

# Approximated Newton Algorithm for the Ising Model Inference Speeds Up Convergence, Performs Optimally and Avoids Over-fitting

Ulisse Ferrari

Institut de la Vision, Sorbonne Universités, UPMC

New Frontiers in Non-equilibrium Physics 2015

## Outlook of the seminar

- 1 Introduction with an application of pairwise Ising Model to Neuroscience
- 2 Maximal Entropy model and the Vanilla (Standard) Learning Algorithm
- 3 Approximate Newton Method
- 4 The Long-Time Limit: Stochastic Dynamics
- 5 Properties of the Stationary Distribution
- 6 Conclusions and Perspectives

# Model Inference:

Finding the probability distribution reproducing  
the data system statistics.

## Model Inference:

Finding the probability distribution reproducing  
the data system statistics.

Useful for characterizing the behavior of  
systems of many, strongly correlated, units:  
*neurons, proteins, virus, species distribution, bird flocks*  
but...

## Model Inference:

Finding the probability distribution reproducing  
the data system statistics.

Useful for characterizing the behavior of  
systems of many, strongly correlated, units:  
*neurons, proteins, virus, species distribution, bird flocks*  
but...

which distribution?

## Model Inference:

Finding the probability distribution reproducing  
the data system statistics.

Useful for characterizing the behavior of  
systems of many, strongly correlated, units:  
*neurons, proteins, virus, species distribution, bird flocks*  
but...

which distribution?

## Maximum Entropy (MaxEnt) Inference:

Search for the largest entropy distribution satisfying a set of  
constraints.

## Example: pairwise Ising Model

Given binary units data-set of  $B$  configurations of  $N$  units:

$$\left\{ \left\{ \sigma_i(b) \right\}_{i=1}^N \right\}_{b=1}^B$$

Find the MaxEnt model reproducing single and pairwise correlations:

$$\langle \sigma_i \rangle_{\text{MODEL}} = \langle \sigma_i \rangle_{\text{DATA}} \equiv \frac{1}{B} \sum_b \sigma_i(b)$$

$$\langle \sigma_i \sigma_j \rangle_{\text{MODEL}} = \langle \sigma_i \sigma_j \rangle_{\text{DATA}} \equiv \frac{1}{B} \sum_b \sigma_i(b) \sigma_j(b)$$

## Example: pairwise Ising Model

Given binary units data-set of  $B$  configurations of  $N$  units:

$$\left\{ \left\{ \sigma_i(b) \right\}_{i=1}^N \right\}_{b=1}^B$$

Find the MaxEnt model reproducing single and pairwise correlations:

$$\langle \sigma_i \rangle_{\text{MODEL}} = \langle \sigma_i \rangle_{\text{DATA}} \equiv \frac{1}{B} \sum_b \sigma_i(b)$$

$$\langle \sigma_i \sigma_j \rangle_{\text{MODEL}} = \langle \sigma_i \sigma_j \rangle_{\text{DATA}} \equiv \frac{1}{B} \sum_b \sigma_i(b) \sigma_j(b)$$

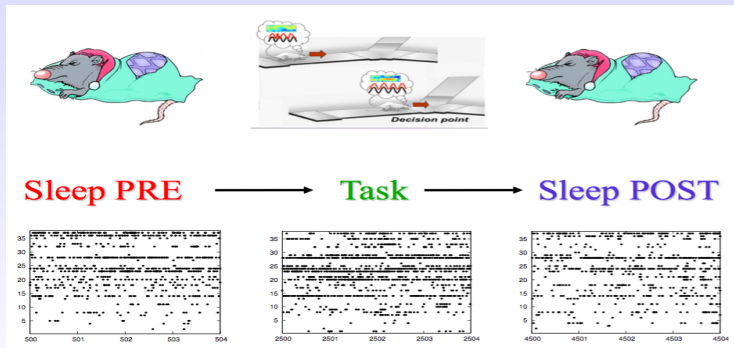
Finely tune the parameters  $\{h, J\}$  of the  
pairwise Ising model:

$$P_{h,j}(\sigma) = \exp \left\{ \sum_i h_i \sigma_i + \sum_{ij} J_{ij} \sigma_i \sigma_j \right\} / Z[h, J]$$



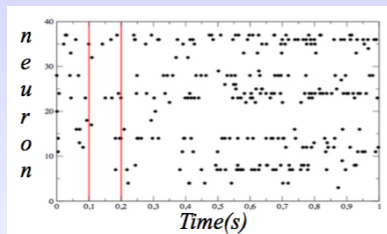
## *In vivo* Pre-Frontal Cortex Recording:

## In vivo Pre-Frontal Cortex Recording: 97 experimental sessions of:



*Peyrache et al. Nat. Neurosci. (2009)*

## Ising Model Inference

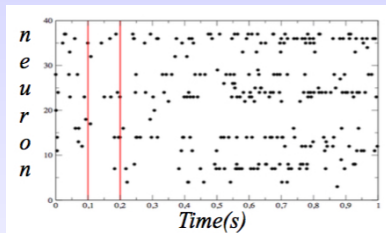


$\sigma_i(b) = 1$  if neuron  $i$  spiked during time-bin  $b$

Ask to reproduce neurons firing rates and correlations.

*Schneidman et al. Nature 2006; Cocco, Monasson, PRL (2011)*

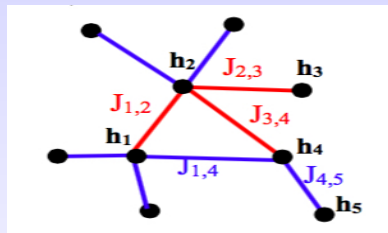
## Ising Model Inference



⇒

⇒

⇒

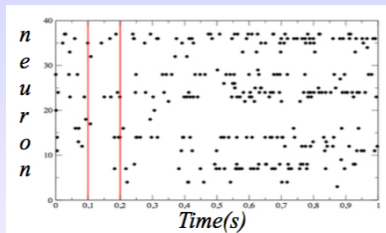


$\sigma_i(b) = 1$  if neuron  $i$  spiked during time-bin  $b$

Ask to reproduce neurons firing rates and correlations.

Schneidman et al. Nature 2006; Cocco, Monasson, PRL (2011)

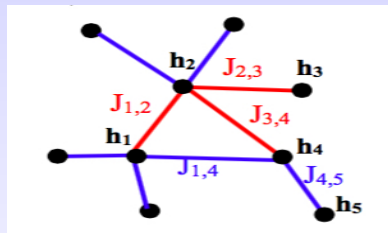
## Ising Model Inference



⇒

⇒

⇒



$\sigma_i(b) = 1$  if neuron  $i$  spiked during time-bin  $b$

Ask to reproduce neurons firing rates and correlations.

$97 \times 3$  couplings network sets ( $97 \times \{\text{PRE, TASK, POST}\}$ )

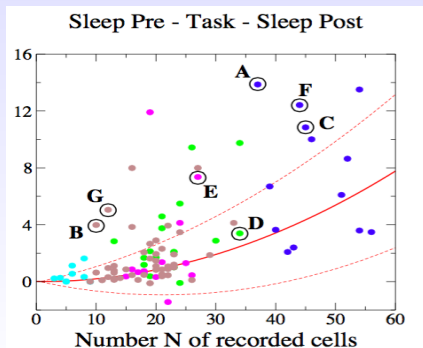
*Schneidman et al. Nature 2006; Cocco, Monasson, PRL (2011)*

## Learning related coupling Adjustment

$$A = \sum_{i,j: J_{ij}^{\text{TASK}}, J_{ij}^{\text{POST}} \neq 0} \text{sign}(J_{ij}^{\text{TASK}} - J_{ij}^{\text{PRE}}) \cdot (J_{ij}^{\text{POST}} - J_{ij}^{\text{PRE}})$$

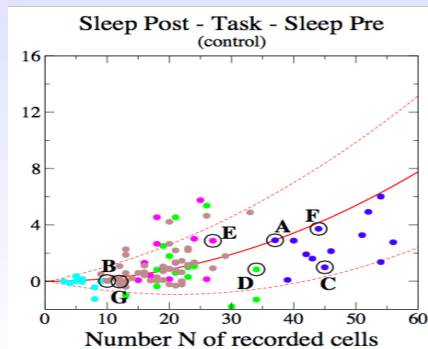
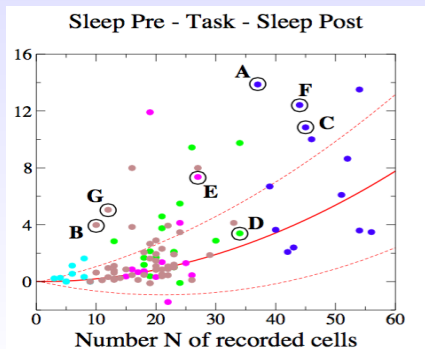
## Learning related coupling Adjustment

$$A = \sum_{i,j: J_{ij}^{\text{TASK}}, J_{ij}^{\text{POST}} \neq 0} \text{sign}(J_{ij}^{\text{TASK}} - J_{ij}^{\text{PRE}}) \cdot (J_{ij}^{\text{POST}} - J_{ij}^{\text{PRE}})$$



## Learning related coupling Adjustment

$$A = \sum_{i,j: J_{ij}^{\text{TASK}}, J_{ij}^{\text{POST}} \neq 0} \text{sign}(J_{ij}^{\text{TASK}} - J_{ij}^{\text{PRE}}) \cdot (J_{ij}^{\text{POST}} - J_{ij}^{\text{PRE}})$$





- 1 Maximal Entropy Models and the Vanilla (standard) Learning Algorithm
- 2 Approximated Newton Method
- 3 The Long-Time Limit: Stochastic Dynamics
- 4 Properties of the Stationary Distribution

## General MaxEnt

Given a list of  $D$  observables to reproduce  $\{\Sigma_a(\sigma)\}_{a=1}^D$   
(generic functions of the system units)

Find the MaxEnt model parameters  $\{X_a\}_{a=1}^D$

$$P_{\mathbf{X}}(\sigma) = \exp \left\{ \sum_a X_a \Sigma_a(\sigma) \right\} / Z[\mathbf{X}]$$

reproducing the observables averages:

$$\langle \Sigma_a \rangle_{\text{DATA}} \equiv P_a = Q_a[\mathbf{X}] \equiv \langle \Sigma_a \rangle_{\mathbf{X}}$$

Equivalent to log-likelihood maximization:

$$\mathbf{X}^* = \arg \max_{\mathbf{X}} [\log L[\mathbf{X}]] \equiv \arg \max_{\mathbf{X}} [\mathbf{X} \cdot \mathbf{P} - \log Z[\mathbf{X}]]$$

Equivalent to log-likelihood maximization:

$$\mathbf{X}^* = \arg \max_{\mathbf{X}} [\log L[\mathbf{X}]] \equiv \arg \max_{\mathbf{X}} [\mathbf{X} \cdot \mathbf{P} - \log Z[\mathbf{X}]]$$

in fact:

$$\nabla_a \log L[\mathbf{X}] = \frac{d}{dX_a} [\mathbf{X} \cdot \mathbf{P} - \log Z[\mathbf{X}]] = P_a - Q_a[\mathbf{X}]$$

Equivalent to log-likelihood maximization:

$$\mathbf{X}^* = \arg \max_{\mathbf{X}} [\log L[\mathbf{X}]] \equiv \arg \max_{\mathbf{X}} [\mathbf{X} \cdot \mathbf{P} - \log Z[\mathbf{X}]]$$

in fact:

$$\nabla_a \log L[\mathbf{X}] = \frac{d}{dX_a} [\mathbf{X} \cdot \mathbf{P} - \log Z[\mathbf{X}]] = P_a - Q_a[\mathbf{X}]$$

Cannot be solved analytically. Ackley, Hinton and Sejnowski  
(Vanilla Gradient):

$$\mathbf{X}_{t+1} = \mathbf{X}_t + \delta \mathbf{X}_t^{\text{VG}}; \quad \delta \mathbf{X}_t^{\text{VG}} = \alpha (\mathbf{P} - \mathbf{Q}[\mathbf{X}_t])$$

Equivalent to log-likelihood maximization:

$$\mathbf{X}^* = \arg \max_{\mathbf{X}} [\log L[\mathbf{X}]] \equiv \arg \max_{\mathbf{X}} [\mathbf{X} \cdot \mathbf{P} - \log Z[\mathbf{X}]]$$

in fact:

$$\nabla_a \log L[\mathbf{X}] = \frac{d}{dX_a} [\mathbf{X} \cdot \mathbf{P} - \log Z[\mathbf{X}]] = P_a - Q_a[\mathbf{X}]$$

Cannot be solved analytically. Ackley, Hinton and Sejnowski  
(Vanilla Gradient):

$$\mathbf{X}_{t+1} = \mathbf{X}_t + \delta \mathbf{X}_t^{\text{VG}}; \quad \delta \mathbf{X}_t^{\text{VG}} = \alpha (\mathbf{P} - \mathbf{Q}[\mathbf{X}_t])$$

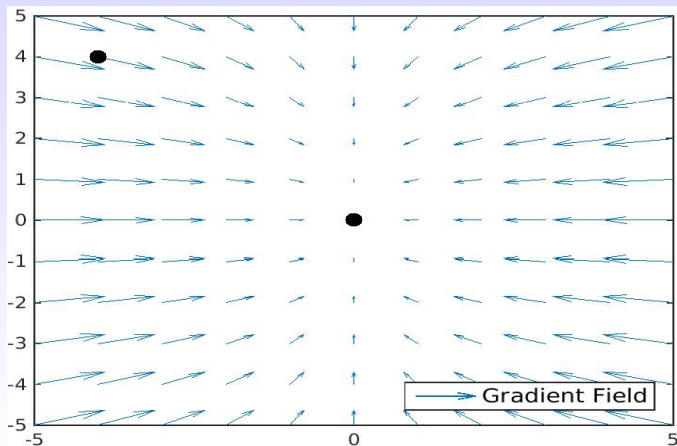
If  $0 < P_a < 1$  for all  $a = 1, \dots, D$ , the problem is well posed:

$\mathbf{X}^*$  exists and is unique and the dynamics converges

(for infinitesimally small  $\alpha$ )

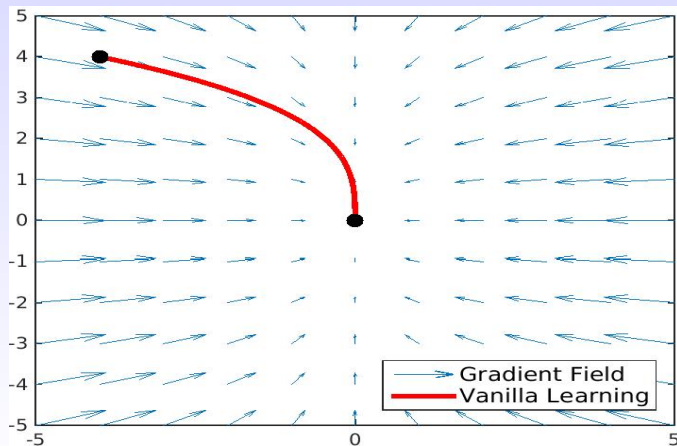
## A 2-dimensional example:

$$\log L[u, v] = -\frac{a}{2}(u - u_\infty)^2 - \frac{b}{2}(v - v_\infty)^2$$



## A 2-dimensional example:

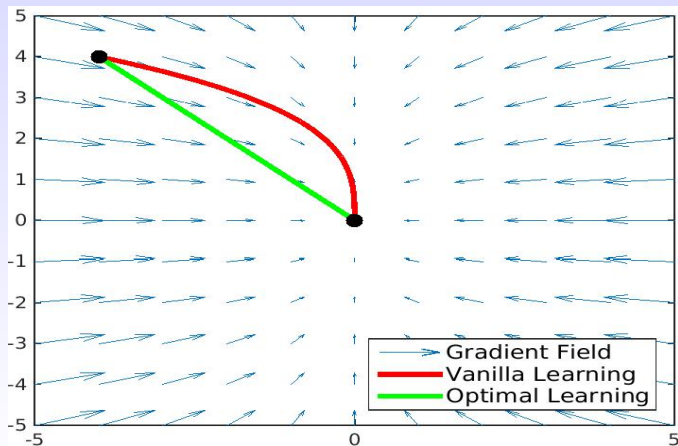
$$\log L[u, v] = -\frac{a}{2}(u - u_\infty)^2 - \frac{b}{2}(v - v_\infty)^2$$





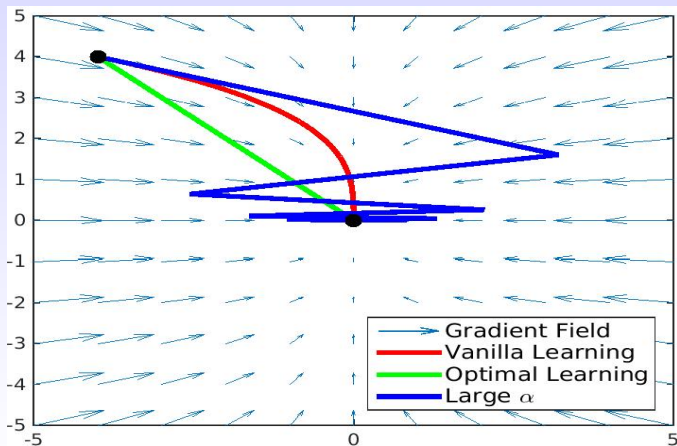
## A 2-dimensional example:

$$\log L[u, v] = -\frac{a}{2}(u - u_\infty)^2 - \frac{b}{2}(v - v_\infty)^2$$



## A 2-dimensional example:

$$\log L[u, v] = -\frac{a}{2}(u - u_\infty)^2 - \frac{b}{2}(v - v_\infty)^2$$



A 2-dimensional example:

$$\log L[u, v] = -\frac{a}{2}(u - u_\infty)^2 - \frac{b}{2}(v - v_\infty)^2$$

**Vanilla Gradient:**

$$\delta u_t^{\text{VG}} \sim (1 - \alpha a)^{-t} \Rightarrow \alpha < 2/a; \quad \delta v_t^{\text{VG}} \sim (1 - \alpha b)^{-t} \Rightarrow \alpha < 2/b$$

A 2-dimensional example:

$$\log L[u, v] = -\frac{a}{2}(u - u_\infty)^2 - \frac{b}{2}(v - v_\infty)^2$$

**Vanilla Gradient:**

$$\delta u_t^{\text{VG}} \sim (1 - \alpha a)^{-t} \Rightarrow \alpha < 2/a; \quad \delta v_t^{\text{VG}} \sim (1 - \alpha b)^{-t} \Rightarrow \alpha < 2/b$$

**Newton Method:**

$$\delta u_t^{\text{VG}} \sim (1 - \alpha)^{-t} \Rightarrow \alpha < 2; \quad \delta v_t^{\text{VG}} \sim (1 - \alpha)^{-t} \Rightarrow \alpha < 2$$

A 2-dimensional example:

$$\log L[u, v] = -\frac{a}{2}(u - u_\infty)^2 - \frac{b}{2}(v - v_\infty)^2$$

**Vanilla Gradient:**

$$\delta u_t^{\text{VG}} \sim (1 - \alpha a)^{-t} \Rightarrow \alpha < 2/a; \quad \delta v_t^{\text{VG}} \sim (1 - \alpha b)^{-t} \Rightarrow \alpha < 2/b$$

**Newton Method:**

$$\delta u_t^{\text{VG}} \sim (1 - \alpha)^{-t} \Rightarrow \alpha < 2; \quad \delta v_t^{\text{VG}} \sim (1 - \alpha)^{-t} \Rightarrow \alpha < 2$$

$$\alpha = 1 \quad \Rightarrow \quad \text{convergence in one step!}$$

- 1 Maximal Entropy Models and the Vanilla (standard) Learning Algorithm
- 2 Approximated Newton Method**
- 3 The Long-Time Limit: Stochastic Dynamics
- 4 Properties of the Stationary Distribution

The same happens for the MaxEnt inference:

$$\log L[\mathbf{X} \approx \mathbf{X}^*] \approx \log L[\mathbf{X}^*] - \frac{1}{2} \sum_{ab} (X_a - X_a^*) \chi[\mathbf{X}^*]_{ab} (X_b - X_b^*)$$

$$\chi_{ab}[\mathbf{X}] \equiv -\frac{\partial^2 \log L[\mathbf{X}]}{\partial X_a \partial X_b} = \langle \Sigma_a \Sigma_b \rangle_{\mathbf{X}} - \langle \Sigma_a \rangle_{\mathbf{X}} \langle \Sigma_b \rangle_{\mathbf{X}}$$

The same happens for the MaxEnt inference:

$$\log L[\mathbf{X} \approx \mathbf{X}^*] \approx \log L[\mathbf{X}^*] - \frac{1}{2} \sum_{ab} (X_a - X_a^*) \chi[\mathbf{X}^*]_{ab} (X_b - X_b^*)$$

$$\chi_{ab}[\mathbf{X}] \equiv -\frac{\partial^2 \log L[\mathbf{X}]}{\partial X_a \partial X_b} = \langle \Sigma_a \Sigma_b \rangle_{\mathbf{X}} - \langle \Sigma_a \rangle_{\mathbf{X}} \langle \Sigma_b \rangle_{\mathbf{X}}$$

**Vanilla Gradient:**  $\delta X_t^{\text{VG}} = \alpha \nabla \log L[\mathbf{X}_{t-1}]$

$$\delta X_t^\mu \equiv \sum_a V_a^\mu \delta X_{a,t} \sim (1 - \alpha \lambda_\mu)^{-t}$$



The same happens for the MaxEnt inference:

$$\log L[\mathbf{X} \approx \mathbf{X}^*] \approx \log L[\mathbf{X}^*] - \frac{1}{2} \sum_{ab} (X_a - X_a^*) \chi[\mathbf{X}^*]_{ab} (X_b - X_b^*)$$

$$\chi_{ab}[\mathbf{X}] \equiv -\frac{\partial^2 \log L[\mathbf{X}]}{\partial X_a \partial X_b} = \langle \Sigma_a \Sigma_b \rangle_{\mathbf{X}} - \langle \Sigma_a \rangle_{\mathbf{X}} \langle \Sigma_b \rangle_{\mathbf{X}}$$

**Vanilla Gradient:**  $\delta X_t^{\text{VG}} = \alpha \nabla \log L[\mathbf{X}_{t-1}]$

$$\delta X_t^\mu \equiv \sum_a V_a^\mu \delta X_{a,t} \sim (1 - \alpha \lambda_\mu)^{-t}$$

**Newton Method**<sup>1</sup>:  $\delta X_t^{\text{NM}} = \alpha \chi^{-1}[\mathbf{X}_{t-1}] \nabla \log L[\mathbf{X}_{t-1}]$

$$\delta X_t^\mu \equiv \sum_a V_a^\mu \delta X_{a,t} \sim (1 - \alpha)^{-t}$$

<sup>1</sup>(here equivalent to Amari98 Natural Gradient)

The same happens for the MaxEnt inference:

$$\log L[\mathbf{X} \approx \mathbf{X}^*] \approx \log L[\mathbf{X}^*] - \frac{1}{2} \sum_{ab} (X_a - X_a^*) \chi[\mathbf{X}^*]_{ab} (X_b - X_b^*)$$

$$\chi_{ab}[\mathbf{X}] \equiv -\frac{\partial^2 \log L[\mathbf{X}]}{\partial X_a \partial X_b} = \langle \Sigma_a \Sigma_b \rangle_{\mathbf{X}} - \langle \Sigma_a \rangle_{\mathbf{X}} \langle \Sigma_b \rangle_{\mathbf{X}}$$

**Vanilla Gradient:**  $\delta X_t^{\text{VG}} = \alpha \nabla \log L[\mathbf{X}_{t-1}]$

$$\delta X_t^\mu \equiv \sum_a V_a^\mu \delta X_{a,t} \sim (1 - \alpha \lambda_\mu)^{-t}$$

**Newton Method**<sup>1</sup>:  $\delta X_t^{\text{NM}} = \alpha \chi^{-1}[\mathbf{X}_{t-1}] \nabla \log L[\mathbf{X}_{t-1}]$

$$\delta X_t^\mu \equiv \sum_a V_a^\mu \delta X_{a,t} \sim (1 - \alpha)^{-t}$$

VERY SLOW: expensive estimation & inversion of  $\chi[\mathbf{X}]$

<sup>1</sup>(here equivalent to Amari98 Natural Gradient)

However, for the Ising model we can approximate:

$$\chi_{ab}[\mathbf{X}^*] \approx \overline{\chi}_{ab} \equiv \langle \Sigma_a \Sigma_b \rangle_{\text{DATA}} - \langle \Sigma_a \rangle_{\text{DATA}} \langle \Sigma_b \rangle_{\text{DATA}}$$

However, for the Ising model we can approximate:

$$\chi_{ab}[\mathbf{X}^*] \approx \overline{\chi}_{ab} \equiv \langle \Sigma_a \Sigma_b \rangle_{\text{DATA}} - \langle \Sigma_a \rangle_{\text{DATA}} \langle \Sigma_b \rangle_{\text{DATA}}$$

Approximated Newton (AN) Method:

$$\delta X_t^{\text{AN}} = \alpha \overline{\chi}^{-1} \nabla \log L[\mathbf{X}_{t-1}]$$

However, for the Ising model we can approximate:

$$\chi_{ab}[\mathbf{X}^*] \approx \overline{\chi}_{ab} \equiv \langle \Sigma_a \Sigma_b \rangle_{\text{DATA}} - \langle \Sigma_a \rangle_{\text{DATA}} \langle \Sigma_b \rangle_{\text{DATA}}$$

Approximated Newton (AN) Method:

$$\delta X_t^{\text{AN}} = \alpha \overline{\chi}^{-1} \nabla \log L[\mathbf{X}_{t-1}]$$

Remarks on  $\chi[\mathbf{X}^*] \approx \overline{\chi}$ :

- equivalent to say that an Ising distribution properly describes data.
- states that the model Fisher is close to the observables co-variance.

As the algorithm works iteratively, it requires an  
early-stop condition

As the algorithm works iteratively, it requires an  
early-stop condition

idea: stop the algorithm when

$Q[\mathbf{X}]$  is statistically compatible with  $\mathbf{P}$

using the  $\mathbf{P}$ -covariance  $\overline{\chi}/B$

As the algorithm works iteratively, it requires an  
early-stop condition

idea: stop the algorithm when

$\mathbf{Q}[\mathbf{X}]$  is statistically compatible with  $\mathbf{P}$

using the  $\mathbf{P}$ -covariance  $\overline{\chi}/B$

$$\epsilon(\mathbf{P}, \mathbf{Q}[\mathbf{X}]) \equiv \frac{B}{2D} \sum_{ab} (P_a - Q_a) (\overline{\chi}^{-1})_{ab} (P_b - Q_b)$$

quantifies the distance between  $\mathbf{Q}[\mathbf{X}]$  and  $\mathbf{P}$  in the  $\overline{\chi}/B$  metric.



As the algorithm works iteratively, it requires an  
early-stop condition

idea: stop the algorithm when

$\mathbf{Q}[\mathbf{X}]$  is statistically compatible with  $\mathbf{P}$

using the  $\mathbf{P}$ -covariance  $\overline{\chi}/B$

$$\epsilon(\mathbf{P}, \mathbf{Q}[\mathbf{X}]) \equiv \frac{B}{2D} \sum_{ab} (P_a - Q_a) (\overline{\chi}^{-1})_{ab} (P_b - Q_b)$$

quantifies the distance between  $\mathbf{Q}[\mathbf{X}]$  and  $\mathbf{P}$  in the  $\overline{\chi}/B$  metric.

For two *i.i.d* data-sets:  $\epsilon(\mathbf{P}, \mathbf{P}') \approx 1$

$\Rightarrow$  we stop the algorithm as soon as  $\epsilon < 1$

## APPROXIMATED NEWTON ALGORITHM:

### 1 Initialization:

- (a) Chose  $\mathbf{X}_0$  and compute  $\mathbf{Q}[\mathbf{X}_0]$  and  $\epsilon_0 = \epsilon(\mathbf{P}, \mathbf{Q}[\mathbf{X}_0])$
- (b) Then set  $\alpha_0 = 1$  and  $M = \min(\frac{2B}{\epsilon_0}, B)$  MCMC samplings

## APPROXIMATED NEWTON ALGORITHM:

### 1 Initialization:

- (a) Chose  $\mathbf{X}_0$  and compute  $\mathbf{Q}[\mathbf{X}_0]$  and  $\epsilon_0 = \epsilon(\mathbf{P}, \mathbf{Q}[\mathbf{X}_0])$
- (b) Then set  $\alpha_0 = 1$  and  $M = \min(\frac{2B}{\epsilon_0}, B)$  MCMC samplings

### 2 Iterate the following step:

- (a) update the  $\mathbf{X}_t$
- (b) estimate  $\mathbf{Q}[\mathbf{X}_t]$  with  $M = \min(\frac{2B}{\epsilon_{t-1}}, B)$  MCMC samplings
- (c) compute  $\epsilon_t = \epsilon(\mathbf{P}, \mathbf{Q}[\mathbf{X}_t])$ ,
  - (d1)  $\epsilon_t < \epsilon_{t-1}$ : accept the update and increase  $\alpha$
  - (d2)  $\epsilon_t > \epsilon_{t-1}$ : discard the update, lower  $\alpha$  and re-estimate  $\mathbf{Q}[\mathbf{X}_t]$ .

## APPROXIMATED NEWTON ALGORITHM:

### 1 Initialization:

- (a) Chose  $\mathbf{X}_0$  and compute  $\mathbf{Q}[\mathbf{X}_0]$  and  $\epsilon_0 = \epsilon(\mathbf{P}, \mathbf{Q}[\mathbf{X}_0])$
- (b) Then set  $\alpha_0 = 1$  and  $M = \min(\frac{2B}{\epsilon_0}, B)$  MCMC samplings

### 2 Iterate the following step:

- (a) update the  $\mathbf{X}_t$
- (b) estimate  $\mathbf{Q}[\mathbf{X}_t]$  with  $M = \min(\frac{2B}{\epsilon_{t-1}}, B)$  MCMC samplings
- (c) compute  $\epsilon_t = \epsilon(\mathbf{P}, \mathbf{Q}[\mathbf{X}_t])$ ,
  - (d1)  $\epsilon_t < \epsilon_{t-1}$ : accept the update and increase  $\alpha$
  - (d2)  $\epsilon_t > \epsilon_{t-1}$ : discard the update, lower  $\alpha$  and re-estimate  $\mathbf{Q}[\mathbf{X}_t]$ .

### 3 stop the algorithm when $\epsilon_t < 1$ .

## Rat retina ganglion cells

Two moving bars.

2.1h of MEA recording

$B = 4.8 \cdot 10^5$  of  $\Delta t = 16ms$

$N = 95$  cells

$D = 4560$  parameters to infer.

## Rat retina ganglion cells

Two moving bars.

2.1h of MEA recording

$$B = 4.8 \cdot 10^5 \text{ of } \Delta t = 16ms$$

$$N = 95 \text{ cells}$$

$D = 4560$  parameters to infer.

Convergence time from  
independent spins model  
with  $8 \times 3.4\text{Ghz}$  CPUs:

$$T^{\text{AN}} = 144 \pm 4s$$

$$T^{\text{VG}}(\alpha = 0.15) = 4.2 \cdot 10^4s$$

## Rat retina ganglion cells

Two moving bars.

2.1h of MEA recording

$B = 4.8 \cdot 10^5$  of  $\Delta t = 16ms$

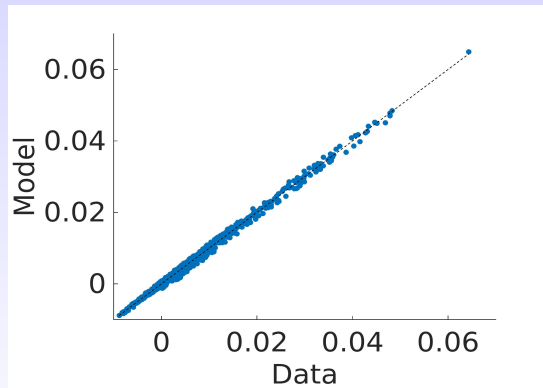
$N = 95$  cells

$D = 4560$  parameters to infer.

Convergence time from  
independent spins model  
with  $8 \times 3.4GHz$  CPUs:

$$T^{AN} = 144 \pm 4s$$

$$T^{VG}(\alpha = 0.15) = 4.2 \cdot 10^4s$$



$$c_{ij} = \langle \sigma_i \sigma_j \rangle - \langle \sigma_i \rangle \langle \sigma_j \rangle$$

## Rat retina ganglion cells

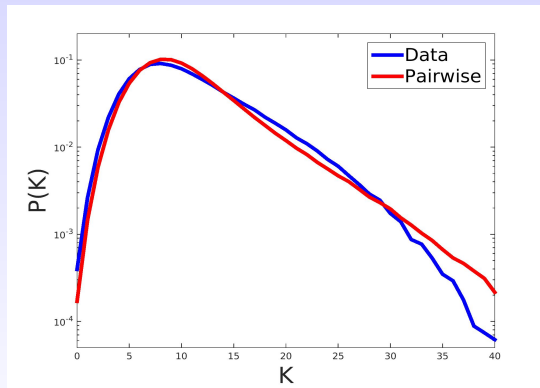
Two moving bars.

2.1h of MEA recording

 $B = 4.8 \cdot 10^5$  of  $\Delta t = 16ms$  $N = 95$  cells $D = 4560$  parameters to infer.Convergence time from  
independent spins model  
with  $8 \times 3.4GHz$  CPUs:

$$T^{AN} = 144 \pm 4s$$

$$T^{VG}(\alpha = 0.15) = 4.2 \cdot 10^4s$$



$$P(K) = \text{Prob}(\sum_i \sigma_i = K)$$



- 1 Maximal Entropy Models and the Vanilla (standard) Learning Algorithm
- 2 Approximated Newton Method
- 3 The Long-Time Limit: Stochastic Dynamics**
- 4 Properties of the Stationary Distribution

$Q[\mathbf{X}]$  is estimated through  $M$  MCMC measurements.

$Q[\mathbf{X}] \Rightarrow Q[\mathbf{X}]^{\text{MC}}$  is random variable!

$Q[\mathbf{X}]$  is estimated through  $M$  MCMC measurements.

$Q[\mathbf{X}] \Rightarrow Q[\mathbf{X}]^{\text{MC}}$  is random variable!

$\nabla \log L_{\mathbf{X}}^{\text{MC}} = \mathbf{P} - Q[\mathbf{X}]^{\text{MC}} \rightarrow 0$  only on average,

$Q[\mathbf{X}]$  is estimated through  $M$  MCMC measurements.

$Q[\mathbf{X}] \Rightarrow Q[\mathbf{X}]^{\text{MC}}$  is random variable!

$\nabla \log L_{\mathbf{X}}^{\text{MC}} = \mathbf{P} - Q[\mathbf{X}]^{\text{MC}} \rightarrow 0$  only on average,

Change of Framework:

$$\mathbf{X}_t \rightarrow P_t(\mathbf{X})$$

$\mathbf{X}$ , rather than converge to a fixed point,

approaches a stationary  $P_{\infty}(\mathbf{X})$

$Q[\mathbf{X}]$  is estimated through  $M$  MCMC measurements.

$Q[\mathbf{X}] \Rightarrow Q[\mathbf{X}]^{\text{MC}}$  is random variable!

$\nabla \log L_{\mathbf{X}}^{\text{MC}} = \mathbf{P} - Q[\mathbf{X}]^{\text{MC}} \rightarrow 0$  only on average,

### Change of Framework:

$$\mathbf{X}_t \rightarrow P_t(\mathbf{X})$$

$\mathbf{X}$ , rather than converge to a fixed point,  
approaches a stationary  $P_{\infty}(\mathbf{X})$

### Master Equation:

$$P_{t+1}(\mathbf{X}') = \int D\mathbf{X} P_t(\mathbf{X}) W_{\mathbf{X} \rightarrow \mathbf{X}'}(\alpha)$$

For  $M \gg 1$  and  $\mathbf{X} \approx \mathbf{X}^*$ :

$$\log L[\mathbf{X}] \simeq \log L[\mathbf{X}^*] - \frac{1}{2} \sum_{ab} (X_a - X_a^*) \chi[\mathbf{X}^*]_{ab} (X_b - X_b^*)$$

For  $M \gg 1$  and  $\mathbf{X} \approx \mathbf{X}^*$ :

$$\log L[\mathbf{X}] \approx \log L[\mathbf{X}^*] - \frac{1}{2} \sum_{ab} (X_a - X_a^*) \chi[\mathbf{X}^*]_{ab} (X_b - X_b^*)$$

$$\mathbf{1} \quad \langle \nabla_a \log L_{\mathbf{X}}^{\text{MC}} \rangle = \sum_b \chi[\mathbf{X}^*]_{ab} (X_b^* - X_b) \approx \sum_b \overline{\chi}_{ab} (X_b^* - X_b)$$

For  $M \gg 1$  and  $\mathbf{X} \approx \mathbf{X}^*$ :

$$\log L[\mathbf{X}] \simeq \log L[\mathbf{X}^*] - \frac{1}{2} \sum_{ab} (X_a - X_a^*) \chi[\mathbf{X}^*]_{ab} (X_b - X_b^*)$$

$$1 \quad \langle \nabla_a \log L_{\mathbf{X}}^{\text{MC}} \rangle = \sum_b \chi[\mathbf{X}^*]_{ab} (X_b^* - X_b) \approx \sum_b \overline{\chi}_{ab} (X_b^* - X_b)$$

$$2 \quad \left\langle \nabla_a \log L_{\mathbf{X}}^{\text{MC}} \nabla_b \log L_{\mathbf{X}}^{\text{MC}} \right\rangle_c = \chi[\mathbf{X}]_{ab} / M \simeq \chi[\mathbf{X}^*]_{ab} / M \approx \overline{\chi}_{ab} / M$$



For  $M \gg 1$  and  $\mathbf{X} \approx \mathbf{X}^*$ :

$$\log L[\mathbf{X}] \simeq \log L[\mathbf{X}^*] - \frac{1}{2} \sum_{ab} (X_a - X_a^*) \chi[\mathbf{X}^*]_{ab} (X_b - X_b^*)$$

$$1 \quad \langle \nabla_a \log L_{\mathbf{X}}^{\text{MC}} \rangle = \sum_b \chi[\mathbf{X}^*]_{ab} (X_b^* - X_b) \approx \sum_b \overline{\chi}_{ab} (X_b^* - X_b)$$

$$2 \quad \left\langle \nabla_a \log L_{\mathbf{X}}^{\text{MC}} \nabla_b \log L_{\mathbf{X}}^{\text{MC}} \right\rangle_c = \chi[\mathbf{X}]_{ab} / M \simeq \chi[\mathbf{X}^*]_{ab} / M \approx \overline{\chi}_{ab} / M$$

a normal approximation gives:

$$P(\nabla \log L_{\mathbf{X}}^{\text{MC}}) \simeq \mathcal{N} \left[ \overline{\chi} \cdot (\mathbf{X}^* - \mathbf{X}) ; \overline{\chi} / M \right] (\nabla \log L_{\mathbf{X}}^{\text{MC}})$$

For  $M \gg 1$  and  $\mathbf{X} \approx \mathbf{X}^*$ :

$$\log L[\mathbf{X}] \simeq \log L[\mathbf{X}^*] - \frac{1}{2} \sum_{ab} (X_a - X_a^*) \chi[\mathbf{X}^*]_{ab} (X_b - X_b^*)$$

$$\mathbf{1} \quad \langle \nabla_a \log L_{\mathbf{X}}^{\text{MC}} \rangle = \sum_b \chi[\mathbf{X}^*]_{ab} (X_b^* - X_b) \approx \sum_b \overline{\chi}_{ab} (X_b^* - X_b)$$

$$\mathbf{2} \quad \left\langle \nabla_a \log L_{\mathbf{X}}^{\text{MC}} \nabla_b \log L_{\mathbf{X}}^{\text{MC}} \right\rangle_c = \chi[\mathbf{X}]_{ab} / M \simeq \chi[\mathbf{X}^*]_{ab} / M \approx \overline{\chi}_{ab} / M$$

a normal approximation gives:

$$P(\nabla \log L_{\mathbf{X}}^{\text{MC}}) \simeq \mathcal{N} \left[ \overline{\chi} \cdot (\mathbf{X}^* - \mathbf{X}); \overline{\chi} / M \right] (\nabla \log L_{\mathbf{X}}^{\text{MC}})$$

$$\mathbf{■} \quad W_{\mathbf{X} \rightarrow \mathbf{X}'}^{\text{VG}}(\alpha) = \text{Prob} \left( \nabla \log L_{\mathbf{X}}^{\text{MC}} = \frac{\mathbf{X}' - \mathbf{X}}{\alpha} \right)$$

$$\mathbf{■} \quad W_{\mathbf{X} \rightarrow \mathbf{X}'}^{\text{AN}}(\alpha) = \text{Prob} \left( \nabla \log L_{\mathbf{X}}^{\text{MC}} = \overline{\chi} \cdot \frac{\mathbf{X}' - \mathbf{X}}{\alpha} \right)$$

Imposing  $P_{t+1}(\mathbf{X}) = P_t(\mathbf{X})$

- $P_{\infty}^{\text{VG}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M} \left(2\delta - \alpha \overline{\chi}\right)^{-1}\right](\mathbf{X})$
- $P_{\infty}^{\text{AN}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M(2-\alpha)} \overline{\chi}^{-1}\right](\mathbf{X})$

Imposing  $P_{t+1}(\mathbf{X}) = P_t(\mathbf{X})$

- $P_{\infty}^{\text{VG}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M} \left(2\delta - \alpha \overline{\chi}\right)^{-1}\right](\mathbf{X}), \quad \alpha \lambda_{\mu} < 2$
- $P_{\infty}^{\text{AN}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M(2-\alpha)} \overline{\chi}^{-1}\right](\mathbf{X}), \quad \alpha < 2$

Imposing  $P_{t+1}(\mathbf{X}) = P_t(\mathbf{X})$

- $P_{\infty}^{\text{VG}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M} \left(2\delta - \alpha \overline{\chi}\right)^{-1}\right](\mathbf{X}), \quad \alpha \lambda_{\mu} < 2$
- $P_{\infty}^{\text{AN}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M(2-\alpha)} \overline{\chi}^{-1}\right](\mathbf{X}), \quad \alpha < 2$

Which self-consistently defines  $\mathbf{X} \approx \mathbf{X}^*$

Imposing  $P_{t+1}(\mathbf{X}) = P_t(\mathbf{X})$

- $P_{\infty}^{\text{VG}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M} \left(2\delta - \alpha \overline{\chi}\right)^{-1}\right](\mathbf{X}), \quad \alpha \lambda_{\mu} < 2$
- $P_{\infty}^{\text{AN}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M(2-\alpha)} \overline{\chi}^{-1}\right](\mathbf{X}), \quad \alpha < 2$

Which self-consistently defines  $\mathbf{X} \approx \mathbf{X}^*$

From  $P(\nabla \log L_{\mathbf{X}}^{\text{MC}}) = P(\mathbf{P} - \mathbf{Q}[\mathbf{X}]^{\text{MC}})$

- $P_{\infty}^{\text{VG}}(\mathbf{Q}^{\text{MC}}) = \mathcal{N}\left[\mathbf{P}; \frac{2}{M} \overline{\chi} \left(2\delta - \alpha \overline{\chi}\right)^{-1}\right](\mathbf{Q}^{\text{MC}})$
- $P_{\infty}^{\text{AN}}(\mathbf{Q}^{\text{MC}}) = \mathcal{N}\left[\mathbf{P}; \frac{2}{M(2-\alpha)} \overline{\chi}\right](\mathbf{Q}^{\text{MC}})$

Imposing  $P_{t+1}(\mathbf{X}) = P_t(\mathbf{X})$

- $P_{\infty}^{\text{VG}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M} \left(2\delta - \alpha \overline{\chi}\right)^{-1}\right](\mathbf{X}), \quad \alpha \lambda_{\mu} < 2$
- $P_{\infty}^{\text{AN}}(\mathbf{X}) = \mathcal{N}\left[\mathbf{X}^*; \frac{\alpha}{M(2-\alpha)} \overline{\chi}^{-1}\right](\mathbf{X}), \quad \alpha < 2$

Which self-consistently defines  $\mathbf{X} \approx \mathbf{X}^*$

From  $P(\nabla \log L_{\mathbf{X}}^{\text{MC}}) = P(\mathbf{P} - \mathbf{Q}[\mathbf{X}]^{\text{MC}})$

- $P_{\infty}^{\text{VG}}(\mathbf{Q}^{\text{MC}}) = \mathcal{N}\left[\mathbf{P}; \frac{2}{M} \overline{\chi} \left(2\delta - \alpha \overline{\chi}\right)^{-1}\right](\mathbf{Q}^{\text{MC}})$
- $P_{\infty}^{\text{AN}}(\mathbf{Q}^{\text{MC}}) = \mathcal{N}\left[\mathbf{P}; \frac{2}{M(2-\alpha)} \overline{\chi}\right](\mathbf{Q}^{\text{MC}})$

Which is better? How to set the parameters?

- 1 Maximal Entropy Models and the Vanilla (standard) Learning Algorithm
- 2 Approximated Newton Method
- 3 The Long-Time Limit: Stochastic Dynamics
- 4 Properties of the Stationary Distribution**



## Algorithm Vs Empirical distributions

An experiment provides

empirical estimates of  $\mathbf{Q}^{\text{EMP}}$ :

$$P^{\text{EMP}}(\mathbf{Q}^{\text{EMP}}) \simeq \mathcal{N}[\mathbf{P}^{\text{TRUE}}, \chi^{\text{EMP}}]$$

- $\mathbf{P}^{\text{TRUE}}$ : result from infinitely long experiment
- $\chi^{\text{EMP}}$  expected co-variance for  $B$  measurements

## Algorithm Vs Empirical distributions

An experiment provides

empirical estimates of  $\mathbf{Q}^{\text{EMP}}$ :

$$P^{\text{EMP}}(\mathbf{Q}^{\text{EMP}}) \simeq \mathcal{N}[\mathbf{P}^{\text{TRUE}}, \chi^{\text{EMP}}]$$

An inference algorithm provides

numerical estimates of  $\mathbf{Q}^{\text{MC}}$ :

$$P_{\mathbf{P}}^{\text{ALG}}(\mathbf{Q}^{\text{MC}}) \simeq \mathcal{N}[\mathbf{P}, \chi^{\text{ALG}}]$$

- $\mathbf{P}^{\text{TRUE}}$ : result from infinitely long experiment
- $\chi^{\text{EMP}}$  expected co-variance for  $B$  measurements
- $\mathbf{P}$  one-shot sampling of  $P^{\text{EMP}}$

## Algorithm Vs Empirical distributions

An experiment provides  
empirical estimates of  $\mathbf{Q}^{\text{EMP}}$ :

$$P^{\text{EMP}}(\mathbf{Q}^{\text{EMP}}) \simeq \mathcal{N}[\mathbf{P}^{\text{TRUE}}, \chi^{\text{EMP}}]$$

An inference algorithm provides  
numerical estimates of  $\mathbf{Q}^{\text{MC}}$ :

$$P_{\mathbf{P}}^{\text{ALG}}(\mathbf{Q}^{\text{MC}}) \simeq \mathcal{N}[\mathbf{P}, \chi^{\text{ALG}}]$$

- $\mathbf{P}^{\text{TRUE}}$ : result from infinitely long experiment
- $\chi^{\text{EMP}}$  expected co-variance for  $B$  measurements
- $\mathbf{P}$  one-shot sampling of  $P^{\text{EMP}}$

An optimal inference algorithm should provide:  
 $P^{\text{ALG}}$  **as close as possible** to  $P^{\text{EMP}}$ .

What is the optimal  $\chi^{\text{ALG}}$  value?

Kullback-Leibler distance between  $P^{\text{EMP}}$  and  $P_{\mathbf{P}}^{\text{ALG}}$ :

$$D_{KL} \left( P^{\text{EMP}}(\cdot) \parallel P_{\mathbf{P}}^{\text{ALG}}(\cdot) \right)$$

Kullback-Leibler distance between  $P^{\text{EMP}}$  and  $P_{\mathbf{P}}^{\text{ALG}}$ :

$$D_{KL} \left( P^{\text{EMP}}(\cdot) \parallel P_{\mathbf{P}}^{\text{ALG}}(\cdot) \right)$$

$$\chi^{\text{OPT}} = \arg \min_{\chi^{\text{ALG}}} \int \mathbf{DP} P^{\text{EMP}}(\mathbf{P}) D_{KL} \left( P^{\text{EMP}}(\cdot) \parallel P_{\mathbf{P}}^{\text{ALG}}(\cdot) \right)$$

Kullback-Leibler distance between  $P^{\text{EMP}}$  and  $P_{\mathbf{P}}^{\text{ALG}}$ :

$$D_{KL} \left( P^{\text{EMP}}(\cdot) \parallel P_{\mathbf{P}}^{\text{ALG}}(\cdot) \right)$$

$$\chi^{\text{OPT}} = \arg \min_{\chi^{\text{ALG}}} \int \mathbf{DP} P^{\text{EMP}}(\mathbf{P}) D_{KL} \left( P^{\text{EMP}}(\cdot) \parallel P_{\mathbf{P}}^{\text{ALG}}(\cdot) \right)$$

The solution and its approximation are:

$$\chi^{\text{OPT}} = 2\chi^{\text{EMP}} \approx 2\overline{\chi} / B$$

Kullback-Leibler distance between  $P^{\text{EMP}}$  and  $P_{\mathbf{P}}^{\text{ALG}}$ :

$$D_{KL} \left( P^{\text{EMP}}(\cdot) \parallel P_{\mathbf{P}}^{\text{ALG}}(\cdot) \right)$$

$$\chi^{\text{OPT}} = \arg \min_{\chi^{\text{ALG}}} \int \mathbf{DP} P^{\text{EMP}}(\mathbf{P}) D_{KL} \left( P^{\text{EMP}}(\cdot) \parallel P_{\mathbf{P}}^{\text{ALG}}(\cdot) \right)$$

The solution and its approximation are:

$$\chi^{\text{OPT}} = 2\chi^{\text{EMP}} \approx 2\overline{\chi} / B$$

to compare with:

$$\chi^{\text{VG}} = \frac{2}{M} \overline{\chi} (2\delta - \alpha \overline{\chi})^{-1}, \quad \chi^{\text{AN}} = \frac{2}{M(2-\alpha)} \overline{\chi}$$

Kullback-Leibler distance between  $P^{\text{EMP}}$  and  $P_{\mathbf{P}}^{\text{ALG}}$ :

$$D_{KL} \left( P^{\text{EMP}}(\cdot) \parallel P_{\mathbf{P}}^{\text{ALG}}(\cdot) \right)$$

$$\chi^{\text{OPT}} = \arg \min_{\chi^{\text{ALG}}} \int \mathbf{DP} P^{\text{EMP}}(\mathbf{P}) D_{KL} \left( P^{\text{EMP}}(\cdot) \parallel P_{\mathbf{P}}^{\text{ALG}}(\cdot) \right)$$

The solution and its approximation are:

$$\chi^{\text{OPT}} = 2\chi^{\text{EMP}} \approx 2\overline{\chi} / B$$

to compare with:

$$\chi^{\text{VG}} = \frac{2}{M} \overline{\chi} (2\delta - \alpha \overline{\chi})^{-1}, \quad \chi^{\text{AN}} = \frac{2}{M(2-\alpha)} \overline{\chi}$$

AN with  $M(2 - \alpha) = B$  reaches the optimum!

VG underfits  $\lambda_{\mu} \gg (2 - B/M)/\alpha$  and overfits  $\lambda_{\mu} \ll (2 - B/M)/\alpha$



## Synthetic data: Theory Vs Simulations

### Bethe Lattice Ising Model

$$N = 10, c = 4$$

$$J_{ij} = \pm 0.53,$$

$$h_i = -0.14 - 2 \sum_j J_{ij}$$

100 independent estimations

of  $\mathbf{P}$  and  $\overline{\chi}$

through  $2^{16}$  sampling of  $P^{\text{EMP}}$

Inference with  $M = B$

## Synthetic data: Theory Vs Simulations

## Bethe Lattice Ising Model

$$N = 10, c = 4$$

$$J_{ij} = \pm 0.53,$$

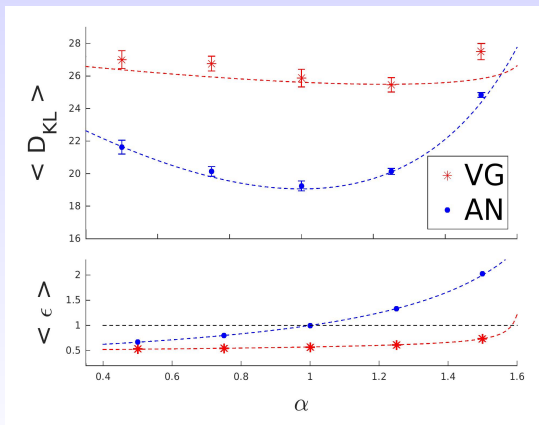
$$h_i = -0.14 - 2 \sum_j J_{ij}$$

100 independent estimations

of  $\mathbf{P}$  and  $\overline{\chi}$

through  $2^{16}$  sampling of  $P^{\text{EMP}}$

Inference with  $M = B$



## Conclusions:

- MaxEnt models are useful to describe multi-units systems
- The AN learning is faster than the VG algorithm.
- Within the large  $B$  approximation is possible to completely characterize the long time behavior
- The AN with  $\alpha = 1$  and  $M = B$  is optimal against overfitting.

## Perspectives:

- Improve the gaussian approximations
- Test the algorithm to non-pairwise models
- Generalize the class of model distributions beyond MaxEnt
- Include hidden variables and the RBM framework

# THANKS

## Collaborators for P.F. Cortex work:

- Francesco Battaglia
- Simona Cocco
- Remi Monasson
- Gaia Tavoni

## Founding

- EU-FP7 FET OPEN project Enlightenment 284801
- Human Brain Project (HBP CLAP)

arXiv:1507.04254