

# TAP equations for invariant coupling matrices

Manfred Opper (TUB),  
joint work with

B Cakmak (Aalborg), L Bachschmid-Romano (TUB), O Winther (DTU)

February 27, 2017



- Cavity method for dense Ising networks
- TAP equations for random coupling matrices
- Dynamics of solving TAP equations
- Learning curves for Inverse Ising model

# Simplest model: Ising

$$\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_N) \in \{\pm 1\}^N$$

$$P(\boldsymbol{\sigma}) = \frac{1}{Z} \exp \left[ \sum_{i < j}^N J_{ij} \sigma_i \sigma_j + \sum_{i=1}^N H_i \sigma_i \right]$$

# Simplest model: Ising

$$\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_N) \in \{\pm 1\}^N$$

$$P(\boldsymbol{\sigma}) = \frac{1}{Z} \exp \left[ \sum_{i < j}^N J_{ij} \sigma_i \sigma_j + \sum_{i=1}^N H_i \sigma_i \right]$$

Try to compute 'exact mean field equations' marginals for  $m_i \doteq \langle S_i \rangle$  when  $\mathbf{J}$  large random matrix  $\rightarrow$  TAP equations (Thouless, Anderson and Palmer 78).

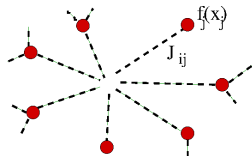
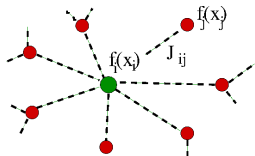
# Cavity method

- The marginal at node  $i$  can be derived from the joint distribution

$$p_i(\sigma, h) \propto e^{\sigma(h+H_i)} p_{\setminus i}(h)$$

with the 'cavity field' distribution

$$p_{\setminus i}(h) = \sum_{\sigma_{\setminus i}} \delta \left( h - \sum_{j \neq i} J_{ij} \sigma_j \right) p_{\setminus i}(\sigma_{\setminus i})$$



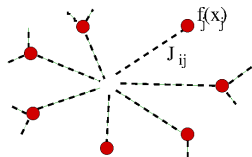
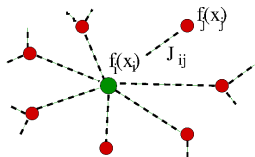
# Cavity method

- The marginal at node  $i$  can be derived from the joint distribution

$$p_i(\sigma, h) \propto e^{\sigma(h+H_i)} p_{\setminus i}(h)$$

with the 'cavity field' distribution

$$p_{\setminus i}(h) = \sum_{\sigma_{\setminus i}} \delta \left( h - \sum_{j \neq i} J_{ij} \sigma_j \right) p_{\setminus i}(\sigma_{\setminus i})$$



- Approximate  $p_{\setminus i}(h)$  by Gaussian (central limit theorem)  
 $p_{\setminus i}(h) = \mathcal{N}(a_i, V_i)$ .

- Magnetisation  $m_i = \langle \sigma_i \rangle$  given by

$$m_i = \tanh \left( H_i + \sum_j J_{ij} m_j - V_i m_i \right)$$

- Susceptibility matrix

$$\chi \doteq \frac{\partial \mathbf{m}}{\partial \mathbf{H}} = (\mathbf{\Lambda} - \mathbf{J})^{-1}$$

where  $\lambda_i = V_i + \frac{1}{\chi_{ii}}$ .

- $V_i$  determined by consistency  $\chi_{ii} = 1 - m_i^2$

# Selfaveraging cavity variance

- Random matrix assumption  $\mathbf{J} \sim \mathbf{O}^\top \mathbf{J} \mathbf{O}$  (Parisi and Potters 95)



# Selfaveraging cavity variance

- Random matrix assumption  $\mathbf{J} \sim \mathbf{O}^\top \mathbf{J} \mathbf{O}$  (Parisi and Potters 95)
- If  $V \equiv R_{\mathbf{J}}(\bar{\chi})$  with  $\bar{\chi} = \frac{1}{N} \sum_i \chi_{ii}$  one can show

$$\left[ (\mathbf{\Lambda} - \mathbf{J})^{-1} \right]_{ii} \rightarrow \chi_{ii}$$

as  $N \rightarrow \infty$ .

# Selfaveraging cavity variance

- Random matrix assumption  $\mathbf{J} \sim \mathbf{O}^\top \mathbf{J} \mathbf{O}$  (Parisi and Potters 95)
- If  $V \equiv R_{\mathbf{J}}(\bar{\chi})$  with  $\bar{\chi} = \frac{1}{N} \sum_i \chi_{ii}$  one can show

$$\left[ (\mathbf{\Lambda} - \mathbf{J})^{-1} \right]_{ii} \rightarrow \chi_{ii}$$

as  $N \rightarrow \infty$ .

- The R-transform is defined as

$$R_{\mathbf{J}}(s) \triangleq S_{\mathbf{J}}^{-1}(s) - 1/s$$

where

$$S_{\mathbf{J}}(z) \triangleq \frac{1}{N} \text{Tr}(z\mathbf{I} - \mathbf{J})^{-1}$$

# Algorithms for solving TAP equations ?

- Analysis of dynamics of random systems usually hard !

# Algorithms for solving TAP equations ?

- Analysis of dynamics of random systems usually hard !
- Exact dynamics converging to TAP For SK model by Bolthausen (2014).
- Other results by Kabashima (2003) in the context of the CDMA and by Donoho, Maleki, Montanari (2009) for compressed sensing.

# Generating functional approach

- Candidate algorithm could be of the form

$$\mathbf{m}(t) = \tanh \left( \{ \mathbf{h}(\tau), \mathbf{m}(\tau) \}_{\tau=0}^{t-1} \right)$$

$$\mathbf{h}(t) = \mathbf{H} + \mathbf{J}\mathbf{m}(t)$$

# Generating functional approach

- Candidate algorithm could be of the form

$$\begin{aligned}\mathbf{m}(t) &= \tanh(\{\mathbf{h}(\tau), \mathbf{m}(\tau)\}_{\tau=0}^{t-1}) \\ \mathbf{h}(t) &= \mathbf{H} + \mathbf{J}\mathbf{m}(t)\end{aligned}$$

- Try average case analysis: use generating functional  $\langle Z(\{\mathbf{l}(t)\}) \rangle_{\mathbf{J}}$  where

$$\begin{aligned}Z(\{\mathbf{l}(t)\}) &= \int \prod_{t=0}^{T-1} \left\{ d\mathbf{m}(t) d\mathbf{h}(t) \delta(\mathbf{m}(t) - \tanh(\{\mathbf{h}(\tau), \mathbf{m}(\tau)\}_{\tau=0}^{t-1})) \right. \\ &\quad \left. \delta(\mathbf{h}(t) - \mathbf{H} - \mathbf{J}\mathbf{m}(t)) e^{i\mathbf{h}(t)^\top \mathbf{l}(t)} \right\}.\end{aligned}$$

$$\mathbf{m}(t) = \tanh \left( \{\mathbf{h}(\tau), \mathbf{m}(\tau)\}_{\tau=0}^{t-1} \right)$$
$$\mathbf{h}(t) = \mathbf{H} + \sum_{s=0}^{t-1} \hat{\mathcal{G}}(t, s) \mathbf{m}(s) + \phi(t) .$$

with  $\phi(t)$  discrete time Gaussian process and order parameters

$$\mathbf{m}(t) = \tanh \left( \{\mathbf{h}(\tau), \mathbf{m}(\tau)\}_{\tau=0}^{t-1} \right)$$
$$\mathbf{h}(t) = \mathbf{H} + \sum_{s=0}^{t-1} \hat{\mathcal{G}}(t, s) \mathbf{m}(s) + \phi(t) .$$

with  $\phi(t)$  discrete time Gaussian process and order parameters

$$\hat{\mathcal{G}} = \mathbf{R}_{\mathbf{J}}(\mathcal{G}) \quad \mathcal{G}(t, \tau) = \left\langle \frac{\partial m_i(t)}{\partial \phi_i(\tau)} \right\rangle_{\phi}$$

and something more complicated for  $C_{\phi}$ .



$$\mathbf{m}(t) = \tanh \left( \{\mathbf{h}(\tau), \mathbf{m}(\tau)\}_{\tau=0}^{t-1} \right)$$
$$\mathbf{h}(t) = \mathbf{H} + \sum_{s=0}^{t-1} \hat{\mathcal{G}}(t, s) \mathbf{m}(s) + \phi(t) .$$

with  $\phi(t)$  discrete time Gaussian process and order parameters

$$\hat{\mathcal{G}} = \mathbf{R}_{\mathbf{J}}(\mathcal{G}) \quad \mathcal{G}(t, \tau) = \left\langle \frac{\partial m_i(t)}{\partial \phi_i(\tau)} \right\rangle_{\phi}$$

and something more complicated for  $C_{\phi}$ .

**Memory could be bad for convergence !**

# Single step memory algorithm

- **Idea:** Try to **erase the memory** terms

$$\hat{G}(t, \tau) = 0 \quad \forall \tau \neq t - 1$$

# Single step memory algorithm

- **Idea:** Try to **erase the memory** terms

$$\hat{G}(t, \tau) = 0 \quad \forall \tau \neq t - 1$$

- and subtract:  $\mathbf{Jm}(\tau) - \hat{G}(\tau, \tau - 1)\mathbf{m}(\tau - 1)$ .

# Single step memory algorithm

- **Idea:** Try to **erase the memory** terms

$$\hat{G}(t, \tau) = 0 \quad \forall \tau \neq t - 1$$

- and subtract:  $\mathbf{Jm}(\tau) - \hat{G}(\tau, \tau - 1)\mathbf{m}(\tau - 1)$ . Consistency with definition of  $G(t, \tau)$  and TAP asymptotics suggests

$$\mathbf{m}(t + 1) = \tanh(\psi(t))$$

$$\psi(t) = Q(t) \sum_{\tau=0}^t a_{t+1-\tau} \frac{\mathbf{H} + \mathbf{Jm}(\tau) - \hat{G}(\tau, \tau - 1)\mathbf{m}(\tau - 1)}{Q(\tau - 1)(1 - q(\tau))}$$

where we define  $\hat{G}(t, t - 1) = \frac{1 - q(t)}{1 - q(t - 1)} \mathbf{R}_{\mathbf{J}}(1 - q(t - 1))$

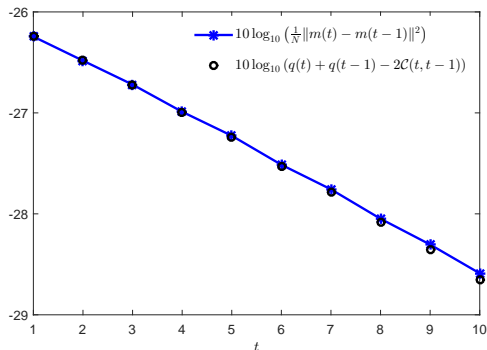
$Q(t) = \prod_{\tau=0}^t \mathbf{R}_{\mathbf{J}}(1 - q(\tau))$  and the coefficients  $a_k$  via

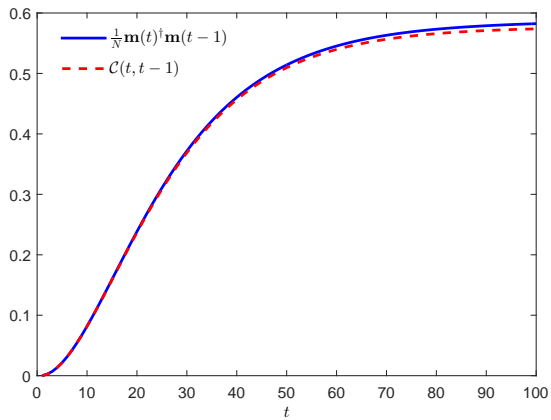
$$\mathbf{R}^{-1}(x) = \sum_{n=1}^{\infty} a_n x^n.$$

- For  $\mathbf{J}$  Gaussian i.i.d. (Sherrington Kirkpatrick model) agrees with Bolthausen's (2014) result
- For  $-\mathbf{J} \sim$  Wishart coincides with AMP algorithm introduced by Kabashima (2003) in the context of the CDMA and by Donoho, Maleki, Montanari (2009) for compressed sensing.

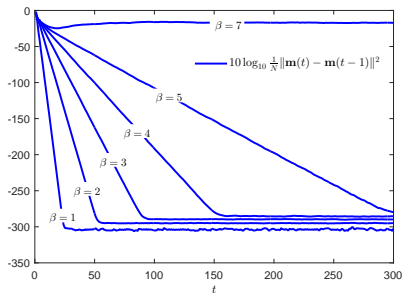
# Random orthogonal ensemble: Analytical results vs Simulation

$R^{-1}(x) = x/(\beta^2 - x^2)$  and  $a_n = \frac{1}{\beta^{n+1}}$  for  $n$  odd and  $a_n = 0$  else.  
 $N = 2^{14}$ ,  $\beta = 20$  and  $h_i = 1$ , single realisation of  $\mathbf{J}$ .



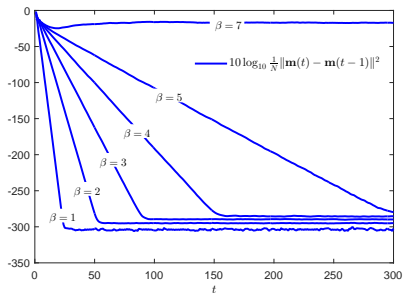


# Convergence for $N = 2^{14}$ , $H_i = 2$ (Simulations)

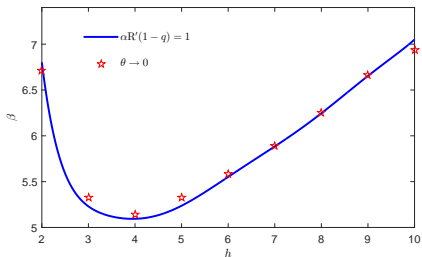




## Convergence for $N = 2^{14}$ , $H_i = 2$ (Simulations)

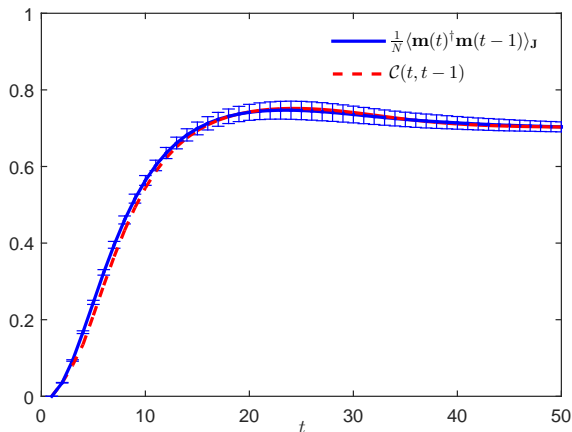


## Stability of fixed point (Almeida–Thouless line)



# Analytical results vs Simulation

Random orthogonal ensemble, region of instability:  $N = 2^{12}$ ,  $\beta = 10$  and  $H_i = 2$ ,  $5 \times 10^3$  Realisations of  $\mathbf{J}$



# Inverse Ising problem

- Learn couplings from Ising data  $\sigma^1, \dots, \sigma^M$  generated from

$$P(\boldsymbol{\sigma}|\mathbf{J}) = Z_{\text{Ising}}^{-1} \exp \left[ \beta \sum_{i < j} J_{ij} \sigma_i \sigma_j \right]$$

# Inverse Ising problem

- Learn couplings from Ising data  $\sigma^1, \dots, \sigma^M$  generated from

$$P(\sigma|\mathbf{J}) = Z_{\text{Ising}}^{-1} \exp \left[ \beta \sum_{i < j} J_{ij} \sigma_i \sigma_j \right]$$

- Local approach:  $\mathbf{W} = (W_1, \dots, W_N) \doteq \sqrt{N}(J_{01}, \dots, J_{0N-1})$  for **spin** 0 using cost function

$$E(\mathbf{W}) = \sum_{k=1}^M \mathcal{E}(\mathbf{W}; \sigma^k)$$

where  $\mathcal{E}(\mathbf{W}; \sigma) = \Phi(\sigma_0, h)$  and  $h \doteq \frac{1}{\sqrt{N}} \sum_{j \neq 0} W_j \sigma_j$

# Inverse Ising problem

- Learn couplings from Ising data  $\sigma^1, \dots, \sigma^M$  generated from

$$P(\sigma|\mathbf{J}) = Z_{\text{Ising}}^{-1} \exp \left[ \beta \sum_{i<j} J_{ij} \sigma_i \sigma_j \right]$$

- Local approach:  $\mathbf{W} = (W_1, \dots, W_N) \doteq \sqrt{N}(J_{01}, \dots, J_{0N-1})$  for **spin** 0 using cost function

$$E(\mathbf{W}) = \sum_{k=1}^M \mathcal{E}(\mathbf{W}; \sigma^k)$$

where  $\mathcal{E}(\mathbf{W}; \sigma) = \Phi(\sigma_0, h)$  and  $h \doteq \frac{1}{\sqrt{N}} \sum_{j \neq 0} W_j \sigma_j$

- Example: **pseudo-likelihood**

$$\mathcal{E}(\mathbf{W}; \sigma) = -\ln P(\sigma_0 | \sigma_{\setminus 0}, \mathbf{W})$$

# Average case learning performance

- Define partition function of student couplings

$$Z = \int d\mathbf{W} \exp[-\nu E(\mathbf{W})],$$

- Order parameters from  $\nu \rightarrow \infty$  limit of quenched free energy

$$F = -N^{-1} \nu^{-1} \overline{\ln Z} = - \lim_{n \rightarrow 0} N^{-1} \nu^{-1} \frac{\partial}{\partial n} \ln \overline{Z^n},$$

# Average case learning performance

- Define partition function of student couplings

$$Z = \int d\mathbf{W} \exp[-\nu E(\mathbf{W})],$$

- Order parameters from  $\nu \rightarrow \infty$  limit of quenched free energy

$$F = -N^{-1} \nu^{-1} \overline{\ln Z} = - \lim_{n \rightarrow 0} N^{-1} \nu^{-1} \frac{\partial}{\partial n} \ln \overline{Z^n},$$

- Quenched averages over the fields performed using **cavity method**:

$$P(\sigma_0, h_*, h_1, \dots, h_n) = \frac{1}{Z_0} e^{\beta \sigma_0 h_*} p_{\text{cav}}(h_*, h_1, \dots, h_n)$$

where  $p_{\text{cav}}(h_*, h_1, \dots, h_n)$  is a **multivariate Gaussian** with covariance

$$\langle h_a h_b \rangle = \frac{1}{N} \sum_{i, j \neq 0} W_j^a C_{ij}^{\setminus 0} W_j^b \quad C_{ij}^{\setminus 0} = \langle \sigma_i \sigma_j \rangle_{\setminus 0}$$

- Reconstruction error

$$\varepsilon = N^{-1} \overline{(\mathbf{W}^* - \mathbf{W})^2} = \left(q - \frac{R^2}{V}\right) T + Y \left(1 - \frac{R}{V}\right)^2.$$

depends on the teacher statistics

$$Y \doteq N^{-1} (\mathbf{W}^*)^2$$

$$V \doteq \frac{1}{N} \sum_{i,j \neq 0} W_j^* C_{ij}^{-1} W_j^*$$

$$T \doteq \frac{1}{N} \text{Tr} \mathbf{C}^{-1} = \beta^2 V + 1$$

- $q$  and  $R$  are derived from the free energy ( $\alpha = M/N$ ):



- Reconstruction error

$$\varepsilon = N^{-1} \overline{(\mathbf{W}^* - \mathbf{W})^2} = \left(q - \frac{R^2}{V}\right) T + Y \left(1 - \frac{R}{V}\right)^2.$$

depends on the teacher statistics

$$Y \doteq N^{-1} (\mathbf{W}^*)^2$$

$$V \doteq \frac{1}{N} \sum_{i,j \neq 0} W_j^* C_{ij}^{\setminus 0} W_j^*$$

$$T \doteq \frac{1}{N} \text{Tr} \mathbf{C}^{-1} = \beta^2 V + 1$$

- $q$  and  $R$  are derived from the free energy ( $\alpha = M/N$ ):

$$F = - \text{extr}_{q,R,x} \max_z \left\{ \frac{1}{2} \frac{q - R^2/V}{x} + \frac{1}{2} \alpha e^{-\frac{\beta^2 R^2}{2q}} \sum_{\sigma_0} \int \mathcal{D}u e^{-\frac{\beta \sigma_0 R}{\sqrt{q}} u} \left[ -\frac{z^2}{2} - \Phi(\sigma_0, \sqrt{x}z + \sqrt{q}u) \right] \right\}.$$

# Maximum Likelihood (mean field)

- Definition

$$J_{ij}^{ML-MF} = -\frac{1}{\beta}(\hat{\mathbf{C}}^{-1})_{ij}$$

# Maximum Likelihood (mean field)

- Definition

$$J_{ij}^{ML-MF} = -\frac{1}{\beta}(\hat{\mathbf{C}}^{-1})_{ij}$$

- Using matrix inversion lemma

$$W_i^{ML-MF} = \sqrt{N}J_{0i}^{ML-MF} = \frac{\sqrt{N}}{\beta}\phi_0 \sum_{j \neq 0} (\hat{\mathbf{C}}^{-1})_{ji} \hat{\mathbf{C}}_{0j}$$

where  $\phi_0 = \frac{1}{1 - \sum_{i,j \neq 0} \hat{\mathbf{C}}_{0i} (\hat{\mathbf{C}}^{-1})_{ij} \hat{\mathbf{C}}_{0j}}$  is selfaveraging.

# Maximum Likelihood (mean field)

- Definition

$$J_{ij}^{ML-MF} = -\frac{1}{\beta}(\hat{\mathbf{C}}^{-1})_{ij}$$

- Using matrix inversion lemma

$$W_i^{ML-MF} = \sqrt{N}J_{0i}^{ML-MF} = \frac{\sqrt{N}}{\beta}\phi_0 \sum_{j \neq 0} (\hat{\mathbf{C}}^{-1})_{ji} \hat{\mathbf{C}}_{0j}$$

where  $\phi_0 = \frac{1}{1 - \sum_{i,j \neq 0} \hat{\mathbf{C}}_{0i} (\hat{\mathbf{C}}^{-1})_{ij} \hat{\mathbf{C}}_{0j}}$  is selfaveraging.

- $W_i^{ML-MF}$  minimises local cost function

$$E(\mathbf{W}) = \frac{\alpha\beta}{2\phi_0} \sum_{i,j \neq 0} W_i \hat{\mathbf{C}}_{ij} W_j - \alpha\sqrt{N} \sum_{j \neq 0} \hat{\mathbf{C}}_{0j} W_j,$$

# Maximum Likelihood (mean field)

- Definition

$$J_{ij}^{ML-MF} = -\frac{1}{\beta} (\hat{\mathbf{C}}^{-1})_{ij}$$

- Using matrix inversion lemma

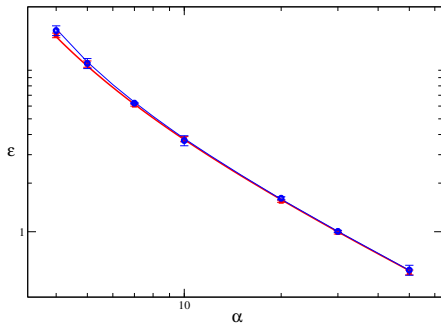
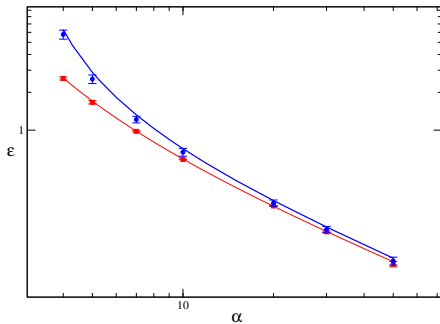
$$W_i^{ML-MF} = \sqrt{N} J_{0i}^{ML-MF} = \frac{\sqrt{N}}{\beta} \phi_0 \sum_{j \neq 0} (\hat{\mathbf{C}}^{-1})_{ji} \hat{\mathbf{C}}_{0j}$$

where  $\phi_0 = \frac{1}{1 - \sum_{i,j \neq 0} \hat{\mathbf{C}}_{0i} (\hat{\mathbf{C}}^{-1})_{ij} \hat{\mathbf{C}}_{0j}}$  is selfaveraging.

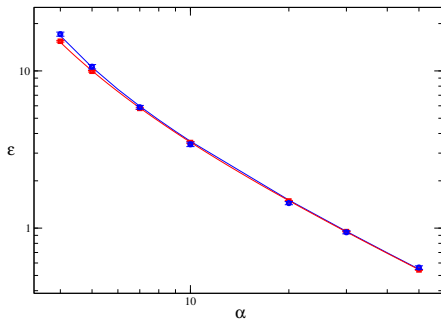
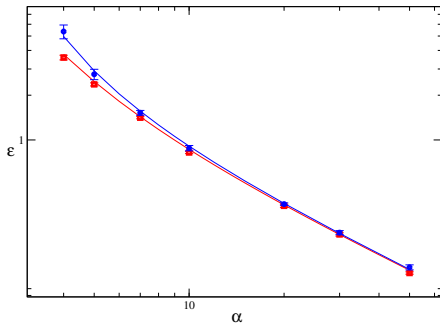
- $W_i^{ML-MF}$  minimises local cost function

$$E(\mathbf{W}) = \frac{\alpha\beta}{2\phi_0} \sum_{i,j \neq 0} W_i \hat{\mathbf{C}}_{ij} W_j - \alpha\sqrt{N} \sum_{j \neq 0} \hat{\mathbf{C}}_{0j} W_j,$$

$$E(\mathbf{W}^{ML-MF}) = -\frac{\alpha N \phi_0}{2\beta} \sum_{i,j \neq 0} \hat{\mathbf{C}}_{0i} (\hat{\mathbf{C}}^{-1})_{ij} \hat{\mathbf{C}}_{0j} = \frac{\alpha N (1 - \phi_0)}{2\beta}.$$

SK:  $\beta=0.2$ SK:  $\beta=0.8$ 

SK-model:  $J_{ij}^* \sim \mathcal{N}(0, 1/N)$

Wishart:  $\rho=0.25, \beta=0.2$ Wishart:  $\rho=0.5, \beta=0.5$ 

Hopfield–Wishart  $J_{ij}^* = \frac{1}{N} \sum_{\mu=1}^{\rho N} \xi_i^{\mu} \xi_j^{\mu}$

$$\varepsilon \simeq \frac{c}{\alpha} \quad \alpha \rightarrow \infty$$

where

$$c_{\text{PLM}} = \frac{1}{\beta^2} \frac{1}{\langle 1 - \tanh^2(\beta h_*) \rangle_{h_*}} \frac{1}{N} \text{Tr} \mathbf{C}^{-1}$$
$$c_{\text{MF-ML}} = \frac{1 + \beta^2 V}{\beta^2} \frac{1}{N} \text{Tr} \mathbf{C}^{-1}$$



- Try to estimate  $R$  transforms for real data.
- Optimise estimators for Inverse Ising problem.



