# Variational inference, spin glasses, and TAP free energy

Song Mei

Stanford University

September 19, 2018

Joint work with Zhou Fan and Andrea Montanari

# General motivation

▶ Bayesian inference: high dimensional integration is hard!

▶ Variational inference: integration/summation → optimization.
  A popular objective function: "mean field free energy".

▶ Applications: topic modeling, stochastic block model, low rank
  matrix estimation, compressed sensing....
  ... within which "MF free energy" is known to be not optimal.

▶ Today: introduce the optimal objective "TAP free energy", and
  provide rigorous results.

# General motivation

▶ Bayesian inference: high dimensional integration is hard!

▶ Variational inference: integration/summation $\rightarrow$ optimization.
  A popular objective function: "mean field free energy".

▶ Applications: topic modeling, stochastic block model, low rank
  matrix estimation, compressed sensing....
  ... within which "MF free energy" is known to be not optimal.

▶ Today: introduce the optimal objective "TAP free energy", and
  provide rigorous results.

# General motivation

- Bayesian inference: high dimensional integration is hard!

- Variational inference: integration/summation $\to$ optimization.
  A popular objective function: "mean field free energy".

- Applications: topic modeling, stochastic block model, low rank matrix estimation, compressed sensing....

  ... within which "MF free energy" is known to be not optimal.

- Today: introduce the optimal objective "TAP free energy", and provide rigorous results.

# General motivation

▶ Bayesian inference: high dimensional integration is hard!

▶ Variational inference: integration/summation $\rightarrow$ optimization.
  A popular objective function: "mean field free energy".

▶ Applications: topic modeling, stochastic block model, low rank matrix estimation, compressed sensing....

  ... within which "MF free energy" is known to be not optimal.

▶ Today: introduce the optimal objective "TAP free energy", and provide rigorous results.

# General motivation

▶ Bayesian inference: high dimensional integration is hard!

▶ Variational inference: integration/summation $\rightarrow$ optimization.
  A popular objective function: "mean field free energy".

▶ Applications: topic modeling, stochastic block model, low rank
  matrix estimation, compressed sensing....
  ... within which "MF free energy" is known to be not optimal.

▶ Today: introduce the optimal objective "TAP free energy", and
  provide rigorous results.

# General motivation

- Bayesian inference: high dimensional integration is hard!

- Variational inference: integration/summation $\rightarrow$ optimization.
  A popular objective function: "mean field free energy".

- Applications: topic modeling, stochastic block model, low rank matrix estimation, compressed sensing....
  ... within which "MF free energy" is known to be not optimal.

- Today: introduce the optimal objective "TAP free energy", and provide rigorous results.

# $\mathbb{Z}_2$ synchronization

- Signal:

$$x = [x_1, \ldots, x_n]^\mathsf{T} \in \mathbb{Z}_2^n, \qquad x_i \overset{i.i.d.}{\sim} \text{Unif}(\mathbb{Z}_2), \qquad \mathbb{Z}_2 = \{+1, -1\}.$$

- 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 $n \to \infty$.

- In matrix notation:

$$Y = \frac{\lambda}{n} x x^\mathsf{T} + W.$$

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

# $\mathbb{Z}_2$ synchronization

- Signal:

$$\boldsymbol{x} = [\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n]^\mathsf{T} \in \mathbb{Z}_2^n, \quad \boldsymbol{x}_i \overset{i.i.d.}{\sim} \texttt{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

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

$$Y_{ij} = \frac{\lambda}{n}\boldsymbol{x}_i\boldsymbol{x}_j + W_{ij}.$$

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

$$\boldsymbol{Y} = \frac{\lambda}{n}\boldsymbol{x}\boldsymbol{x}^\mathsf{T} + \boldsymbol{W}.$$

- Task: given $\boldsymbol{Y} = (Y_{ij})$, estimate $\boldsymbol{x}$ (or say $\boldsymbol{X} = \boldsymbol{x}\boldsymbol{x}^\mathsf{T}$).

# $\mathbb{Z}_2$ synchronization

- Signal:

$$\boldsymbol{x} = [\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n]^\mathsf{T} \in \mathbb{Z}_2^n, \quad \boldsymbol{x}_i \overset{i.i.d.}{\sim} \text{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

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

$$Y_{ij} = \frac{\lambda}{n} \boldsymbol{x}_i \boldsymbol{x}_j + W_{ij}.$$

- Noise $W_{ij} \sim \mathcal{N}(0, 1/n)$.

- SNR $\lambda \in [0, \infty)$ fixed, dimension $n \to \infty$.

- In matrix notation:

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

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

# $\mathbb{Z}_2$ synchronization

- Signal:

$$\boldsymbol{x} = [\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n]^\mathsf{T} \in \mathbb{Z}_2^n, \quad \boldsymbol{x}_i \overset{i.i.d.}{\sim} \mathtt{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

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

$$Y_{ij} = \frac{\lambda}{n}\boldsymbol{x}_i\boldsymbol{x}_j + W_{ij}.$$

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

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

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

# $\mathbb{Z}_2$ synchronization

▶ Signal:

$$x = [x_1, \ldots, x_n]^\mathsf{T} \in \mathbb{Z}_2^n, \quad x_i \overset{i.i.d.}{\sim} \mathrm{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

▶ 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 $n \to \infty$.
▶ In matrix notation:

$$Y = \frac{\lambda}{n} xx^\mathsf{T} + W.$$

▶ Task: given $Y = (Y_{ij})$, estimate $x$ (or say $X = xx^\mathsf{T}$).

# $\mathbb{Z}_2$ synchronization

- Signal:

$$\boldsymbol{x} = [\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n]^\mathsf{T} \in \mathbb{Z}_2^n, \quad \boldsymbol{x}_i \overset{i.i.d.}{\sim} \texttt{Unif}(\mathbb{Z}_2), \quad \mathbb{Z}_2 = \{+1, -1\}.$$

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

$$Y_{ij} = \frac{\lambda}{n} \boldsymbol{x}_i \boldsymbol{x}_j + W_{ij}.$$

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

$$\boldsymbol{Y} = \frac{\lambda}{n} \boldsymbol{x} \boldsymbol{x}^\mathsf{T} + \boldsymbol{W}.$$

- Task: given $\boldsymbol{Y} = (Y_{ij})$, estimate $\boldsymbol{x}$ (or say $\boldsymbol{X} = \boldsymbol{x}\boldsymbol{x}^\mathsf{T}$).

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

- Settings:
$$x \sim \text{Unif}(\mathbb{Z}_2^n), \qquad Y = (\lambda/n)xx^\top + W.$$

- Estimate $X = xx^\top$ with loss:

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

- For $\lambda < 1$, estimation is impossible.

- For $\lambda > 1$, estimation is possible and efficient, e.g., spectral estimator (Baik, Ben Arous, Peche phase transition).

- The optimal estimator is the Bayes estimator (also minimax estimator):

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

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

- Settings:
$$\boldsymbol{x} \sim \mathrm{Unif}(\mathbb{Z}_2^n), \qquad \boldsymbol{Y} = (\lambda/n)\boldsymbol{x}\boldsymbol{x}^\top + \boldsymbol{W}.$$

- Estimate $\boldsymbol{X} = \boldsymbol{x}\boldsymbol{x}^\top$ with loss:
$$\ell(\boldsymbol{X}, \widehat{\boldsymbol{X}}) = (1/n^2)\|\boldsymbol{X} - \widehat{\boldsymbol{X}}\|_F^2.$$

- For $\lambda < 1$, estimation is impossible.
- For $\lambda > 1$, estimation is possible and efficient, e.g., spectral estimator (Baik, Ben Arous, Peche phase transition).
- The optimal estimator is the Bayes estimator (also minimax estimator):
$$\widehat{\boldsymbol{X}}_{\mathrm{Bayes}} = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^\top | \boldsymbol{Y}].$$

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

▶ Settings:
$$x \sim \text{Unif}(\mathbb{Z}_2^n), \qquad Y = (\lambda/n)xx^\top + W.$$

▶ Estimate $X = xx^\top$ with loss:
$$\ell(X, \widehat{X}) = (1/n^2)\|X - \widehat{X}\|_F^2.$$

▶ For $\lambda < 1$, estimation is impossible.

▶ For $\lambda > 1$, estimation is possible and efficient, e.g., spectral estimator (Baik, Ben Arous, Peche phase transition).

▶ The optimal estimator is the Bayes estimator (also minimax estimator):
$$\widehat{X}_{\text{Bayes}} = \mathbb{E}[xx^\top | Y].$$

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

- Settings:
$$\boldsymbol{x} \sim \text{Unif}(\mathbb{Z}_2^n), \qquad \boldsymbol{Y} = (\lambda/n)\boldsymbol{x}\boldsymbol{x}^\top + \boldsymbol{W}.$$

- Estimate $\boldsymbol{X} = \boldsymbol{x}\boldsymbol{x}^\top$ with loss:
$$\ell(\boldsymbol{X}, \widehat{\boldsymbol{X}}) = (1/n^2)\|\boldsymbol{X} - \widehat{\boldsymbol{X}}\|_F^2.$$

- For $\lambda < 1$, estimation is impossible.

- For $\lambda > 1$, estimation is possible and efficient, e.g., spectral estimator (Baik, Ben Arous, Peche phase transition).

- The optimal estimator is the Bayes estimator (also minimax estimator):
$$\widehat{\boldsymbol{X}}_{\text{Bayes}} = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^\top | \boldsymbol{Y}].$$

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

- Settings:
$$x \sim \mathrm{Unif}(\mathbb{Z}_2^n), \qquad Y = (\lambda/n)xx^\top + W.$$

- Estimate $X = xx^\top$ with loss:
$$\ell(X, \widehat{X}) = (1/n^2)\|X - \widehat{X}\|_F^2.$$

- For $\lambda < 1$, estimation is impossible.

- For $\lambda > 1$, estimation is possible and efficient, e.g., spectral estimator (Baik, Ben Arous, Peche phase transition).

- The optimal estimator is the Bayes estimator (also minimax estimator):
$$\widehat{X}_{\mathrm{Bayes}} = \mathbb{E}[xx^\top | Y].$$

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

- Settings:
$$\boldsymbol{x} \sim \mathrm{Unif}(\mathbb{Z}_2^n), \qquad \boldsymbol{Y} = (\lambda/n)\boldsymbol{x}\boldsymbol{x}^\top + \boldsymbol{W}.$$

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

# Compute the Bayesian estimator

▶ The Bayesian estimator:

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

▶ The posterior distribution:

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

# Compute the Bayesian estimator

▶ The Bayesian estimator:

$$\widehat{X}_{\mathrm{Bayes}} = \mathbb{E}[xx^\mathsf{T}|Y] = \sum_{\sigma \in \mathbb{Z}_2^n} \sigma\sigma^\mathsf{T} p(\sigma|Y).$$

▶ The posterior distribution:

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

# Mean field variational inference

▶ The posterior distribution:

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

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

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

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

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

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

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

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

# Mean field variational inference

▶ The posterior distribution:

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

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

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

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

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

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

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

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

# Mean field variational inference

▶ The posterior distribution:

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

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

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

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

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

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

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

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

# Mean field variational inference

▶ The posterior distribution:

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

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

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

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

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

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

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

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

# Mean field variational inference

▶ Mean field free energy:

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

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

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

▶ It was shown that $\boldsymbol{m}_\star \boldsymbol{m}_\star^\top \napprox \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^\top | \boldsymbol{Y}]$ [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:

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

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

$$\boldsymbol{m}_\star\boldsymbol{m}_\star^\mathsf{T} \approx \widehat{\boldsymbol{X}}_{\mathrm{Bayes}} = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^\mathsf{T}|\boldsymbol{Y}].$$

▶ It was shown that $\boldsymbol{m}_\star\boldsymbol{m}_\star^\mathsf{T} \not\approx \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^\mathsf{T}|\boldsymbol{Y}]$ [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:

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

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

$$\boldsymbol{m}_\star \boldsymbol{m}_\star^\mathsf{T} \approx \widehat{\boldsymbol{X}}_{\mathrm{Bayes}} = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^\mathsf{T}|\boldsymbol{Y}].$$

- It was shown that $\boldsymbol{m}_\star \boldsymbol{m}_\star^\mathsf{T} \not\approx \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^\mathsf{T}|\boldsymbol{Y}]$ [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:

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

► For $\boldsymbol{m}_\star = \arg\min_{\boldsymbol{m}} \mathcal{F}_{\mathrm{MF}}(\boldsymbol{m})$, we hope

$$\boldsymbol{m}_\star \boldsymbol{m}_\star^{\mathsf{T}} \approx \widehat{\boldsymbol{X}}_{\mathrm{Bayes}} = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^{\mathsf{T}}|\boldsymbol{Y}].$$

► It was shown that $\boldsymbol{m}_\star \boldsymbol{m}_\star^{\mathsf{T}} \not\approx \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^{\mathsf{T}}|\boldsymbol{Y}]$ [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

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

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

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

# 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

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

  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

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

# 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

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

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

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

# The TAP free energy

- The TAP free energy

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

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

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

- Our main theorem shows that this is correct.

# The TAP free energy

- The TAP free energy

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

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

$$\boldsymbol{m}_\star \boldsymbol{m}_\star^{\mathsf{T}} \approx \widehat{\boldsymbol{X}}_{\mathrm{Bayes}} = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^{\mathsf{T}}|\boldsymbol{Y}].$$

- Our main theorem shows that this is correct.

# The TAP free energy

- The TAP free energy

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

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

$$\boldsymbol{m}_\star \boldsymbol{m}_\star^{\mathsf{T}} \approx \widehat{\boldsymbol{X}}_{\mathrm{Bayes}} = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^{\mathsf{T}}|\boldsymbol{Y}].$$

- Our main theorem shows that this is correct.

# Proof of the main theorem

**Theorem (Fan, M., Montanari, 2018)**

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

$$\lim_{n \to \infty} \mathbb{E}\Big[ \sup_{m \in \mathcal{C}_{\lambda,n}} \frac{1}{n^2} \|mm^\intercal - \widehat{X}_{\text{Bayes}}\|_F^2 \wedge 1 \Big] = 0. \tag{1}$$

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

# Proof of the main theorem

**Theorem (Fan, M., Montanari, 2018)**

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

$$\lim_{n \to \infty} \mathbb{E}\Big[\sup_{m \in \mathcal{C}_{\lambda,n}} \frac{1}{n^2} \|mm^\top - \widehat{X}_{\mathrm{Bayes}}\|_F^2 \wedge 1\Big] = 0. \tag{1}$$

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

# Relationship with AMP

▶ Another way to construct the Bayes estimator is approximate message passing [Donoho, Maleki, and Montanari, 2009], [Bolthausen, 2014]:

$$\boldsymbol{m}^{k+1} = \tanh(\lambda \boldsymbol{Y} \boldsymbol{m}^k - \lambda^2 [1 - \|\boldsymbol{m}^k\|_2^2/n] \boldsymbol{m}^{k-1}).$$

▶ Fixed point of AMP is a critical point of the TAP free energy.

▶ The risk of AMP iterations converge to the Bayes risk [Deshpande, Abbes, and Montanari, 2016], [Montanari and Venkataramanan, 2017]:

$$\lim_{k \to \infty} \lim_{n \to \infty} \frac{1}{n^2} \|\boldsymbol{m}^k (\boldsymbol{m}^k)^\top - \boldsymbol{x} \boldsymbol{x}^\top\|_F^2 = \lim_{n \to \infty} \mathrm{MSE}_n(\widehat{\boldsymbol{X}}_{\mathrm{Bayes}}).$$

▶ But it is not known if AMP will converge to a fixed point (It is still an open problem).

# Relationship with AMP

▶ Another way to construct the Bayes estimator is approximate message passing [Donoho, Maleki, and Montanari, 2009], [Bolthausen, 2014]:

$$\boldsymbol{m}^{k+1} = \tanh(\lambda \boldsymbol{Y} \boldsymbol{m}^k - \lambda^2 [1 - \|\boldsymbol{m}^k\|_2^2 / n] \boldsymbol{m}^{k-1}).$$

▶ Fixed point of AMP is a critical point of the TAP free energy.

▶ The risk of AMP iterations converge to the Bayes risk [Deshpande, Abbes, and Montanari, 2016], [Montanari and Venkataramanan, 2017]:

$$\lim_{k \to \infty} \lim_{n \to \infty} \frac{1}{n^2} \|\boldsymbol{m}^k (\boldsymbol{m}^k)^\top - \boldsymbol{x} \boldsymbol{x}^\top\|_F^2 = \lim_{n \to \infty} \mathrm{MSE}_n(\widehat{\boldsymbol{X}}_{\mathrm{Bayes}}).$$

▶ But it is not known if AMP will converge to a fixed point (It is still an open problem).

# Relationship with AMP

▶ Another way to construct the Bayes estimator is approximate message passing [Donoho, Maleki, and Montanari, 2009], [Bolthausen, 2014]:

$$\boldsymbol{m}^{k+1} = \tanh(\lambda \boldsymbol{Y} \boldsymbol{m}^k - \lambda^2 [1 - \|\boldsymbol{m}^k\|_2^2/n] \boldsymbol{m}^{k-1}) \,.$$

▶ Fixed point of AMP is a critical point of the TAP free energy.

▶ The risk of AMP iterations converge to the Bayes risk [Deshpande, Abbes, and Montanari, 2016], [Montanari and Venkataramanan, 2017]:

$$\lim_{k \to \infty} \lim_{n \to \infty} \frac{1}{n^2} \|\boldsymbol{m}^k (\boldsymbol{m}^k)^\mathsf{T} - \boldsymbol{x} \boldsymbol{x}^\mathsf{T}\|_F^2 = \lim_{n \to \infty} \mathrm{MSE}_n(\widehat{\boldsymbol{X}}_{\mathrm{Bayes}}).$$

▶ But it is not known if AMP will converge to a fixed point (It is still an open problem).

# Relationship with AMP

▶ Another way to construct the Bayes estimator is approximate message passing [Donoho, Maleki, and Montanari, 2009], [Bolthausen, 2014]:

$$m^{k+1} = \tanh(\lambda Y m^k - \lambda^2 [1 - \|m^k\|_2^2/n] m^{k-1}) \,.$$

▶ Fixed point of AMP is a critical point of the TAP free energy.

▶ The risk of AMP iterations converge to the Bayes risk [Deshpande, Abbes, and Montanari, 2016], [Montanari and Venkataramanan, 2017]:

$$\lim_{k \to \infty} \lim_{n \to \infty} \frac{1}{n^2} \|m^k (m^k)^\mathsf{T} - xx^\mathsf{T}\|_F^2 = \lim_{n \to \infty} \mathrm{MSE}_n(\widehat{X}_{\mathrm{Bayes}}).$$

▶ But it is not known if AMP will converge to a fixed point (It is still an open problem).

# 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].

# Proof of the main theorem

**Theorem (Fan, M., Montanari, 2018)**

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

$$\lim_{n \to \infty} \mathbb{E}\Big[ \sup_{\boldsymbol{m} \in \mathcal{C}_{\lambda,n}} \frac{1}{n^2} \| \boldsymbol{m}\boldsymbol{m}^\mathsf{T} - \widehat{\boldsymbol{X}}_{\mathrm{Bayes}} \|_F^2 \wedge 1 \Big] = 0. \tag{1}$$

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

# Proof idea - Count the number of critical points

- Recall

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

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

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

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

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

## Proposition

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

# Proof idea - Count the number of critical points

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

where

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

$$- q\mu - \varphi\nu - a\tau - \left[ -\frac{\beta^2}{4}(1-q^2) + \frac{a}{2} - e \right] \gamma + \log I,$$

and

$$I = \int_{-\infty}^{\infty} \frac{1}{(2\pi \beta^2 q)^{1/2}} \exp \left\{ -\frac{(x - \beta \lambda \varphi)^2}{2\beta^2 q} \right.$$

$$\left. + \mu \tanh^2(x) + \nu \tanh(x) + \tau x \tanh(x) + \gamma \log[2\cosh(x)] \right\} dx.$$

# Proof idea - Count the number of critical points

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

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

▶ For any $U$ such that $T(U) > 0$, there could potentially be critical points of $\mathcal{F}_{\mathrm{TAP}}$ in $U$.

▶ For any $U$ such that $T(U) < 0$, there is no critical points of $\mathcal{F}_{\mathrm{TAP}}$ 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 $U \subseteq \mathbb{R}^3$,

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

▶ For any $U$ such that $T(U) > 0$, there could potentially be critical points of $\mathcal{F}_{\mathrm{TAP}}$ in $U$.

▶ For any $U$ such that $T(U) < 0$, there is no critical points of $\mathcal{F}_{\mathrm{TAP}}$ 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 $U \subseteq \mathbb{R}^3$,

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

▶ For any $U$ such that $T(U) > 0$, there could potentially be critical points of $\mathcal{F}_{\mathrm{TAP}}$ in $U$.

▶ For any $U$ such that $T(U) < 0$, there is no critical points of $\mathcal{F}_{\mathrm{TAP}}$ 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.
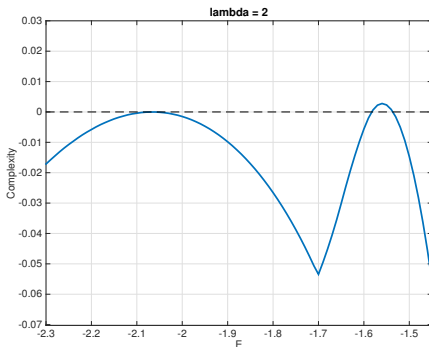
# Proof idea - Count the number of critical points

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

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

▶ For any $U$ such that $T(U) > 0$, there could potentially be critical points of $\mathcal{F}_{\mathrm{TAP}}$ in $U$.

▶ For any $U$ such that $T(U) < 0$, there is no critical points of $\mathcal{F}_{\mathrm{TAP}}$ 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.
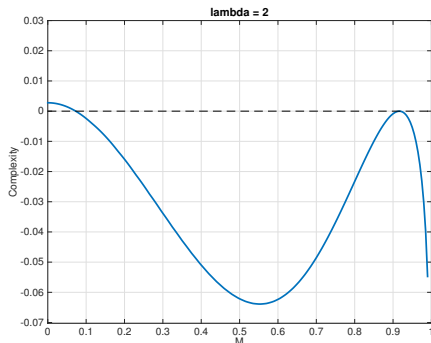
# Proof idea - the complexity function $S_\star$

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



- 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 $S_\star$

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



lambda = 2

- 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 $(q_\star, \varphi_\star, e_\star) \approx (Q(m_\star), M(m_\star), E(m_\star))$ for $\widehat{X}_{\text{Bayes}} \approx m_\star m_\star^{\mathsf{T}}$.
- $S_\star(q, \varphi, e) < 0$ 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 $S_\star$

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

- $S_\star(q_\star, \varphi_\star, e_\star) = 0$, where $(q_\star, \varphi_\star, e_\star) \approx (Q(m_\star), M(m_\star), E(m_\star))$ for $\widehat{X}_{\mathrm{Bayes}} \approx m_\star m_\star^\mathsf{T}$.
- $S_\star(q, \varphi, e) < 0$ 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.

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

Now suffice to show the key inequality.

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

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

- $S_\star(q_\star, \varphi_\star, e_\star) = 0$, where $(q_\star, \varphi_\star, e_\star) \approx (Q(\boldsymbol{m}_\star), M(\boldsymbol{m}_\star), E(\boldsymbol{m}_\star))$ for $\widehat{\boldsymbol{X}}_{\text{Bayes}} \approx \boldsymbol{m}_\star \boldsymbol{m}_\star^\mathsf{T}$.
- $S_\star(q, \varphi, e) < 0$ 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.

# Calculating the Crit: Kac-Rice formula

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

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

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

Then

$$\mathbb{E}[\mathrm{Crit}(T)] = \mathbb{E}\Big[ \int_T \big| \det \nabla^2 f(\boldsymbol{m}) \big| \cdot \delta(\nabla f(\boldsymbol{m})) \cdot \mathrm{d}\boldsymbol{m} \Big]$$

$$= \int_T \mathbb{E}\Big[ \big| \det \nabla^2 f(\boldsymbol{m}) \big| \,\Big|\, \nabla f(\boldsymbol{m}) = \boldsymbol{0} \Big] p_{\boldsymbol{m}}(\boldsymbol{0}) \mathrm{d}\boldsymbol{m}.$$

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

# Dealing with determinant of Hessian

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

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

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

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

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

# Dealing with determinant of Hessian

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

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

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

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

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

# Dealing with determinant of Hessian

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

$$\frac{1}{n} \log \mathbb{E}[|\det(\boldsymbol{H})|] = \frac{1}{n} \log \mathbb{E}\Big[ \prod_{i=1}^{n} |\lambda_i(\boldsymbol{H})| \Big] \approx \frac{1}{n} \log \Big[ \prod_{i=1}^{n} |\lambda_i(\boldsymbol{H})| \Big]$$

$$= \frac{1}{n} \sum_{i=1}^{n} \log |\lambda_i(\boldsymbol{H})| = \int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x) \approx \mathbb{E}\Big[ \int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x) \Big].$$

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

- ▶ Approximate equalities are due to concentration.
- ▶ The Stieltjes transform of $\mu_H$ can be approximately calculated using free probability theory.
- ▶ Once the Stieltjes transform of $\mu_H$ is known, the quantity $\mathbb{E}\Big[ \int_{\mathbb{R}}(\log |x|)\mu_H(\mathrm{d}x) \Big]$ can be computed.

# Dealing with determinant of Hessian

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

$$\frac{1}{n} \log \mathbb{E}[|\det(\boldsymbol{H})|] = \frac{1}{n} \log \mathbb{E}\Big[ \prod_{i=1}^{n} |\lambda_i(\boldsymbol{H})|\Big] \approx \frac{1}{n} \log \Big[ \prod_{i=1}^{n} |\lambda_i(\boldsymbol{H})|\Big]$$

$$= \frac{1}{n} \sum_{i=1}^{n} \log |\lambda_i(\boldsymbol{H})| = \int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x) \approx \mathbb{E}\Big[ \int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x)\Big].$$

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

▶ Approximate equalities are due to concentration.

▶ The Stieltjes transform of $\mu_{\boldsymbol{H}}$ can be approximately calculated using free probability theory.

▶ Once the Stieltjes transform of $\mu_{\boldsymbol{H}}$ is known, the quantity $\mathbb{E}\Big[ \int_{\mathbb{R}} (\log |x|) \mu_{\boldsymbol{H}}(\mathrm{d}x)\Big]$ can be computed.

# Free convolution of two distribution

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

$$g_A(z) = \int_{\mathbb{R}} \frac{1}{x - z} \mu_A(\mathrm{d}x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\lambda_i(A) - z}.$$

## Lemma (Due to free probability theory)

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

$$\mathbb{E}g_H(z) = \frac{1}{n} \sum_{i=1}^n \frac{1}{d_i - z - \mathbb{E}g_H(z)} + o_n(1). \tag{3}$$

# Free convolution of two distribution

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

$$g_A(z) = \int_{\mathbb{R}} \frac{1}{x - z} \mu_A(\mathrm{d}x) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{\lambda_i(A) - z}.$$

## Lemma (Due to free probability theory)

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

$$\mathbb{E}g_H(z) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{d_i - z - \mathbb{E}g_H(z)} + o_n(1). \tag{3}$$

# Dealing with determinant of Hessian

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

$$\frac{1}{n} \log \mathbb{E}[|\det(\boldsymbol{H})|] = \frac{1}{n} \log \mathbb{E}\Big[ \prod_{i=1}^{n} |\lambda_i(\boldsymbol{H})| \Big] \approx \frac{1}{n} \log \Big[ \prod_{i=1}^{n} |\lambda_i(\boldsymbol{H})| \Big]$$

$$= \frac{1}{n} \sum_{i=1}^{n} \log |\lambda_i(\boldsymbol{H})| = \int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x) \approx \mathbb{E}\Big[ \int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x) \Big].$$

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

- ▶ Approximate equalities are due to concentration.
- ▶ The Stieltjes transform of $\mu_{\boldsymbol{H}}$ can be approximately calculated using free probability theory.
- ▶ Once the Stieltjes transform of $\mu_{\boldsymbol{H}}$ is known, the quantity $\mathbb{E}\Big[ \int_{\mathbb{R}} (\log |x|) \mu_{\boldsymbol{H}}(\mathrm{d}x) \Big]$ can be computed.

# Calculate $\mathbb{E}[\int_{\mathbb{R}} \log|x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x)]$

▶ Define
$$B(t) = \mathbb{E} \int_{\mathbb{R}} \log(x - it)\mu_{\boldsymbol{H}}(\mathrm{d}x).$$

▶ We have
$$\Re B(0+) = \mathbb{E} \int_{\mathbb{R}} \log|x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x),$$
$$B'(t) = -i\mathbb{E} \int_{\mathbb{R}} [1/(x - it)]\mu_{\boldsymbol{H}}(\mathrm{d}x) = -i\mathbb{E}[g_{\boldsymbol{H}}(it)].$$

▶ We guess a formula
$$\tilde{B}(t) = \frac{1}{n} \sum_{i=1}^{n} \log(d_i - it - \mathbb{E}g_{\boldsymbol{H}}(it)) + \frac{1}{2}[\mathbb{E}g_{\boldsymbol{H}}(it)]^2.$$

Then $\tilde{B}(t)$ satisfy all the conditions that $B(t)$ approximately satisfy, so that $\tilde{B}(t) = B(t) + o_n(1)$.

▶ Hence
$$\frac{1}{n} \log \mathbb{E}[|\det(\boldsymbol{H})|] = \tilde{B}(0) + o_n(1).$$

# Calculate $\mathbb{E}[\int_{\mathbb{R}} \log|x| \cdot \mu_H(\mathrm{d}x)]$

- Define
$$B(t) = \mathbb{E} \int_{\mathbb{R}} \log(x - it) \mu_H(\mathrm{d}x).$$

- We have
$$\Re B(0+) = \mathbb{E} \int_{\mathbb{R}} \log|x| \cdot \mu_H(\mathrm{d}x),$$
$$B'(t) = -i\mathbb{E} \int_{\mathbb{R}} [1/(x - it)] \mu_H(\mathrm{d}x) = -i\mathbb{E}[g_H(it)].$$

- We guess a formula
$$\tilde{B}(t) = \frac{1}{n} \sum_{i=1}^{n} \log(d_i - it - \mathbb{E}g_H(it)) + \frac{1}{2}[\mathbb{E}g_H(it)]^2.$$

Then $\tilde{B}(t)$ satisfy all the conditions that $B(t)$ approximately satisfy, so that $\tilde{B}(t) = B(t) + o_n(1)$.

- Hence
$$\frac{1}{n} \log \mathbb{E}[|\det(H)|] = \tilde{B}(0) + o_n(1).$$

# Calculate $\mathbb{E}[\int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x)]$

▶ Define
$$B(t) = \mathbb{E} \int_{\mathbb{R}} \log(x - it) \mu_{\boldsymbol{H}}(\mathrm{d}x).$$

▶ We have
$$\Re B(0+) = \mathbb{E} \int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x),$$
$$B'(t) = -i\mathbb{E} \int_{\mathbb{R}} [1/(x - it)] \mu_{\boldsymbol{H}}(\mathrm{d}x) = -i\mathbb{E}[g_{\boldsymbol{H}}(it)].$$

▶ We guess a formula
$$\tilde{B}(t) = \frac{1}{n} \sum_{i=1}^{n} \log(d_i - it - \mathbb{E}g_{\boldsymbol{H}}(it)) + \frac{1}{2}[\mathbb{E}g_{\boldsymbol{H}}(it)]^2.$$

Then $\tilde{B}(t)$ satisfy all the conditions that $B(t)$ approximately satisfy, so that $\tilde{B}(t) = B(t) + o_n(1)$.

▶ Hence
$$\frac{1}{n} \log \mathbb{E}[|\det(\boldsymbol{H})|] = \tilde{B}(0) + o_n(1).$$

# Calculate $\mathbb{E}[\int_{\mathbb{R}} \log|x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x)]$

► Define
$$B(t) = \mathbb{E}\int_{\mathbb{R}} \log(x - it)\mu_{\boldsymbol{H}}(\mathrm{d}x).$$

► We have
$$\Re B(0+) = \mathbb{E}\int_{\mathbb{R}} \log|x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x),$$
$$B'(t) = -i\mathbb{E}\int_{\mathbb{R}}[1/(x - it)]\mu_{\boldsymbol{H}}(\mathrm{d}x) = -i\mathbb{E}[g_{\boldsymbol{H}}(it)].$$

► We guess a formula
$$\tilde{B}(t) = \frac{1}{n}\sum_{i=1}^{n} \log(d_i - it - \mathbb{E}g_{\boldsymbol{H}}(it)) + \frac{1}{2}[\mathbb{E}g_{\boldsymbol{H}}(it)]^2.$$

Then $\tilde{B}(t)$ satisfy all the conditions that $B(t)$ approximately satisfy, so that $\tilde{B}(t) = B(t) + o_n(1)$.

► Hence
$$\frac{1}{n}\log\mathbb{E}[|\det(\boldsymbol{H})|] = \tilde{B}(0) + o_n(1).$$

# Dealing with determinant of Hessian

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

$$\frac{1}{n} \log \mathbb{E}[|\det(\boldsymbol{H})|] = \frac{1}{n} \log \mathbb{E}\Big[\prod_{i=1}^{n} |\lambda_i(\boldsymbol{H})|\Big] \approx \frac{1}{n} \log \Big[\prod_{i=1}^{n} |\lambda_i(\boldsymbol{H})|\Big]$$

$$= \frac{1}{n} \sum_{i=1}^{n} \log |\lambda_i(\boldsymbol{H})| = \int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x) \approx \mathbb{E}\Big[\int_{\mathbb{R}} \log |x| \cdot \mu_{\boldsymbol{H}}(\mathrm{d}x)\Big].$$

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

- Approximate equalities are due to concentration.
- The Stieltjes transform of $\mu_{\boldsymbol{H}}$ can be approximately calculated using free probability theory.
- Once the Stieltjes transform of $\mu_{\boldsymbol{H}}$ is known, the quantity $\mathbb{E}\Big[\int_{\mathbb{R}}(\log |x|)\mu_{\boldsymbol{H}}(\mathrm{d}x)\Big]$ can be computed.

# Summary

▶ TAP free energy is accurate for $\mathbb{Z}_2$ synchronization.

▶ Can be generalized to topic modeling, low rank matrix estimation, compressed sensing, etc...

▶ It is interesting to study and apply variational inference beyond 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...

▶ It is interesting to study and apply variational inference beyond 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...

▶ It is interesting to study and apply variational inference beyond mean field.