

# Mean Field Asymptotics in High-dimensional Statistics

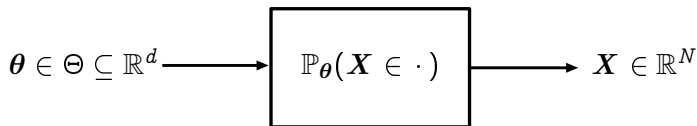
Andrea Montanari

Stanford University

August 8, 2018

## The big picture

# Statistical estimation

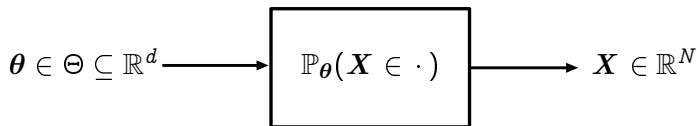


**Estimator:**  $\hat{\theta} : \mathbb{R}^N \rightarrow \mathbb{R}^d \quad \mathbf{X} \mapsto \hat{\theta}(\mathbf{X})$

Amount of data :  $N$

Model complexity :  $d$

# Statistical estimation



**Estimator:**  $\hat{\theta} : \mathbb{R}^N \rightarrow \mathbb{R}^d \quad \mathbf{X} \mapsto \hat{\theta}(\mathbf{X})$

Amount of data :  $N$

Model complexity :  $d$

Regimes:  $\theta \in \Theta \subseteq \mathbb{R}^d$ ,  $\mathbf{X} \in \mathbb{R}^N$

Classical – Data rich

$$N \gg d \Rightarrow \|\hat{\theta}(\mathbf{X}) - \theta_0\| \xrightarrow{\mathbb{P}} 0$$

[Fisher, Le Cam, Huber, ... 1920-...]

Modern – Data poor – High-dimensional

$$d \gg N \gg \text{eff.dim}(\Theta) \Rightarrow \|\hat{\theta}(\mathbf{X}) - \theta_0\| \xrightarrow{\mathbb{P}} 0$$

[Donoho, Candés, Tao, ... 2006-...]

Noisy high-dimensional

$$N \asymp d \asymp \text{eff.dim}(\Theta) \Rightarrow (\hat{\theta}(\mathbf{X}), \theta_0) \xrightarrow{d} \text{Non-trivial limit}$$

Regimes:  $\theta \in \Theta \subseteq \mathbb{R}^d$ ,  $\mathbf{X} \in \mathbb{R}^N$

Classical – Data rich

$$N \gg d \Rightarrow \|\hat{\theta}(\mathbf{X}) - \theta_0\| \xrightarrow{\mathbb{P}} 0$$

[Fisher, Le Cam, Huber, ... 1920-...]

Modern – Data poor – High-dimensional

$$d \gg N \gg \text{eff.dim}(\Theta) \Rightarrow \|\hat{\theta}(\mathbf{X}) - \theta_0\| \xrightarrow{\mathbb{P}} 0$$

[Donoho, Candés, Tao, ... 2006-...]

Noisy high-dimensional

$$N \asymp d \asymp \text{eff.dim}(\Theta) \Rightarrow (\hat{\theta}(\mathbf{X}), \theta_0) \xrightarrow{d} \text{Non-trivial limit}$$

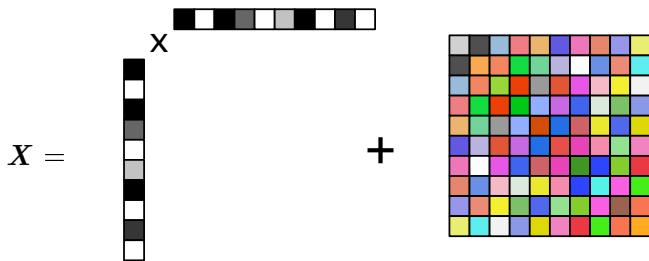
# Outline

- 1 A simple model
- 2 Bayes optimal estimation
- 3 Efficient algorithms
- 4 Generic/robust algorithms
- 5 Conclusion

## A simple model



# Rank-one matrix estimation



$$X = \frac{\lambda}{n} \theta_0 \theta_0^T + W$$

# Rank-one matrix estimation

$$\mathbf{X} = \frac{\lambda}{n} \boldsymbol{\theta}_0 \boldsymbol{\theta}_0^\top + \mathbf{W}$$

▶ **Parameter:**  $\boldsymbol{\theta}_0 \in \mathbb{R}^n, \|\boldsymbol{\theta}_0\|_2^2 = n + o(n)$

▶ **Data:**  $\mathbf{X} \in \mathbb{R}^{n \times n}$

▶ **Noise:**  $\mathbf{W} \sim \text{GOE}(n)$

$$(\mathbf{W} = \mathbf{W}^\top, W_{ii} \sim \text{N}(0, 2/n), W_{ij} \sim \text{N}(0, 1/n))$$

## Structured parameter

$$\mathbf{X} = \frac{\lambda}{n} \boldsymbol{\theta}_0 \boldsymbol{\theta}_0^\top + \mathbf{W}$$

$$\hat{p}_\theta \equiv \frac{1}{n} \sum_{i=1}^n \delta_{\theta_i},$$

$$\Theta(\varepsilon) \equiv \left\{ \boldsymbol{\theta} \in \mathbb{R}^n : \hat{p}_\theta = p_\varepsilon \equiv \varepsilon \delta_{a_+} + (1 - \varepsilon) \delta_{-a_-} \right\},$$

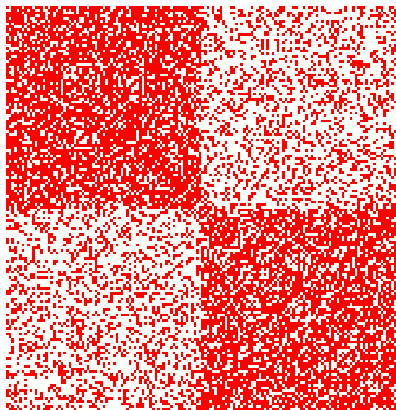
$$\int x p_\varepsilon(dx) = 0 \quad \int x^2 p_\varepsilon(dx) = 1.$$

$\varepsilon = 0.5$ : ' $\mathbb{Z}_2$ -synchronization'

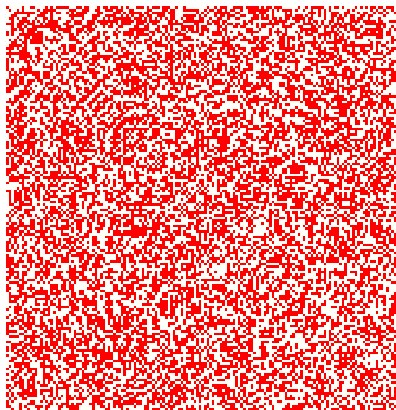
$$X = \frac{\lambda}{n} \theta_0 \theta_0^\top + W$$

$$\theta_0 \in \{+1, -1\}^n; \quad \sum_{i=1}^n \theta_{0,i} = 0.$$

# Example



## The same example



# Goal

$$\mathbf{X} = \frac{\lambda}{n} \boldsymbol{\theta}_0 \boldsymbol{\theta}_0^\top + \mathbf{W}$$

**Estimator:**

$$\hat{\boldsymbol{\theta}} : \mathbf{X} \mapsto \hat{\boldsymbol{\theta}}(\mathbf{X}) \in \mathbb{R}^n$$

**Accuracy:**

Minimax: 
$$\min_{\boldsymbol{\theta}_0 \in \Theta(\varepsilon)} \mathbb{E}_{\mathbf{X}} \left\{ \frac{|\langle \hat{\boldsymbol{\theta}}(\mathbf{X}), \boldsymbol{\theta}_0 \rangle|}{\|\hat{\boldsymbol{\theta}}(\mathbf{X})\|_2 \|\boldsymbol{\theta}_0\|_2} \right\},$$

Bayesian: 
$$\mathbb{E}_{\boldsymbol{\theta}_0 \sim \text{Unif}(\Theta(\varepsilon))} \mathbb{E}_{\mathbf{X}} \left\{ \frac{|\langle \hat{\boldsymbol{\theta}}(\mathbf{X}), \boldsymbol{\theta}_0 \rangle|}{\|\hat{\boldsymbol{\theta}}(\mathbf{X})\|_2 \|\boldsymbol{\theta}_0\|_2} \right\}$$

Minimax  $\approx$  Bayesian (by symmetry)

# Goal

$$\mathbf{X} = \frac{\lambda}{n} \boldsymbol{\theta}_0 \boldsymbol{\theta}_0^\top + \mathbf{W}$$

**Estimator:**

$$\hat{\boldsymbol{\theta}} : \mathbf{X} \mapsto \hat{\boldsymbol{\theta}}(\mathbf{X}) \in \mathbb{R}^n$$

**Accuracy:**

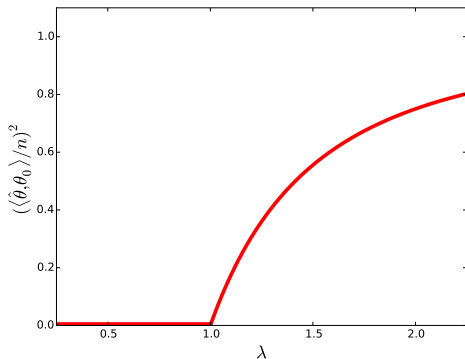
Minimax: 
$$\min_{\boldsymbol{\theta}_0 \in \Theta(\varepsilon)} \mathbb{E}_{\mathbf{X}} \left\{ \frac{|\langle \hat{\boldsymbol{\theta}}(\mathbf{X}), \boldsymbol{\theta}_0 \rangle|}{\|\hat{\boldsymbol{\theta}}(\mathbf{X})\|_2 \|\boldsymbol{\theta}_0\|_2} \right\},$$

Bayesian: 
$$\mathbb{E}_{\boldsymbol{\theta}_0 \sim \text{Unif}(\Theta(\varepsilon))} \mathbb{E}_{\mathbf{X}} \left\{ \frac{|\langle \hat{\boldsymbol{\theta}}(\mathbf{X}), \boldsymbol{\theta}_0 \rangle|}{\|\hat{\boldsymbol{\theta}}(\mathbf{X})\|_2 \|\boldsymbol{\theta}_0\|_2} \right\}$$

Minimax  $\approx$  Bayesian (by symmetry)



## Example: Principal Component Analysis



$$\hat{\theta}^{\text{PCA}}(\mathbf{X}) = \sqrt{n} v_1(\mathbf{X}) \quad (\text{principal eigenvector})$$

## Example: Principal Component Analysis

$$\mathbf{X} = \frac{\lambda}{n} \boldsymbol{\theta}_0 \boldsymbol{\theta}_0^\top + \mathbf{W}, \quad \hat{\boldsymbol{\theta}}^{\text{PCA}}(\mathbf{X}) = \mathbf{v}_1(\mathbf{X}).$$

Theorem (Baik, Ben Arous, Pécché, 2005)

$$\lim_{n \rightarrow \infty} \frac{|\langle \hat{\boldsymbol{\theta}}^{\text{PCA}}(\mathbf{X}), \boldsymbol{\theta}_0 \rangle|}{\|\hat{\boldsymbol{\theta}}^{\text{PCA}}(\mathbf{X})\|_2 \|\boldsymbol{\theta}_0\|_2} = \sqrt{\left(1 - \frac{1}{\lambda^2}\right)_+}$$

Can we do better?

## Example: Principal Component Analysis

$$\mathbf{X} = \frac{\lambda}{n} \boldsymbol{\theta}_0 \boldsymbol{\theta}_0^\top + \mathbf{W}, \quad \hat{\boldsymbol{\theta}}^{\text{PCA}}(\mathbf{X}) = \mathbf{v}_1(\mathbf{X}).$$

Theorem (Baik, Ben Arous, P  ch  , 2005)

$$\lim_{n \rightarrow \infty} \frac{|\langle \hat{\boldsymbol{\theta}}^{\text{PCA}}(\mathbf{X}), \boldsymbol{\theta}_0 \rangle|}{\|\hat{\boldsymbol{\theta}}^{\text{PCA}}(\mathbf{X})\|_2 \|\boldsymbol{\theta}_0\|_2} = \sqrt{\left(1 - \frac{1}{\lambda^2}\right)_+}$$

Can we do better?

# Agenda

- ▶ Bayes optimal estimation
- ▶ Efficient algorithms
- ▶ Generic/robust algorithms

## Bayes optimal estimation

# Bayesian estimation ( $\approx$ minimax)

## Posterior

$$p_{\text{Bayes}}(d\boldsymbol{\theta} | \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \exp \left\{ \frac{\lambda}{2} \langle \boldsymbol{\theta}, \mathbf{X} \boldsymbol{\theta} \rangle - \frac{\lambda^2}{4n} \|\boldsymbol{\theta}\|_2^4 \right\} \prod_{i=1}^n p_{\varepsilon}(d\theta_i).$$

## Estimator

$$\hat{\boldsymbol{\theta}}^{\text{Bayes}}(\mathbf{X}) = \mathbb{E}\{\boldsymbol{\theta} | \mathbf{X}\}.$$

[For  $\varepsilon = 0.5$  need to break symmetry...]

# Bayesian estimation: Analysis

## Posterior

$$p_{\text{Bayes}}(d\boldsymbol{\theta} | \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \exp \left\{ \frac{\lambda}{2} \langle \boldsymbol{\theta}, \mathbf{X} \boldsymbol{\theta} \rangle - \frac{\lambda^2}{4n} \|\boldsymbol{\theta}\|_2^4 \right\} \prod_{i=1}^n p_{\varepsilon}(d\theta_i).$$

- ▶ Generalized spin-glass model
- ▶ Replica symmetric [Korada, Macris 2007; M. 2008]

$$\mathbf{X} \sim p_{\text{Bayes}}(\cdot), \quad (\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) \sim p_{\text{Bayes}}(\cdot | \mathbf{X}) \otimes p_{\text{Bayes}}(\cdot | \mathbf{X})$$

$$\frac{1}{n} \langle \boldsymbol{\theta}^1, \boldsymbol{\theta}^2 \rangle \xrightarrow{\text{P}} q_* \quad (\text{nonrandom})$$

Intuition: for  $\boldsymbol{\theta} \sim p_{\text{Bayes}}(\cdot | \mathbf{X})$ ,  $\theta_i, \theta_j$ ,  $i \neq j$  approx independent

# Bayesian estimation: Analysis

## Posterior

$$p_{\text{Bayes}}(d\boldsymbol{\theta} | \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \exp \left\{ \frac{\lambda}{2} \langle \boldsymbol{\theta}, \mathbf{X} \boldsymbol{\theta} \rangle - \frac{\lambda^2}{4n} \|\boldsymbol{\theta}\|_2^4 \right\} \prod_{i=1}^n p_{\varepsilon}(d\theta_i).$$

- ▶ Generalized spin-glass model
- ▶ Replica symmetric [Korada, Macris 2007; M. 2008]

$$\mathbf{X} \sim p_{\text{Bayes}}(\cdot), \quad (\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) \sim p_{\text{Bayes}}(\cdot | \mathbf{X}) \otimes p_{\text{Bayes}}(\cdot | \mathbf{X})$$

$$\frac{1}{n} \langle \boldsymbol{\theta}^1, \boldsymbol{\theta}^2 \rangle \xrightarrow{\text{P}} q_* \quad (\text{nonrandom})$$

Intuition: for  $\boldsymbol{\theta} \sim p_{\text{Bayes}}(\cdot | \mathbf{X})$ ,  $\theta_i, \theta_j$ ,  $i \neq j$  approx independent



# Bayesian estimation: Analysis

## Posterior

$$p_{\text{Bayes}}(d\boldsymbol{\theta} | \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \exp \left\{ \frac{\lambda}{2} \langle \boldsymbol{\theta}, \mathbf{X} \boldsymbol{\theta} \rangle - \frac{\lambda^2}{4n} \|\boldsymbol{\theta}\|_2^4 \right\} \prod_{i=1}^n p_{\varepsilon}(d\theta_i).$$

- ▶ Generalized spin-glass model
- ▶ Replica symmetric [Korada, Macris 2007; M. 2008]

$$\mathbf{X} \sim p_{\text{Bayes}}(\cdot), \quad (\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) \sim p_{\text{Bayes}}(\cdot | \mathbf{X}) \otimes p_{\text{Bayes}}(\cdot | \mathbf{X})$$

$$\frac{1}{n} \langle \boldsymbol{\theta}^1, \boldsymbol{\theta}^2 \rangle \xrightarrow{\text{P}} q_* \quad (\text{nonrandom})$$

Intuition: for  $\boldsymbol{\theta} \sim p_{\text{Bayes}}(\cdot | \mathbf{X})$ ,  $\theta_i, \theta_j$ ,  $i \neq j$  approx independent

# Bayesian estimation: Analysis

## Posterior

$$p_{\text{Bayes}}(d\boldsymbol{\theta} | \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \exp \left\{ \frac{\lambda}{2} \langle \boldsymbol{\theta}, \mathbf{X} \boldsymbol{\theta} \rangle - \frac{\lambda^2}{4n} \|\boldsymbol{\theta}\|_2^4 \right\} \prod_{i=1}^n p_{\varepsilon}(d\theta_i).$$

- ▶ Generalized spin-glass model
- ▶ Replica symmetric [Korada, Macris 2007; M. 2008]

$$\mathbf{X} \sim p_{\text{Bayes}}(\cdot), \quad (\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) \sim p_{\text{Bayes}}(\cdot | \mathbf{X}) \otimes p_{\text{Bayes}}(\cdot | \mathbf{X})$$

$$\frac{1}{n} \langle \boldsymbol{\theta}^1, \boldsymbol{\theta}^2 \rangle \xrightarrow{\text{P}} q_* \quad (\text{nonrandom})$$

Intuition: for  $\boldsymbol{\theta} \sim p_{\text{Bayes}}(\cdot | \mathbf{X})$ ,  $\theta_i, \theta_j$ ,  $i \neq j$  approx independent

# Bayesian estimation: Analysis

## Posterior

$$p_{\text{Bayes}}(d\boldsymbol{\theta} | \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \exp \left\{ \frac{\lambda}{2} \langle \boldsymbol{\theta}, \mathbf{X} \boldsymbol{\theta} \rangle - \frac{\lambda^2}{4n} \|\boldsymbol{\theta}\|_2^4 \right\} \prod_{i=1}^n p_{\varepsilon}(d\theta_i).$$

- ▶ Generalized spin-glass model
- ▶ Replica symmetric [Korada, Macris 2007; M. 2008]

$$\mathbf{X} \sim p_{\text{Bayes}}(\cdot), \quad (\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) \sim p_{\text{Bayes}}(\cdot | \mathbf{X}) \otimes p_{\text{Bayes}}(\cdot | \mathbf{X})$$

$$\frac{1}{n} \langle \boldsymbol{\theta}^1, \boldsymbol{\theta}^2 \rangle \xrightarrow{\text{P}} q_* \quad (\text{nonrandom})$$

Intuition: for  $\boldsymbol{\theta} \sim p_{\text{Bayes}}(\cdot | \mathbf{X})$ ,  $\theta_i, \theta_j$ ,  $i \neq j$  approx independent

# Bayesian estimation: Analysis

Theorem (Lelarge, Miolane, 2017; Barbier et al. 2016)

Consider the rank-one model with  $\theta_{0,i} \sim p$  (+ conditions). Let  $I(\gamma) = \mathbb{E} \log \frac{dp_{Y|X_0}}{dp_Y}(Y, X_0)$  for  $(X_0, G) \sim p \otimes \mathcal{N}(0, 1)$ ,  $Y = \sqrt{\gamma}X_0 + G$ , and define

$$\Psi(\gamma, \lambda) \equiv \frac{\lambda^2}{4} + \frac{\gamma^2}{4\lambda} - \frac{\gamma}{2} + I(\gamma),$$
$$\gamma_{\text{Bayes}}(\lambda) \equiv \arg \max_{\gamma \geq 0} \Psi(\gamma, \lambda).$$

Then

$$\lim_{n \rightarrow \infty} \frac{|\langle \hat{\theta}^{\text{Bayes}}(\mathbf{X}), \theta_0 \rangle|}{\|\hat{\theta}^{\text{Bayes}}(\mathbf{X})\|_2 \|\theta_0\|_2} = \frac{\sqrt{\gamma_{\text{Bayes}}(\lambda)}}{\lambda}.$$

[see also: Deshpande, M., 2014; Deshpande, Abbe, M., 2016; Barbier, Dia, Macris, Krzakala, Lesieur, Zdeborová, 2016; Krzakala, Xu, Zdeborová, 2016; El Alaoui, Jordan, 2018; ...]

# Bayesian estimation: Analysis

Theorem (Lelarge, Miolane, 2017; Barbier et al. 2016)

Consider the rank-one model with  $\theta_{0,i} \sim p$  (+ conditions). Let  $I(\gamma) = \mathbb{E} \log \frac{dp_{Y|X_0}}{dp_Y}(Y, X_0)$  for  $(X_0, G) \sim p \otimes \mathcal{N}(0, 1)$ ,  $Y = \sqrt{\gamma}X_0 + G$ , and define

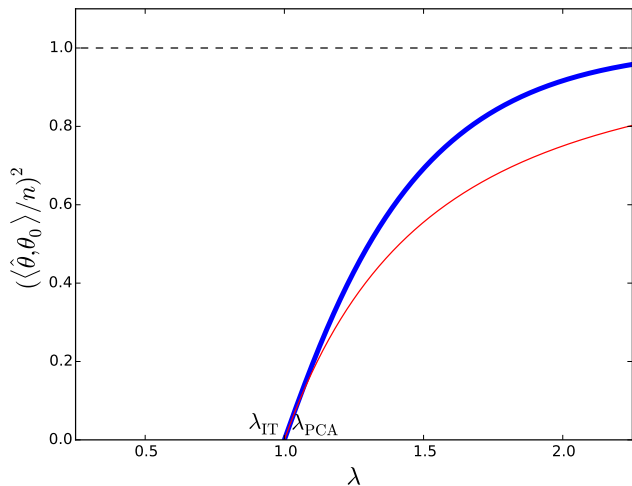
$$\Psi(\gamma, \lambda) \equiv \frac{\lambda^2}{4} + \frac{\gamma^2}{4\lambda} - \frac{\gamma}{2} + I(\gamma),$$
$$\gamma_{\text{Bayes}}(\lambda) \equiv \arg \max_{\gamma \geq 0} \Psi(\gamma, \lambda).$$

Then

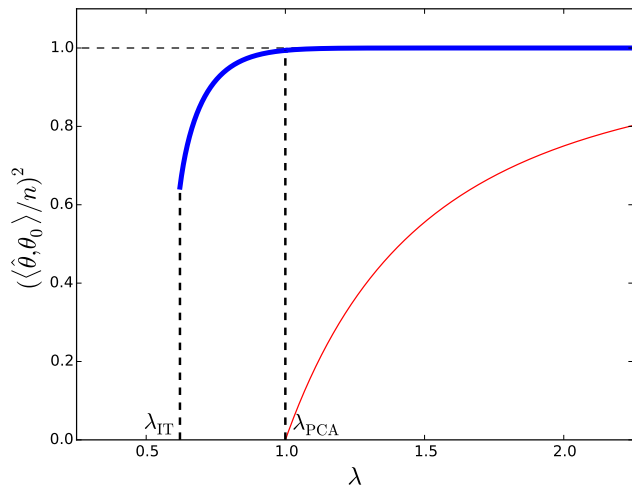
$$\lim_{n \rightarrow \infty} \frac{|\langle \hat{\theta}^{\text{Bayes}}(\mathbf{X}), \theta_0 \rangle|}{\|\hat{\theta}^{\text{Bayes}}(\mathbf{X})\|_2 \|\theta_0\|_2} = \frac{\sqrt{\gamma_{\text{Bayes}}(\lambda)}}{\lambda}.$$

[see also: Deshpande, M., 2014; Deshpande, Abbe, M., 2016; Barbier, Dia, Macris, Krzakala, Lesieur, Zdeborová, 2016; Krzakala, Xu, Zdeborová, 2016; El Alaoui, Jordan, 2018; ...]

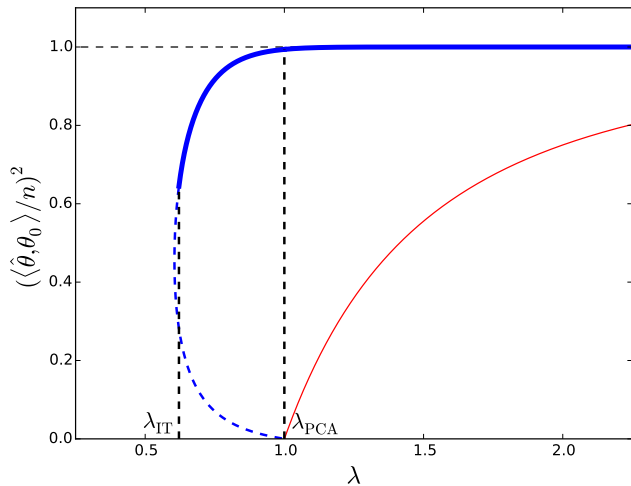
$$\varepsilon = 0.5$$



$$\varepsilon = 0.025$$

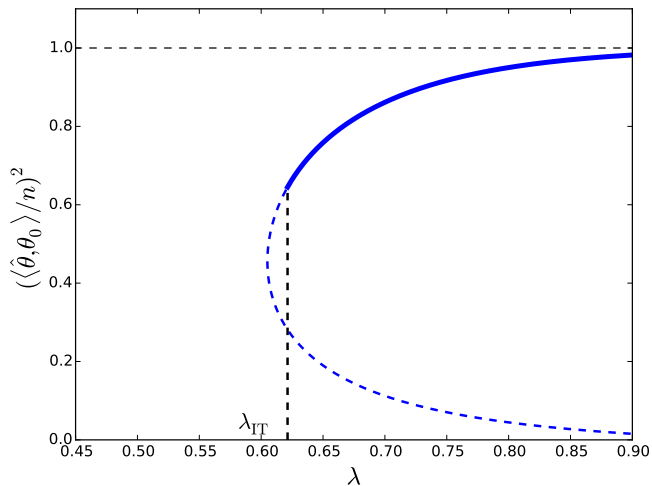


$$\varepsilon = 0.025$$



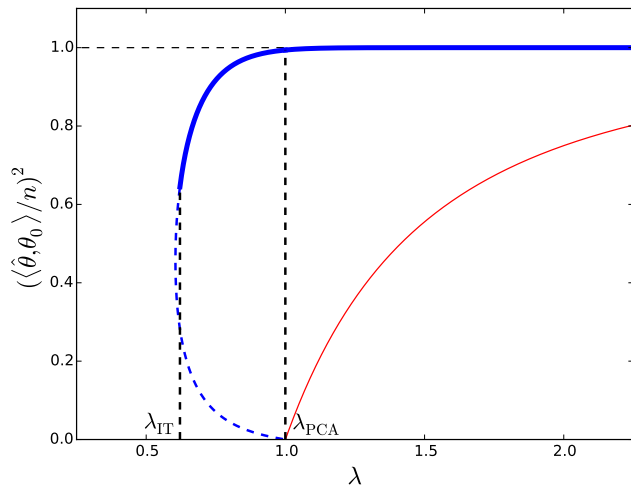


$$\varepsilon = 0.025$$



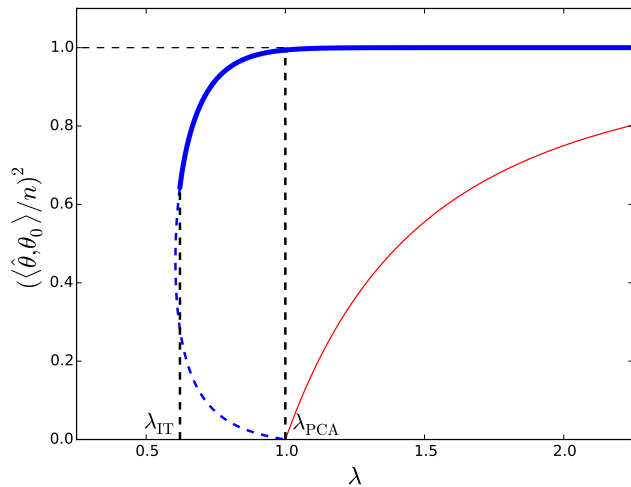
## Efficient algorithms

Are we happy about this?



Not really...

Are we happy about this?



Not really...

# Approximate Message Passing (AMP)

**Data:**  $\mathbf{X} \in \mathbb{R}^{n \times n}$

**Initialization:**  $\hat{\boldsymbol{\theta}}^0 \in \mathbb{R}^n$

**Iteration:**

$$\begin{aligned}\hat{\boldsymbol{\theta}}^{t+1} &= \mathbf{X} \mathbf{f}_t(\hat{\boldsymbol{\theta}}^t) - \mathbf{b}_t \mathbf{f}_{t-1}(\hat{\boldsymbol{\theta}}^{t-1}), \\ \mathbf{b}_t &\equiv \frac{1}{n} \operatorname{div} \mathbf{f}_t(\hat{\boldsymbol{\theta}}^t).\end{aligned}$$

# Approximate Message Passing (AMP)

$$\hat{\theta}^{t+1} = \mathbf{X} \mathbf{f}_t(\hat{\theta}^t) - \mathbf{b}_t \mathbf{f}_{t-1}(\hat{\theta}^{t-1}),$$
$$\mathbf{b}_t \equiv \frac{1}{n} \operatorname{div} \mathbf{f}_t(\hat{\theta}^t).$$

- ▶ Generalization of iterative scheme to solve TAP equations for mean field spin glasses  
[Thouless, Anderson, Palmer, 1977; Kabashima 2003]
- ▶ Closely related to loopy belief propagation  
[Gallager, 1962; Pearl, 1986; ...]

# Analysis (heuristic!!!)

$$\hat{\theta}^t \approx \mu_t \theta_0 + \sigma_t g, \quad g = \mathbf{N}(\mathbf{0}, \mathbf{I}_{n \times n})$$

$$X f_t(\hat{\theta}^t) = \frac{\lambda}{n} \langle \theta_0, f_t(\hat{\theta}^t) \rangle \theta_0 + W f_t(\hat{\theta}^t)$$

$$\approx \mu_{t+1} \theta_0 + \sigma_{t+1} g'$$

Separable  $f_t(x) = (f_t(x_1), \dots, f_t(x_n))$ ,  $((X_0, G) \sim \hat{p}_{\theta_0}(\cdot) \otimes \mathbf{N}(0, 1))$

$$\mu_{t+1} = \lambda \mathbb{E}\{X_0 f_t(\mu_t X_0 + \sigma_t G)\},$$

$$\sigma_{t+1}^2 = \mathbb{E}\{f_t(\mu_t X_0 + \sigma_t G)^2\}.$$

# Analysis (heuristic!!!)

$$\hat{\theta}^t \approx \mu_t \theta_0 + \sigma_t g, \quad g = \mathbf{N}(\mathbf{0}, \mathbf{I}_{n \times n})$$

$$\mathbf{X} f_t(\hat{\theta}^t) = \frac{\lambda}{n} \langle \theta_0, f_t(\hat{\theta}^t) \rangle \theta_0 + \mathbf{W} f_t(\hat{\theta}^t)$$

$$\approx \mu_{t+1} \theta_0 + \sigma_{t+1} g'$$

Separable  $f_t(x) = (f_t(x_1), \dots, f_t(x_n))$ ,  $((X_0, G) \sim \hat{p}_{\theta_0}(\cdot) \otimes \mathbf{N}(0, 1))$

$$\mu_{t+1} = \lambda \mathbb{E}\{X_0 f_t(\mu_t X_0 + \sigma_t G)\},$$

$$\sigma_{t+1}^2 = \mathbb{E}\{f_t(\mu_t X_0 + \sigma_t G)^2\}.$$



## Analysis (heuristic!!!)

$$\hat{\theta}^t \approx \mu_t \theta_0 + \sigma_t g, \quad g = \mathbf{N}(0, \mathbf{I}_{n \times n})$$

$$\mathbf{X} f_t(\hat{\theta}^t) = \frac{\lambda}{n} \langle \theta_0, f_t(\hat{\theta}^t) \rangle \theta_0 + \mathbf{W} f_t(\hat{\theta}^t)$$

$$\approx \mu_{t+1} \theta_0 + \sigma_{t+1} g'$$

Separable  $f_t(\mathbf{x}) = (f_t(x_1), \dots, f_t(x_n))$ ,  $((X_0, G) \sim \hat{p}_{\theta_0}(\cdot) \otimes \mathbf{N}(0, 1))$

$$\mu_{t+1} = \lambda \mathbb{E}\{X_0 f_t(\mu_t X_0 + \sigma_t G)\},$$

$$\sigma_{t+1}^2 = \mathbb{E}\{f_t(\mu_t X_0 + \sigma_t G)^2\}.$$

# Consequence

$$\begin{aligned}\mu_{t+1} &= \lambda \mathbb{E}\{X_0 f_t(\mu_t X_0 + \sigma_t G)\}, \\ \sigma_{t+1}^2 &= \mathbb{E}\{f_t(\mu_t X_0 + \sigma_t G)^2\}.\end{aligned}$$

## Bayes-AMP algorithm

- ▶ Optimal denoiser:  $f_t(y) = \mathbb{E}\{X_0 | \mu_t X_0 + \sigma_t G = y\}$ .
- ▶ Initialization  $\hat{\theta}^0 = c \cdot \hat{\theta}^{\text{PCA}}(\mathbf{X})$

# AMP estimation

Theorem (M., Venkatarmanan, 2018)

Consider the rank-one model with  $\theta_{0,i} \sim p$  (+ conditions),  $\lambda > 1$  and define

$$\Psi(\gamma, \lambda) \equiv \frac{\lambda^2}{4} + \frac{\gamma^2}{4\lambda} - \frac{\gamma}{2} + I(\gamma),$$
$$\gamma_{\text{ALG}}(\lambda) \equiv \min \{ \gamma > 0 : \partial_{\gamma} \Psi(\gamma, \lambda) = 0 \}.$$

(Recall that  $\gamma_{\text{Bayes}}(\lambda) \equiv \arg \max_{\gamma \geq 0} \Psi(\gamma, \lambda)$ ).

Then

$$\lim_{t \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{|\langle \hat{\theta}^t(\mathbf{X}), \theta_0 \rangle|}{\|\hat{\theta}^t(\mathbf{X})\|_2 \|\theta_0\|_2} = \frac{\sqrt{\gamma_{\text{ALG}}(\lambda)}}{\lambda}.$$

[see also: Bolthausen, 2014; Bayati, M., 2011; Bayati, Lelarge, M., 2015; Berthier, M., Nguyen 2018; ...]

# AMP estimation

Theorem (M., Venkatarmanan, 2018)

Consider the rank-one model with  $\theta_{0,i} \sim p$  (+ conditions),  $\lambda > 1$  and define

$$\Psi(\gamma, \lambda) \equiv \frac{\lambda^2}{4} + \frac{\gamma^2}{4\lambda} - \frac{\gamma}{2} + I(\gamma),$$
$$\gamma_{\text{ALG}}(\lambda) \equiv \min \{ \gamma > 0 : \partial_{\gamma} \Psi(\gamma, \lambda) = 0 \}.$$

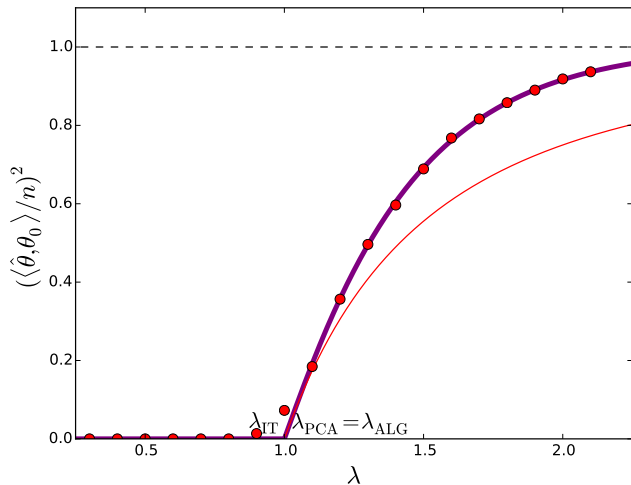
(Recall that  $\gamma_{\text{Bayes}}(\lambda) \equiv \arg \max_{\gamma \geq 0} \Psi(\gamma, \lambda)$ ).

Then

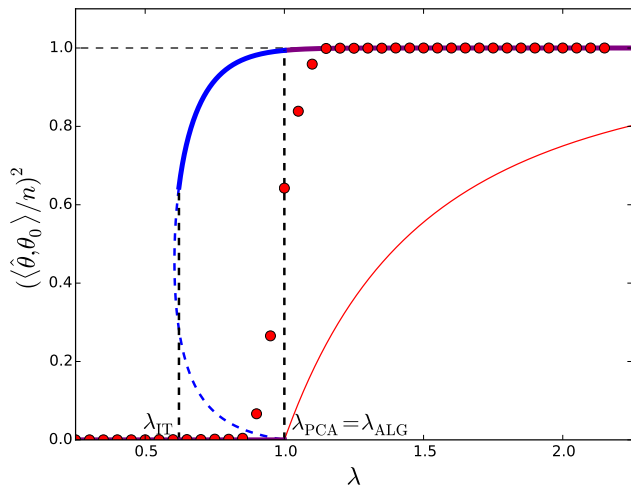
$$\lim_{t \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{|\langle \hat{\theta}^t(\mathbf{X}), \theta_0 \rangle|}{\|\hat{\theta}^t(\mathbf{X})\|_2 \|\theta_0\|_2} = \frac{\sqrt{\gamma_{\text{ALG}}(\lambda)}}{\lambda}.$$

[see also: **Bolthausen, 2014**; Bayati, M., 2011; Bayati, Lelarge, M., 2015; Berthier, M., Nguyen 2018; ...]

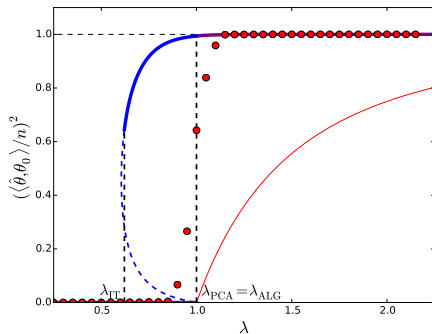
$$\varepsilon = 0.5$$



$$\varepsilon = 0.025$$



# A bold conjecture



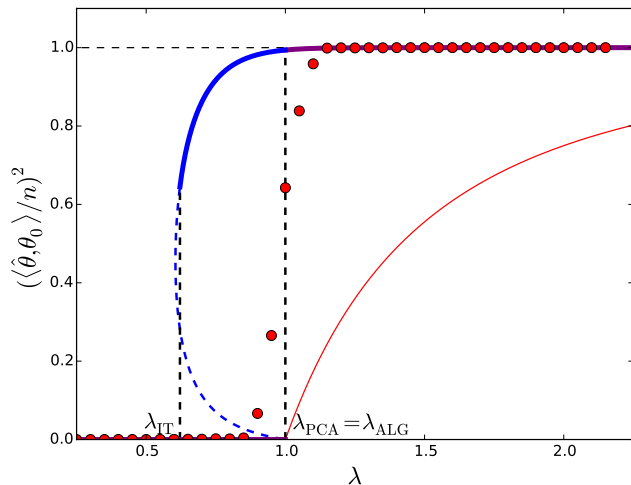
## Conjecture

*No polytime algorithm can beat AMP.*

## Generic/robust algorithms

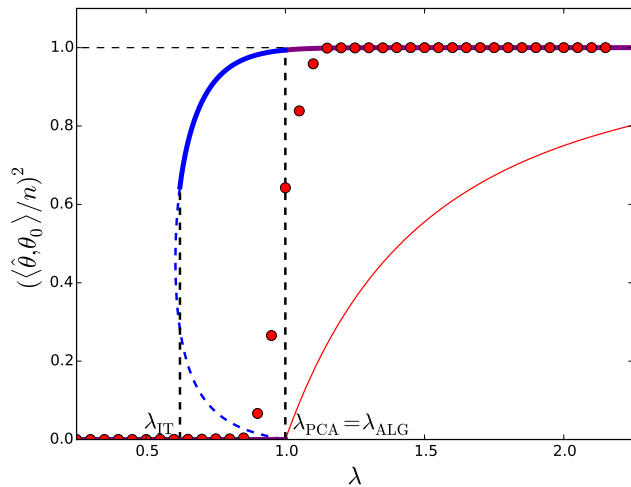


# Are we happy about this?



Not really...

Are we happy about this?



Not really...

Hereafter  $\varepsilon = 0.5$  (generalizable?)

## Starting point: Maximum likelihood

$\text{OPT}(\mathbf{X})$

$$\begin{aligned} & \text{maximize} && \langle \boldsymbol{\sigma}, \mathbf{X} \boldsymbol{\sigma} \rangle, \\ & \text{subject to} && \boldsymbol{\sigma} \in \{+1, -1\}^n. \end{aligned}$$

- ▶ Meaningful irrespective of the model
- ▶ Robust to ‘small’ changes in  $\mathbf{X}$
- ▶ NP hard

# Semi-Definite Programming relaxation

OPT( $X$ )

$$\begin{aligned} & \text{maximize} && \langle X, \sigma\sigma^T \rangle, \\ & \text{subject to} && \sigma \in \{+1, -1\}^n. \end{aligned}$$

**Equivalent formulation**

$$\begin{aligned} & \text{maximize} && \langle X, Z \rangle, \\ & \text{subject to} && Z \in \mathbb{R}^{n \times n}, \quad Z \succeq 0, \quad \text{rank}(Z) = 1 \\ & && Z_{ii} = 1 \quad \forall i \in \{1, \dots, n\}. \end{aligned}$$

# Semi-Definite Programming relaxation

OPT( $X$ )

$$\begin{aligned} & \text{maximize} && \langle X, \sigma\sigma^T \rangle, \\ & \text{subject to} && \sigma \in \{+1, -1\}^n. \end{aligned}$$

## Relaxation

$$\begin{aligned} & \text{maximize} && \langle X, Z \rangle, \\ & \text{subject to} && Z \in \mathbb{R}^{n \times n}, \quad Z \succeq 0, \quad \text{rank}(Z) = 1 \\ & && Z_{ii} = 1 \quad \forall i \in \{1, \dots, n\}. \end{aligned}$$

# Semi-Definite Programming relaxation

OPT( $\mathbf{X}$ )

$$\begin{aligned} & \text{maximize} && \langle \mathbf{X}, \boldsymbol{\sigma} \boldsymbol{\sigma}^T \rangle, \\ & \text{subject to} && \boldsymbol{\sigma} \in \{+1, -1\}^n. \end{aligned}$$

SDP( $\mathbf{X}$ )

$$\begin{aligned} & \text{maximize} && \langle \mathbf{X}, \mathbf{Z} \rangle, \\ & \text{subject to} && \mathbf{Z} \in \mathbb{R}^{n \times n}, \quad \mathbf{Z} \succeq \mathbf{0}, \\ & && Z_{ii} = 1 \quad \forall i \in \{1, \dots, n\}. \end{aligned}$$

$$\hat{\boldsymbol{\theta}}^{\text{SDP}}(\mathbf{X}) = c \cdot \mathbf{v}_1(\mathbf{Z}_*), \quad (\text{principal eigenvector})$$

$$\hat{\boldsymbol{\theta}}^{\text{SDP, theor.}}(\mathbf{X}) \sim \mathcal{N}(0, c \mathbf{Z}_*).$$

# Semi-Definite Programming relaxation

OPT( $\mathbf{X}$ )

$$\begin{aligned} & \text{maximize} && \langle \mathbf{X}, \boldsymbol{\sigma} \boldsymbol{\sigma}^T \rangle, \\ & \text{subject to} && \boldsymbol{\sigma} \in \{+1, -1\}^n. \end{aligned}$$

SDP( $\mathbf{X}$ )

$$\begin{aligned} & \text{maximize} && \langle \mathbf{X}, \mathbf{Z} \rangle, \\ & \text{subject to} && \mathbf{Z} \in \mathbb{R}^{n \times n}, \quad \mathbf{Z} \succeq \mathbf{0}, \\ & && Z_{ii} = 1 \quad \forall i \in \{1, \dots, n\}. \end{aligned}$$

$$\hat{\boldsymbol{\theta}}^{\text{SDP}}(\mathbf{X}) = c \cdot \mathbf{v}_1(\mathbf{Z}_*), \quad (\text{principal eigenvector})$$

$$\hat{\boldsymbol{\theta}}^{\text{SDP, theor.}}(\mathbf{X}) \sim \mathcal{N}(0, c \mathbf{Z}_*).$$



# Optimal weak recovery threshold

## Theorem (M., Sen, 2016)

Consider the  $\mathbb{Z}_2$ -synchronization model with  $\theta_0 \in \{+1, -1\}^n$ .  
For any  $\lambda > 1$ , there exist  $\varepsilon(\lambda) > 0$  such that

$$\liminf_{n \rightarrow \infty} \mathbb{E} \left\{ \frac{|\langle \hat{\theta}^{\text{SDP, theor.}}(\mathbf{X}), \theta_0 \rangle|}{\|\hat{\theta}^{\text{SDP, theor.}}(\mathbf{X})\|_2 \|\theta_0\|_2} \right\} \geq \varepsilon(\lambda).$$

[see also: Guédon, Vershynin, 2016; Hajek, Wu, Xu, 2016; Perry, Wein, Moitra, 2016; Bandeira, 2018; ...]

# Optimal weak recovery threshold

## Theorem (M., Sen, 2016)

Consider the  $\mathbb{Z}_2$ -synchronization model with  $\theta_0 \in \{+1, -1\}^n$ .  
For any  $\lambda > 1$ , there exist  $\varepsilon(\lambda) > 0$  such that

$$\liminf_{n \rightarrow \infty} \mathbb{E} \left\{ \frac{|\langle \hat{\theta}^{\text{SDP, theor.}}(\mathbf{X}), \theta_0 \rangle|}{\|\hat{\theta}^{\text{SDP, theor.}}(\mathbf{X})\|_2 \|\theta_0\|_2} \right\} \geq \varepsilon(\lambda).$$

[see also: Guédon, Vershynin, 2016; Hajek, Wu, Xu, 2016; Perry, Wein, Moitra, 2016; Bandeira, 2018; ...]

# Exact limit: Statistical physics prediction

Conjecture (Javanmard, M., Ricci-Tersenghi, 2016)

*In the  $\mathbb{Z}_2$  synchronization model*

$$\lim_{n \rightarrow \infty} \frac{|\langle \hat{\boldsymbol{\theta}}^{\text{SDP}}(\mathbf{X}), \boldsymbol{\theta}_0 \rangle|}{\|\hat{\boldsymbol{\theta}}^{\text{SDP}}(\mathbf{X})\|_2 \|\boldsymbol{\theta}_0\|_2} = \text{ExplicitFormula}(\lambda).$$

## Explicit Formula( $\lambda$ )

For  $G \sim N(0, 1)$ , define  $\rho = \rho(G; \mu, q, r)$  via

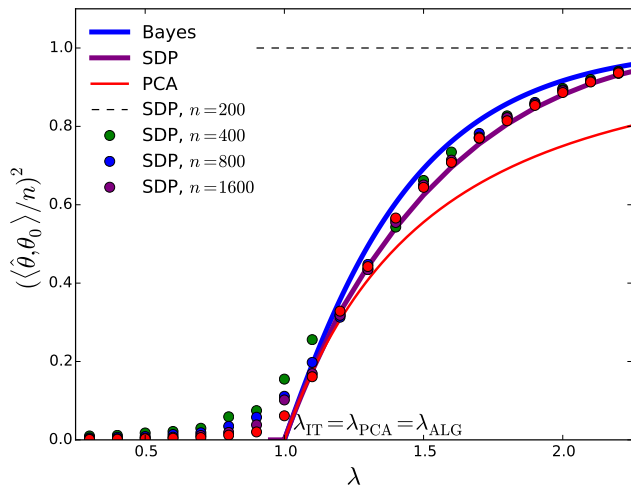
$$1 = \frac{(\mu + \sqrt{q}G)^2}{(\rho + r)^2} + \frac{1 - q}{\rho^2}.$$

Let  $\mu, q, r \in \mathbb{R}$  be the solution of

$$\begin{aligned}\mu &= \lambda \mathbb{E} \left\{ \frac{\mu + \sqrt{q}G}{\rho + r} \right\}, & q &= \mathbb{E} \left\{ \frac{(\mu + \sqrt{q}G)^2}{(\rho + r)^2} \right\}, \\ r &= \mathbb{E} \left\{ \frac{1}{\rho} - \frac{\mu}{\sqrt{q}} \frac{G}{\rho + r} - \frac{G^2}{\rho + r} \right\}.\end{aligned}$$

$$\lim_{n \rightarrow \infty} \frac{|\langle \hat{\theta}^{\text{SDP}}(X), \theta_0 \rangle|}{\|\hat{\theta}^{\text{SDP}}(X)\|_2 \|\theta_0\|_2} = 1 - 2\Phi\left(-\frac{\mu(\lambda)}{\sqrt{q(\lambda)}}\right).$$

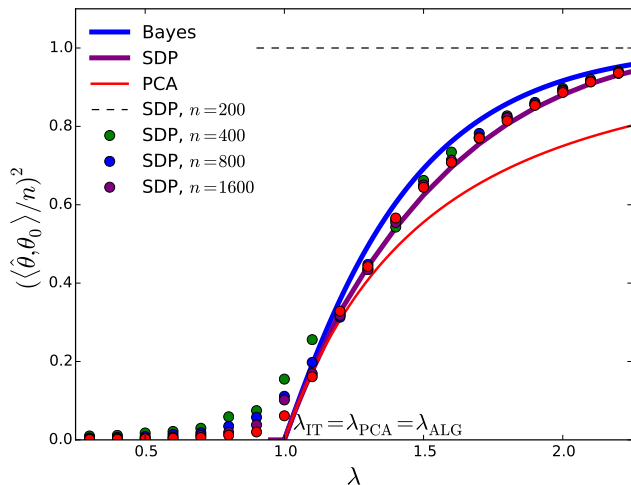
# Bayes vs SDP vs PCA



# SDP

- ▶ Nearly Bayes-optimal
- ▶ Computationally practical
- ▶ Generic/robust
- ▶ Analysis still incomplete

# Can we close the gap?



# TAP free energy

$$\mathcal{F}_{\lambda, \mathbf{X}} : [-1, 1]^n \rightarrow \mathbb{R}$$

$$\mathcal{F}_{\lambda, \mathbf{X}}(\mathbf{m}) = -\frac{1}{n} \sum_{i=1}^n h(m_i) - \frac{\lambda}{2n} \langle \mathbf{m}, \mathbf{X} \mathbf{m} \rangle - \frac{\lambda^2}{4} \left( 1 - \frac{\|\mathbf{m}\|_2^2}{n} \right)^2,$$

[Thouless, Anderson, Palmer, 1977]

## Statistical physics prediction

$$\hat{\boldsymbol{\theta}}^{\text{Bayes}}(\mathbf{X}) = \mathbb{E}\{\boldsymbol{\theta} | \mathbf{X}\} \approx \arg \min_{\mathbf{m} \in [-1, +1]^n} \mathcal{F}_{\lambda, \mathbf{X}}(\mathbf{m})$$

[Technical nuisance: need to break  $\mathbf{m} \leftrightarrow -\mathbf{m}$  symmetry]



# TAP free energy

$$\mathcal{F}_{\lambda, \mathbf{X}}(\mathbf{m}) = -\frac{1}{n} \sum_{i=1}^n h(m_i) - \frac{\lambda}{2n} \langle \mathbf{m}, \mathbf{X} \mathbf{m} \rangle - \frac{\lambda^2}{4} \left( 1 - \frac{\|\mathbf{m}\|_2^2}{n} \right)^2,$$

## Interpretation

$$\mathcal{F}_{\lambda, \mathbf{X}}(\mathbf{m}) \approx \min \left\{ \mathcal{D}_{\text{KL}}(p \| p_{\text{Bayes}}(\cdot | \mathbf{X})) \right. \\ \left. \text{such that } p \in \mathcal{P}(\{+1, -1\}^n), \quad \mathbb{E}_p(\boldsymbol{\sigma}) = \mathbf{m} \right\}$$

[see also: Chatterjee 2010; Talagrand 2010; ...]

# TAP estimator

$$\hat{\boldsymbol{\theta}}^{\text{TAP}}(\mathbf{X}) \equiv \arg \min_{\mathbf{m} \in [-1, +1]^n} \mathcal{F}_{\lambda, \mathbf{X}}(\mathbf{m}).$$

- ▶ Non-convex  $\quad$  :-)
- ▶ Stationary points  $\leftrightarrow$  AMP fixed points  $\quad$  :-)
- ▶ More robust than AMP  $\quad$  :-)
- ▶ Can we prove anything?

# TAP estimator

$$\hat{\boldsymbol{\theta}}^{\text{TAP}}(\mathbf{X}) \equiv \arg \min_{\mathbf{m} \in [-1, +1]^n} \mathcal{F}_{\lambda, \mathbf{X}}(\mathbf{m}).$$

- ▶ Non-convex                    :-)
- ▶ Stationary points  $\leftrightarrow$  AMP fixed points    :-)
- ▶ More robust than AMP        :-)
- ▶ Can we prove anything?

# TAP estimator

$$\hat{\boldsymbol{\theta}}^{\text{TAP}}(\mathbf{X}) \equiv \arg \min_{\mathbf{m} \in [-1, +1]^n} \mathcal{F}_{\lambda, \mathbf{X}}(\mathbf{m}).$$

- ▶ Non-convex                    :-)
- ▶ Stationary points  $\leftrightarrow$  AMP fixed points    :-)
- ▶ More robust than AMP        :-)
- ▶ Can we prove anything?

## What would we like to prove?

$$\hat{\theta}^{\text{TAP}}(\mathbf{X}) \equiv \arg \min_{m \in [-1, +1]^n} \mathcal{F}_{\lambda, \mathbf{X}}(m).$$

- ▶ Landscape is ‘simple’ (can be minimized efficiently)
- ▶ Global minimum close to  $\hat{\theta}^{\text{Bayes}}(\mathbf{X})$

## A recent result

### Theorem (Fan, Mei, M, 2018)

Denote  $\mathcal{C}_{\lambda,n}$  be the set of critical points of  $\mathcal{F}_{\lambda,\mathbf{X}}$  with  $\mathcal{F}_{\lambda,\mathbf{X}}(\mathbf{m}) \leq -\lambda^2/3$ . There exists  $\lambda_0 > 0$ , such that for any  $\lambda > \lambda_0$ , we have

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[ \sup_{\mathbf{m} \in \mathcal{C}_{\lambda,n}} \frac{1}{n^2} \|\mathbf{m}\mathbf{m}^\top - \mathbb{E}\{\boldsymbol{\theta}\boldsymbol{\theta}^\top | \mathbf{X}\}\|_F^2 \right] = 0.$$

## Conclusion

## Conclusion: A general strategy

- ▶ High-dimensional Bayes posterior  $\rightarrow$  Fundamental statistical limits
- ▶ AMP  $\rightarrow$  Computational limits
- ▶ Convex relaxations  $\rightarrow$  Robust/generic algorithms
- ▶ Free energy approximations  $\rightarrow$  Better cost functions

Thanks!



## Conclusion: A general strategy

- ▶ High-dimensional Bayes posterior  $\rightarrow$  Fundamental statistical limits
- ▶ AMP  $\rightarrow$  Computational limits
- ▶ Convex relaxations  $\rightarrow$  Robust/generic algorithms
- ▶ Free energy approximations  $\rightarrow$  Better cost functions

Thanks!