Crit}_n(U)] \leq \exp \left\{ n \overbrace{\sup_{(q, \varphi, e) \in U} S_\star(q, \varphi, e)}^{T(U)} + o(n) \right\},$$

- ▶ For any  $U$  such that  $T(U) > 0$ , there could potentially be critical points of  $\mathcal{F}_{\text{TAP}}$  in  $U$ .
- ▶ For any  $U$  such that  $T(U) < 0$ , there is no critical points of  $\mathcal{F}_{\text{TAP}}$  in  $U$ , with high probability.
- ▶ If we admit the key proposition, suffice to show that  $T(U) < 0$  unless  $U$  contains a neighborhood of the Bayes estimator.

# Proof idea - Count the number of critical points

- ▶ Key proposition: for  $U \subseteq \mathbb{R}^3$ ,

$$\mathbb{E}[\text{Crit}_n(U)] \leq \exp \left\{ n \overbrace{\sup_{(q, \varphi, e) \in U} S_\star(q, \varphi, e)}^{T(U)} + o(n) \right\},$$

- ▶ For any  $U$  such that  $T(U) > 0$ , there could potentially be critical points of  $\mathcal{F}_{\text{TAP}}$  in  $U$ .
- ▶ For any  $U$  such that  $T(U) < 0$ , there is no critical points of  $\mathcal{F}_{\text{TAP}}$  in  $U$ , with high probability.
- ▶ If we admit the key proposition, suffice to show that  $T(U) < 0$  unless  $U$  contains a neighborhood of the Bayes estimator.

# Proof idea - Count the number of critical points

- ▶ Key proposition: for  $U \subseteq \mathbb{R}^3$ ,

$$\mathbb{E}[\text{Crit}_n(U)] \leq \exp \left\{ n \overbrace{\sup_{(q, \varphi, e) \in U} S_\star(q, \varphi, e)}^{T(U)} + o(n) \right\},$$

- ▶ For any  $U$  such that  $T(U) > 0$ , there could potentially be critical points of  $\mathcal{F}_{\text{TAP}}$  in  $U$ .
- ▶ For any  $U$  such that  $T(U) < 0$ , there is no critical points of  $\mathcal{F}_{\text{TAP}}$  in  $U$ , with high probability.
- ▶ If we admit the key proposition, suffice to show that  $T(U) < 0$  unless  $U$  contains a neighborhood of the Bayes estimator.

# Proof idea - Count the number of critical points

- ▶ Key proposition: for  $U \subseteq \mathbb{R}^3$ ,

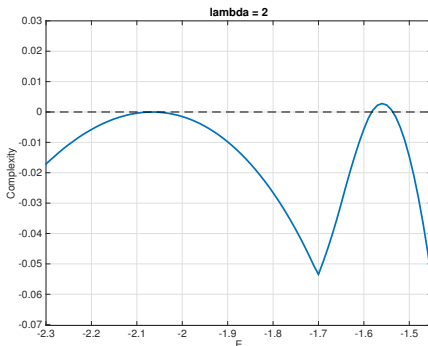
$$\mathbb{E}[\text{Crit}_n(U)] \leq \exp \left\{ n \overbrace{\sup_{(q, \varphi, e) \in U} S_\star(q, \varphi, e)}^{T(U)} + o(n) \right\},$$

- ▶ For any  $U$  such that  $T(U) > 0$ , there could potentially be critical points of  $\mathcal{F}_{\text{TAP}}$  in  $U$ .
- ▶ For any  $U$  such that  $T(U) < 0$ , there is no critical points of  $\mathcal{F}_{\text{TAP}}$  in  $U$ , with high probability.
- ▶ If we admit the key proposition, suffice to show that  $T(U) < 0$  unless  $U$  contains a neighborhood of the Bayes estimator.



## Proof idea - the complexity function $S_\star$

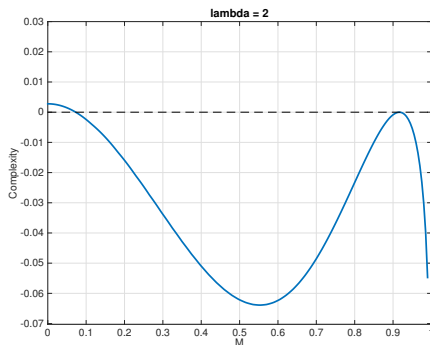
►  $S_\star(e) = \sup_{q, \varphi} S_\star(q, \varphi, e).$



► At  $e_\star$ ,  $S_\star(e_\star) = 0.$

# Proof idea - the complexity function $S_\star$

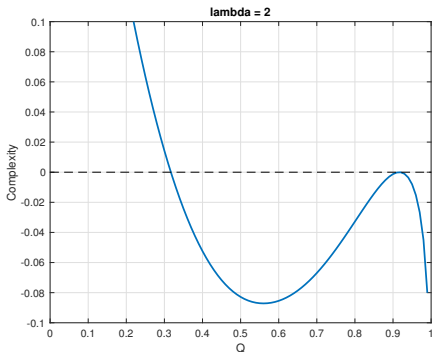
►  $S_\star(\varphi) = \sup_{q,e} S_\star(q, \varphi, e).$



► At  $\varphi_\star$ ,  $S_\star(\varphi_\star) = 0.$

# Proof idea - the complexity function $S_\star$

►  $S_\star(q) = \sup_{\varphi, e} S_\star(q, \varphi, e).$



► At  $q_\star$ ,  $S_\star(q_\star) = 0$ .

## Proof idea - the complexity function $S_\star$

There exists  $\lambda_0$ , for  $\lambda \geq \lambda_0$ ,

- ▶  $S_\star(q_\star, \varphi_\star, e_\star) = 0$ , where  $(q_\star, \varphi_\star, e_\star) \approx (Q(\mathbf{m}_\star), M(\mathbf{m}_\star), E(\mathbf{m}_\star))$  for  $\widehat{\mathbf{X}}_{\text{Bayes}} \approx \mathbf{m}_\star \mathbf{m}_\star^\top$ .
- ▶  $S_\star(q, \varphi, e) < 0$  for any  $e \leq -\lambda^2/3$  and  $(q, \varphi, e) \neq (q_\star, \varphi_\star, e_\star)$ .

The proof of these two properties is more than calculus. It requires bounds using concentration inequalities.

Combining with the key inequality it is easy to show the main theorem.

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

Now suffice to show the key inequality.

## Proof idea - the complexity function $S_*$

There exists  $\lambda_0$ , for  $\lambda \geq \lambda_0$ ,

- ▶  $S_*(q_*, \varphi_*, e_*) = 0$ , where  $(q_*, \varphi_*, e_*) \approx (Q(\mathbf{m}_*), M(\mathbf{m}_*), E(\mathbf{m}_*))$  for  $\widehat{\mathbf{X}}_{\text{Bayes}} \approx \mathbf{m}_* \mathbf{m}_*^\top$ .
- ▶  $S_*(q, \varphi, e) < 0$  for any  $e \leq -\lambda^2/3$  and  $(q, \varphi, e) \neq (q_*, \varphi_*, e_*)$ .

The proof of these two properties is more than calculus. It requires bounds using concentration inequalities.

Combining with the key inequality it is easy to show the main theorem.

$$\mathbb{E}[\text{Crit}_n(U)] \leq \exp \left\{ n \sup_{(q, \varphi, e) \in U} S_*(q, \varphi, e) + o(n) \right\}.$$

Now suffice to show the key inequality.

## Proof idea - the complexity function $S_\star$

There exists  $\lambda_0$ , for  $\lambda \geq \lambda_0$ ,

- ▶  $S_\star(q_\star, \varphi_\star, e_\star) = 0$ , where  $(q_\star, \varphi_\star, e_\star) \approx (Q(\mathbf{m}_\star), M(\mathbf{m}_\star), E(\mathbf{m}_\star))$  for  $\widehat{\mathbf{X}}_{\text{Bayes}} \approx \mathbf{m}_\star \mathbf{m}_\star^\top$ .
- ▶  $S_\star(q, \varphi, e) < 0$  for any  $e \leq -\lambda^2/3$  and  $(q, \varphi, e) \neq (q_\star, \varphi_\star, e_\star)$ .

The proof of these two properties is more than calculus. It requires bounds using concentration inequalities.

Combining with the key inequality it is easy to show the main theorem.

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

Now suffice to show the key inequality.

# Dealing with determinant of Hessian

- ▶ The conditional Hessian is distributed as (up to some scaling)

$$[\nabla^2 \mathcal{F}_{\text{TAP}}(\mathbf{m}) | \nabla \mathcal{F}_{\text{TAP}}(\mathbf{m}) = \mathbf{0}] \stackrel{d}{=} \mathbf{D} + \mathbf{W} + \text{low rank perturbation},$$

where  $\mathbf{D} = \text{diag}(d_i)$ , and  $\mathbf{W} \sim \text{GOE}(n)$ .

- ▶ The low rank perturbation has vanishing effects. Therefore, we just need to calculate  $\mathbb{E}[|\det(\mathbf{H})|]$ , with

$$\mathbf{H} = \mathbf{D} + \mathbf{W}.$$

# Determinant of Hessian: "replica trick"

$$\mathbf{H} = \mathbf{D} + \mathbf{W} = \text{diagonal} + \text{GOE}.$$

We have

$$\begin{aligned}\mathbb{E}[|\det(\mathbf{H})|] &\approx \mathbb{E}[|\det(\mathbf{H})|^{-1/2}]^2 \\ &= \mathbb{E}\left[\int_{\mathbb{R}^d} \exp\{-\langle \mathbf{x}, \mathbf{H}\mathbf{x} \rangle/2\} d\mathbf{x}\right]^2 \\ &= \mathbb{E}\left[\int_{\mathbb{R}^d} \exp\{-\langle \mathbf{x}, \mathbf{W}\mathbf{x} \rangle/2 - \langle \mathbf{x}, \mathbf{D}\mathbf{x} \rangle\} d\mathbf{x}\right]^2\end{aligned}$$



# Determinant of Hessian: Stieltjes transform

$$\mathbf{H} = \mathbf{D} + \mathbf{W} = \text{diagonal} + \text{GOE}.$$

$$\frac{1}{n} \log |\det(\mathbf{H})| = \frac{1}{n} \log \prod_{i=1}^n |\lambda_i(\mathbf{H})| = \int_{\mathbb{R}} \log |x| \cdot \mu_{\mathbf{H}}(\mathrm{d}x),$$

where  $\mu_{\mathbf{H}} = (1/n) \sum_{i=1}^n \delta(\lambda_i(\mathbf{H}))$ .

- ▶ The Stieltjes transform of  $\mu_{\mathbf{H}}$  can be approximately calculated using free probability theory.
- ▶ Once the Stieltjes transform of  $\mu_{\mathbf{H}}$  is known, the quantity  $\mathbb{E} \left[ \int_{\mathbb{R}} (\log |x|) \mu_{\mathbf{H}}(\mathrm{d}x) \right]$  can be computed.

# Determinant of Hessian: Stieltjes transform

$$\mathbf{H} = \mathbf{D} + \mathbf{W} = \text{diagonal} + \text{GOE}.$$

$$\frac{1}{n} \log |\det(\mathbf{H})| = \frac{1}{n} \log \prod_{i=1}^n |\lambda_i(\mathbf{H})| = \int_{\mathbb{R}} \log |x| \cdot \mu_{\mathbf{H}}(\mathrm{d}x),$$

where  $\mu_{\mathbf{H}} = (1/n) \sum_{i=1}^n \delta(\lambda_i(\mathbf{H}))$ .

- ▶ The Stieltjes transform of  $\mu_{\mathbf{H}}$  can be approximately calculated using free probability theory.
- ▶ Once the Stieltjes transform of  $\mu_{\mathbf{H}}$  is known, the quantity  $\mathbb{E} \left[ \int_{\mathbb{R}} (\log |x|) \mu_{\mathbf{H}}(\mathrm{d}x) \right]$  can be computed.

# Free convolution of two distribution

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , and  $\mu_{\mathbf{A}} = (1/n) \sum_{i=1}^n \delta(\lambda_i(\mathbf{A}))$ . For any  $z \in \mathbb{C}_+$ , the Stieltjes transform of  $\mu_{\mathbf{A}}$  is defined as

$$g_{\mathbf{A}}(z) = \int_{\mathbb{R}} \frac{1}{x - z} \mu_{\mathbf{A}}(dx) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\lambda_i(\mathbf{A}) - z}.$$

Lemma (Free probability theory heuristics/Leave one out)

Let  $\mathbf{D} = \text{diag}(d_i)$  be a diagonal matrix, and let  $\mathbf{H} = \mathbf{D} + \mathbf{W}$ . Then

$$\mathbb{E} g_{\mathbf{H}}(z) = \frac{1}{n} \sum_{i=1}^n \frac{1}{d_i - z - \mathbb{E} g_{\mathbf{H}}(z)} + o_n(1). \quad (3)$$

# Free convolution of two distribution

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , and  $\mu_{\mathbf{A}} = (1/n) \sum_{i=1}^n \delta(\lambda_i(\mathbf{A}))$ . For any  $z \in \mathbb{C}_+$ , the Stieltjes transform of  $\mu_{\mathbf{A}}$  is defined as

$$g_{\mathbf{A}}(z) = \int_{\mathbb{R}} \frac{1}{x - z} \mu_{\mathbf{A}}(dx) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\lambda_i(\mathbf{A}) - z}.$$

Lemma (Free probability theory heuristics/Leave one out)

Let  $\mathbf{D} = \text{diag}(d_i)$  be a diagonal matrix, and let  $\mathbf{H} = \mathbf{D} + \mathbf{W}$ . Then

$$\mathbb{E} g_{\mathbf{H}}(z) = \frac{1}{n} \sum_{i=1}^n \frac{1}{d_i - z - \mathbb{E} g_{\mathbf{H}}(z)} + o_n(1). \quad (3)$$

# Summary

- ▶ TAP free energy is accurate for  $\mathbb{Z}_2$  synchronization.
- ▶ Can be generalized to topic modeling, low rank matrix estimation, compressed sensing, etc...
- ▶ More sophisticated variational inference beyond naive mean field.

# Summary

- ▶ **TAP free energy** is accurate for  $\mathbb{Z}_2$  synchronization.
- ▶ Can be generalized to topic modeling, low rank matrix estimation, compressed sensing, etc...
- ▶ More sophisticated variational inference beyond **naive mean field**.

# Summary

- ▶ **TAP free energy** is accurate for  $\mathbb{Z}_2$  synchronization.
- ▶ Can be generalized to topic modeling, low rank matrix estimation, compressed sensing, etc...
- ▶ More sophisticated variational inference beyond **naive mean field**.