



Revisiting the training of RBMs

Beatriz Seoane

Theoretical Physics, UCM Madrid



Beatriz Seoane

Theoretical Physics, UCM Madrid

In collaboration with



Aurélien Decelle
(UCM)



Cyril Furtlehner
(Tau team, Université Paris-Saclay)

Nicolas Bereux

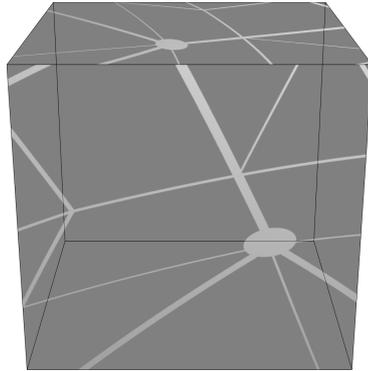


Student Paris-Saclay

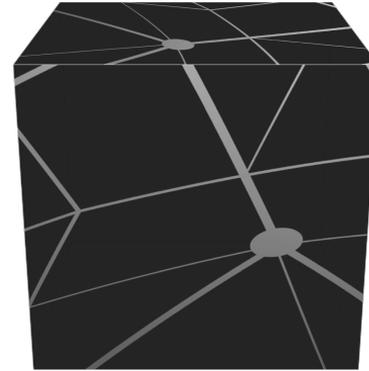
Introduction : generative approach

0
1
2
3
4
5
6
7
8
9

training



generating

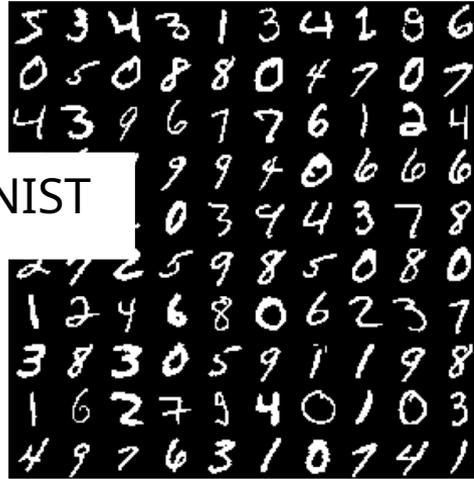


1
6
9
5
4
7
8
3
7
5

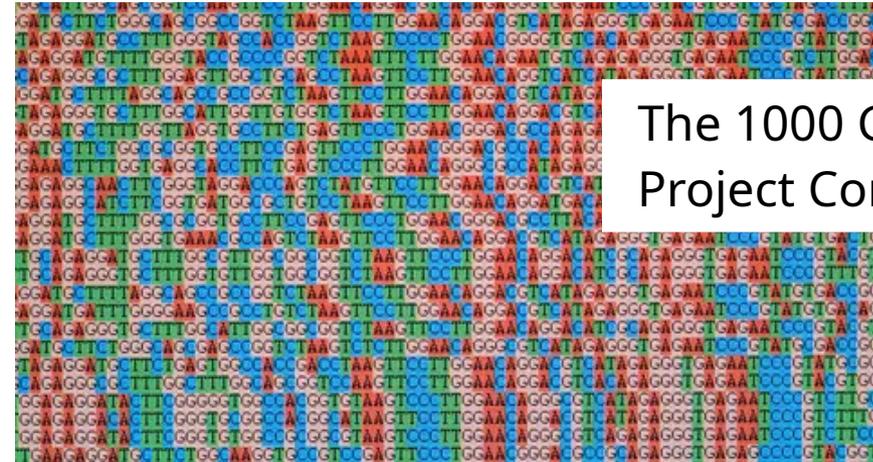
- **Energy based models (RBMs, Generative Convnets)**
- Generative Adverarial Network (GAN)
- Variational AutoEncoder (VAE)

Introduction : generative approach

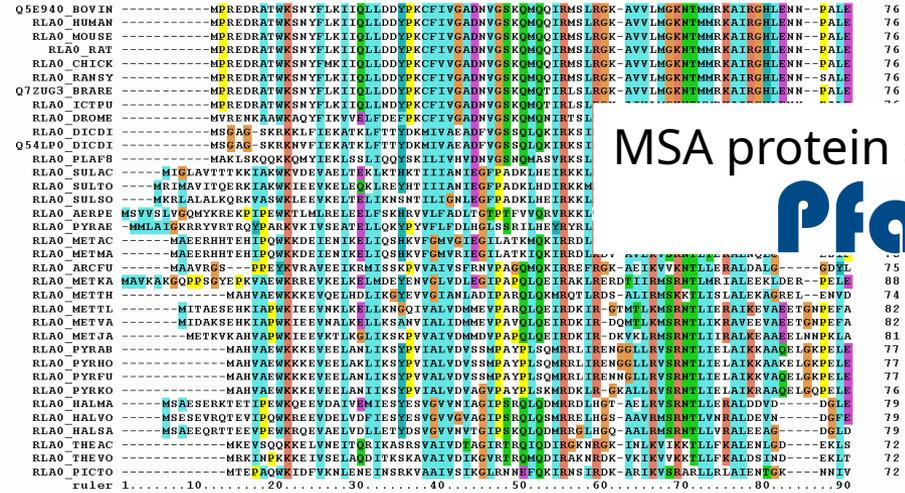
MNIST



CELEBA

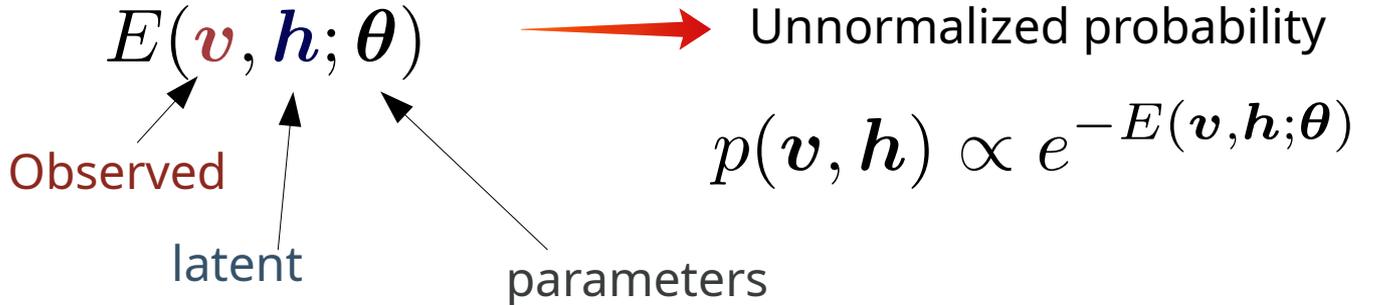


The 1000 Genomes Project Consortium



MSA protein sequences
Pfam

Energy based models (EBMs)

- Assign an energy $E(\mathbf{v}, \mathbf{h}; \theta)$ 

Unnormalized probability

$$p(\mathbf{v}, \mathbf{h}) \propto e^{-E(\mathbf{v}, \mathbf{h}; \theta)}$$

- Learning**: adjust the parameters so that the dataset configurations are typical configurations of the model.

- Maximize the **likelihood**:

$$L = \prod_{m=1}^M p\left(v^{(m)}\right)$$

Gradient ascent

$$\nabla_{\theta} L$$

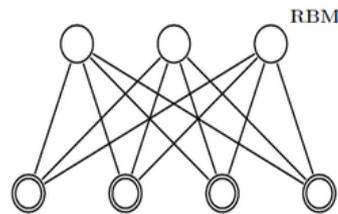
- Problem $Z = \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h})$ Is in generally impossible to compute exactly

We need Monte Carlo Sampling

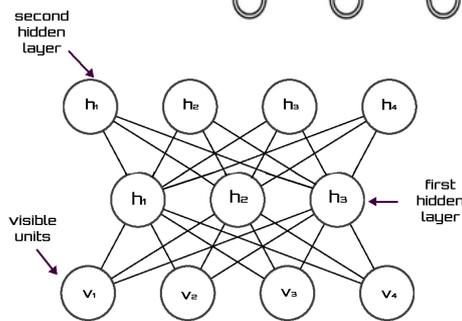
Energy based models (EBMs)

$$E(\mathbf{v}, \mathbf{h}; \theta)$$

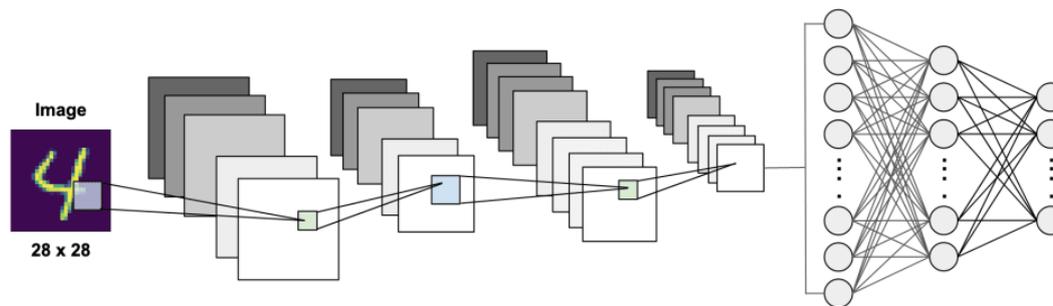
- **Restricted Boltzmann Machines**



- Deep Boltzmann Machines



- Generative ConvNets



Why Restricted Boltzmann Machines (RBMs)?

In this talk we focus on RBM for the following reason:

- It is an *Ising model*: a canonical model for statistical physicists
- We can write explicitly the probability distribution which give us tools to study it
- It relies on a simple shallow neural network that can be looked into: can it be used to **extract dataset features (interpretability)** ?
- It can model complex dataset : it can be shown in the binary case that it can overfit anything...

Reasons to not use it:

- The training can be (very) capricious for complicated reasons
- Generate new data can be long ... (How long?)
- So far, no implementation of a convolutional RBM has been made usable in practice
- ...

RBM (biased) Historical events

Machine Learning aspects

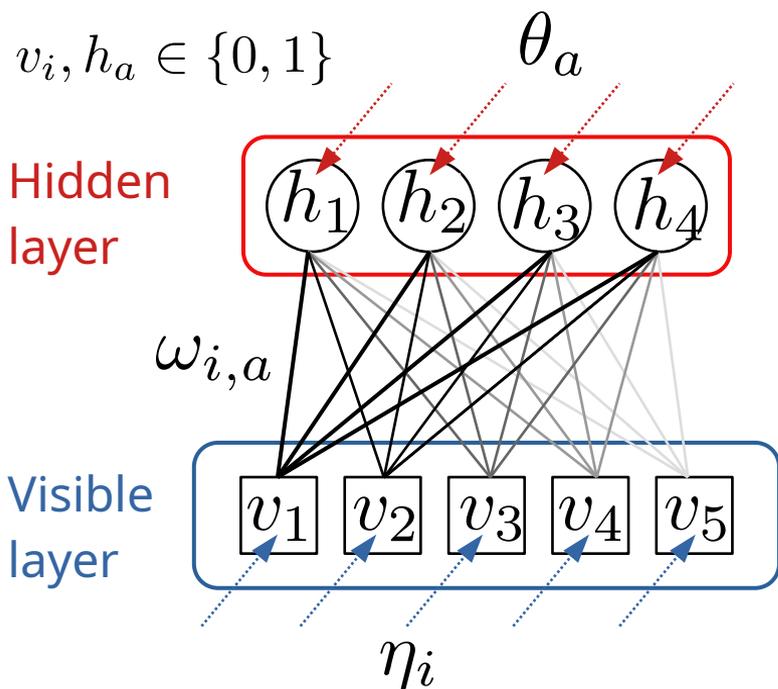
- It was introduced by Smolenski and popularized by Hinton ~80/90
→ introducing the Contrastive Divergence it was proved to be “practical”
- It was getting popular as a pre-training tools for deep-network
- Then, the interest for RBMs slowly decreases around ~2010
- The rise of GAN/VAE alternative achieve to out-faschion RBM

Statistical Physics aspects

- Many works on RBM in ~2010 analyzing the phase diagram of RBM
- Works on the learning dynamics
- Works on message-passing learning algorithm
- ...

Restricted Boltzmann Machine

Smolensky (1989)



Energy of a configuration

$$E[v, h; w, \eta, \theta] = - \sum_{ia} v_i w_{ia} h_a - \sum_i \eta_i v_i - \sum_a \theta_a h_a$$

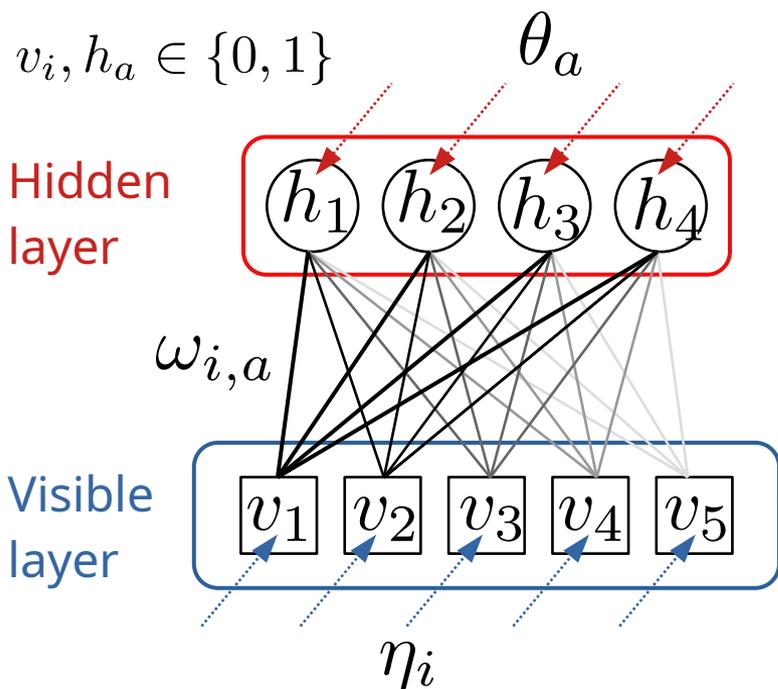
Visible : **data**



Hidden : "Neurons" → **features extracted**
(interpretability!)

Restricted Boltzmann Machine

Smolensky (1989)



Energy of a configuration

$$E[v, h; w, \eta, \theta] = - \sum_{ia} v_i w_{ia} h_a - \sum_i \eta_i v_i - \sum_a \theta_a h_a$$

Visible : **data**

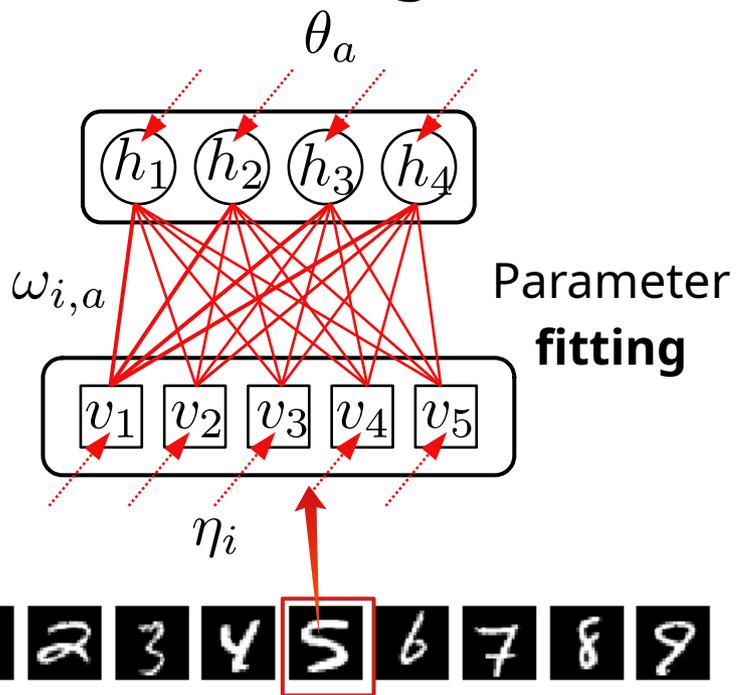


Hidden : "Neurons" → **features extracted**
(interpretability!)

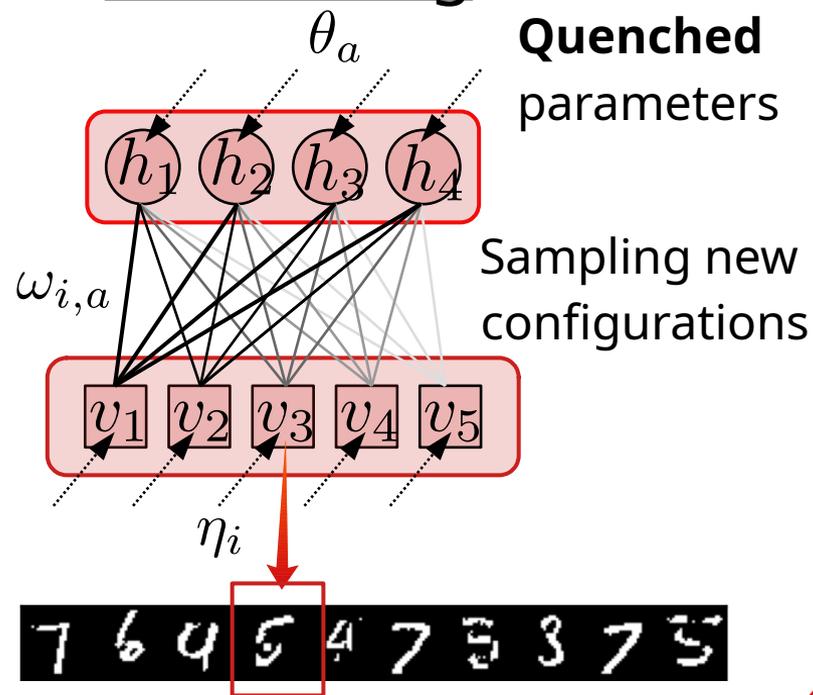
Sequencing context: Tubiana, Cocco, Monasson, *Elife* (2019); Shimagaki, Weigt, *PRE* (2019), Yelmen *et al.*, *PLoS genetics* (2021); Bravi *et al.* *PloS CB* (2021); Bravi *et al.* *Cell Systems* (2021); ⇒ *Tubiana's talk*

Training vs. Sampling and RBM

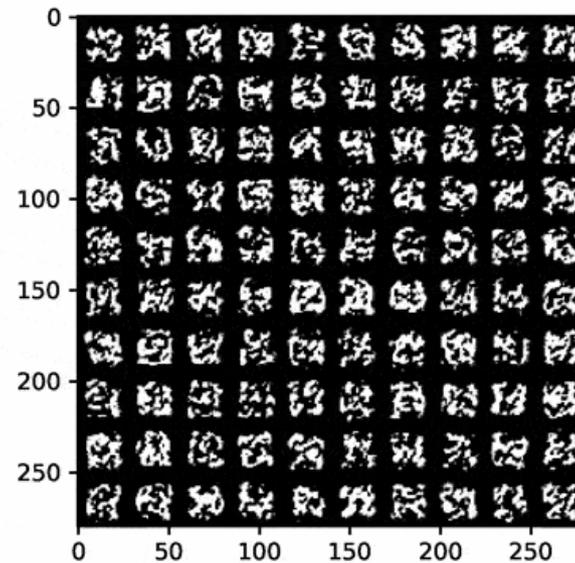
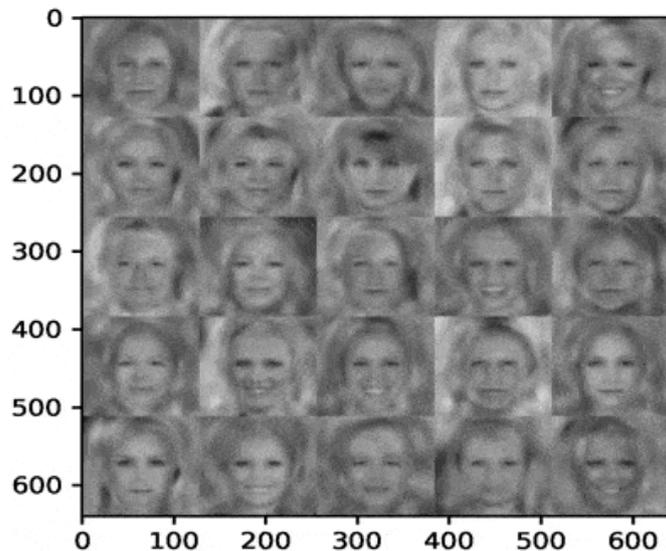
Learning



Generating



Gibbs sampling of a trained RBM



Learning an RBM

- Gibbs equilibrium distribution

$$p[\mathbf{v}, \mathbf{h} | \mathbf{w}, \boldsymbol{\eta}, \boldsymbol{\theta}] = \frac{\exp(-E[\mathbf{v}, \mathbf{h}; \mathbf{w}, \boldsymbol{\eta}, \boldsymbol{\theta}])}{Z} \quad \text{with } Z = \sum_{\{\mathbf{v}, \mathbf{h}\}} e^{-E[\mathbf{v}, \mathbf{h}]}$$

- Dataset $S = \{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(M)}\}$ are the **typical samples** of $p(\mathbf{v})$

- Maximize the log-likelihood $\mathcal{L}(\mathbf{w}, \boldsymbol{\eta}, \boldsymbol{\theta} | S) = \sum_{m=1}^M \ln p(\mathbf{v} = \mathbf{v}^{(m)} | \mathbf{w}, \boldsymbol{\eta}, \boldsymbol{\theta})$

- Gradient ascent
$$\frac{\partial \mathcal{L}}{\partial w_{ia}} = \langle v_i h_a \rangle_{\mathcal{D}} - \langle v_i h_a \rangle_{\mathcal{H}}$$
$$\frac{\partial \mathcal{L}}{\partial \eta_i} = \langle v_i \rangle_{\mathcal{D}} - \langle v_i \rangle_{\mathcal{H}} \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial \theta_a} = \langle h_a \rangle_{\mathcal{D}} - \langle h_a \rangle_{\mathcal{H}}$$

On the interpretability (I)

- One can extract the **point correlations** of our data up to any order !

Once we integrate out the hidden variables: $E(v; \theta) = -\log \left(\sum_{\mathbf{h}} e^{-E[\mathbf{v}, \mathbf{h}; \mathbf{w}, \boldsymbol{\eta}, \boldsymbol{\theta}]}$

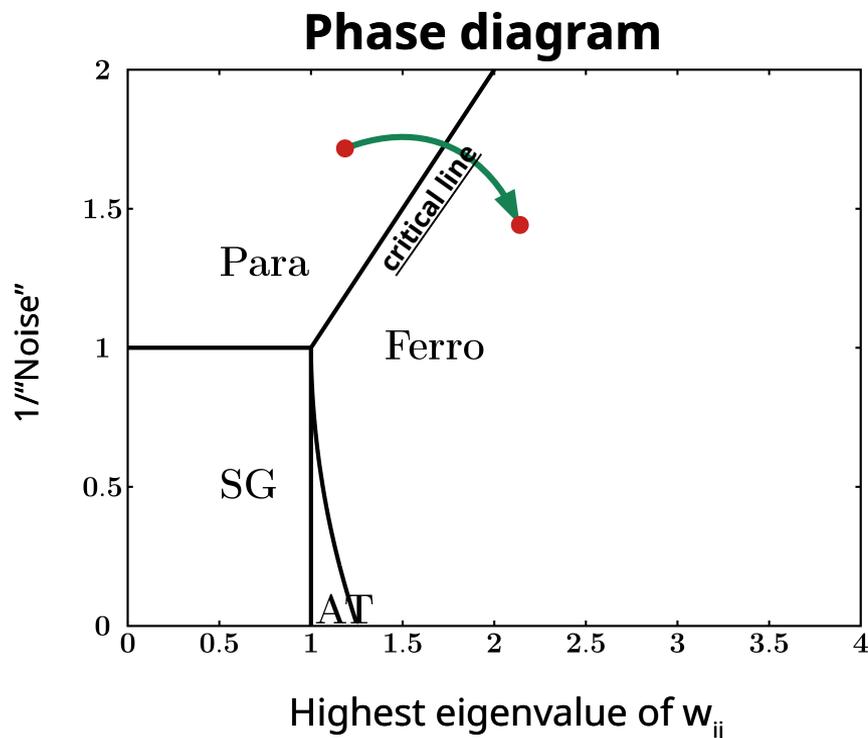
$$E(v) = -E_0 - \sum_i H_i v_i - \sum_{i,j} J_{i,j}^{(2)} v_i v_j - \sum_{i,j,k} J_{ijk}^{(3)} v_i v_j v_k \cdots - \sum_{j_1 \cdots j_n} J_{j_1 \cdots j_n}^{(n)} v_{j_1} \cdots v_{j_n} - \cdots$$

Effective model for the problem

On the interpretability (II)

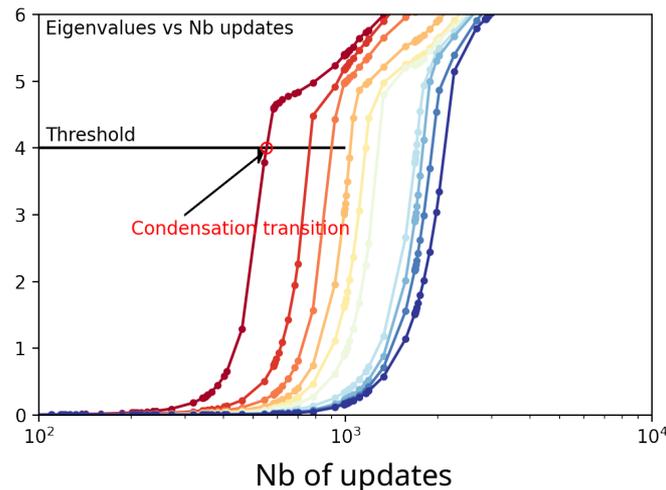
- It can be shown that when the RBM starts to learn features of the data, the system suffers a **phase order transition** from a paramagnetic phase to a ferromagnetic phase

[Decelle, Fissore, Furtlehner J. of stat phys (2018)]



Para: high temperature (low variance)

Ferro: strong eigenmode - low noise

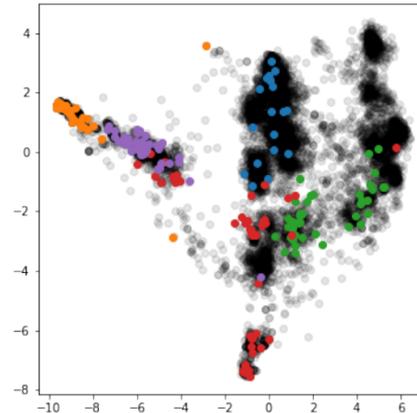
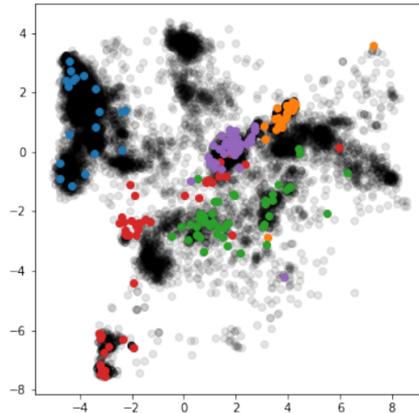
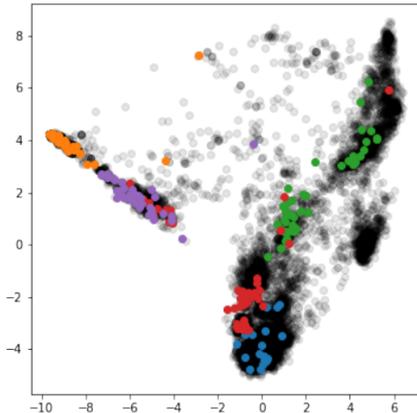


On the interpretability (II)

- The eigenvectors \mathbf{w}_i of matrix $\boldsymbol{\omega}$ align with the important directions of the dataset:

$$m_i = \langle \mathbf{v} \cdot \mathbf{w}_i \rangle_{\mathcal{D}} \neq 0$$

[Decelle, Fissore, Furtlehner J.
of stat phys (2018)]



CPF protein

Work in collaboration with
R. Vanderhaegen, A. Decelle,
A. Carbone

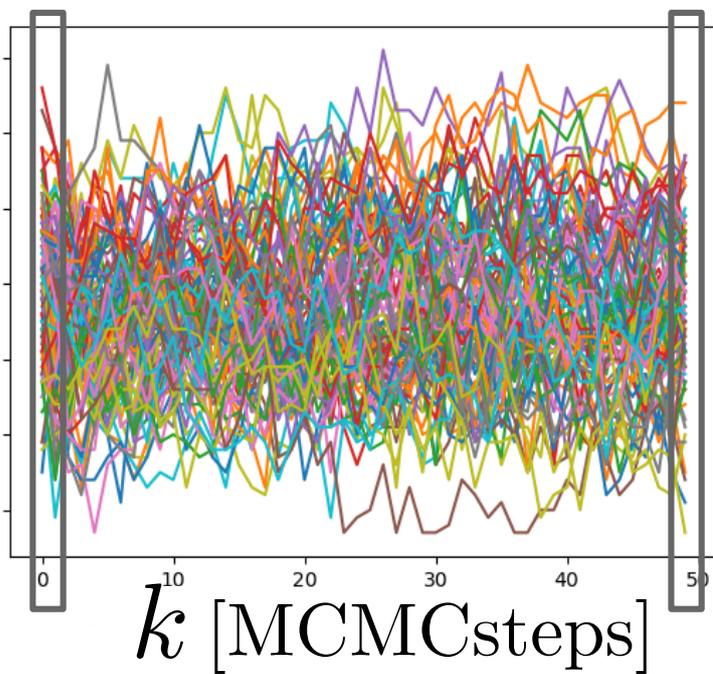
On the gradient

$$\langle f(\mathbf{v}, \mathbf{h}) \rangle_{\mathcal{D}} = M^{-1} \sum_m \sum_{\{\mathbf{h}\}} f(\mathbf{v}^{(m)}, \boldsymbol{\tau}) p(\mathbf{h} | \mathbf{v}^{(m)})$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_{ia}} &= \langle v_i h_a \rangle_{\mathcal{D}} - \langle v_i h_a \rangle_{\mathcal{H}} \\ \frac{\partial \mathcal{L}}{\partial \eta_i} &= \langle v_i \rangle_{\mathcal{D}} - \langle v_i \rangle_{\mathcal{H}} \\ \frac{\partial \mathcal{L}}{\partial \theta_a} &= \langle h_a \rangle_{\mathcal{D}} - \langle h_a \rangle_{\mathcal{H}} \end{aligned}$$

Monte Carlo

N_s parallel
Markov chains
initialization



Measure \Rightarrow r.h.s gradient

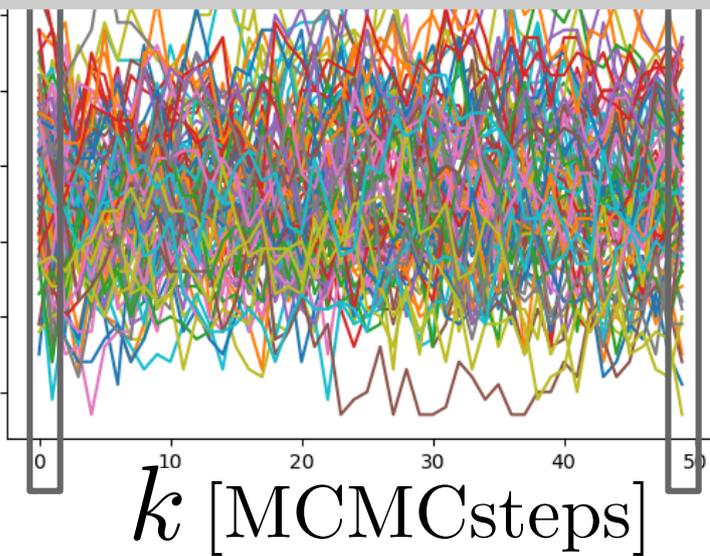
On the gradient

 $\langle \dots \rangle_{\mathcal{H}}$

R.h.s gradient will be correctly computed if the simulations thermalize

$$k \sim n\tau_{\text{mixing}}$$

N_s parallel
Markov chains
initialization



$$\frac{\partial \mathcal{L}}{\partial w_{ia}} = \langle v_i h_a \rangle_{\mathcal{D}} - \langle v_i h_a \rangle_{\mathcal{H}}$$
$$\frac{\partial \mathcal{L}}{\partial \eta_i} = \langle v_i \rangle_{\mathcal{D}} - \langle v_i \rangle_{\mathcal{H}}$$
$$\frac{\partial \mathcal{L}}{\partial \theta_a} = \langle h_a \rangle_{\mathcal{D}} - \langle h_a \rangle_{\mathcal{H}}$$

Monte Carlo

Measure \Rightarrow r.h.s gradient

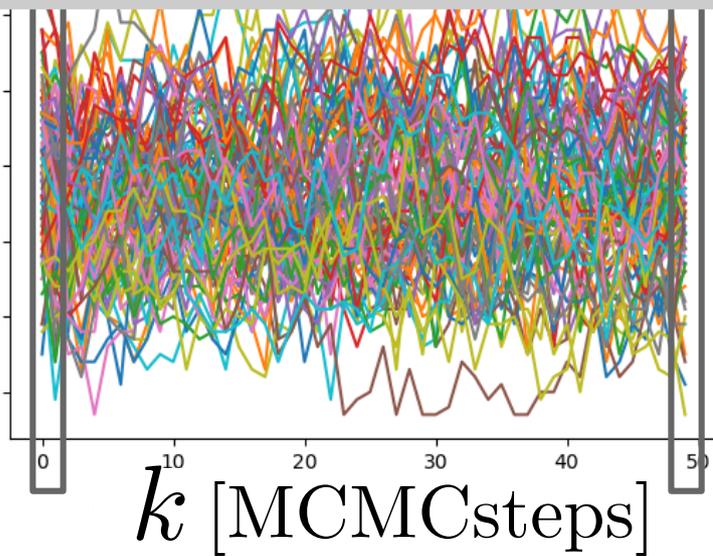
On the gradient

$\langle \dots \rangle_{\mathcal{H}}$

R.h.s gradient will be correctly computed if the simulations thermalize

$$k \sim n\tau_{\text{mixing}}$$

N_s parallel
Markov chains
initialization



$$\frac{\partial \mathcal{L}}{\partial w_{ia}} = \langle v_i h_a \rangle_{\mathcal{D}} - \langle v_i h_a \rangle_{\mathcal{H}}$$
$$\frac{\partial \mathcal{L}}{\partial \eta_i} = \langle v_i \rangle_{\mathcal{D}} - \langle v_i \rangle_{\mathcal{H}}$$
$$\frac{\partial \mathcal{L}}{\partial \dots} = \langle h \dots \rangle_{\mathcal{D}} - \langle h \dots \rangle_{\mathcal{H}}$$

Monte Carlo

Let's choose good starting point and approximate with $k \sim O(10)$ steps

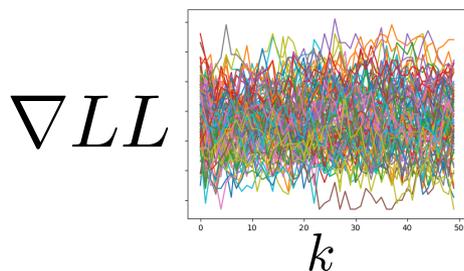
- Contrastive divergence (CD) [Hinton (2002)] **Init - dataset**
- Persistence contrastive divergence (PCD) [Tieleman (2008)] **Init - previous last point**
- Mean field (TAP) [Gabrié, Tramel, and Krzakala (2015)]

Or

- Simulated annealing, Parallel Tempering, ...

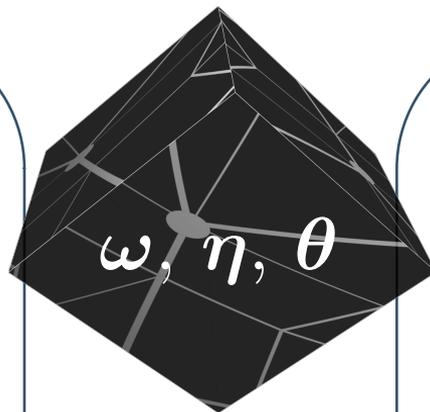
Equilibrium vs. Non-eq. regimes

Training

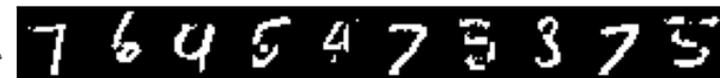


- **Non-equilibrium** $k < t_{\text{therm}}$

- **Equilibrium** $k > t_{\text{therm}}$



Sampling

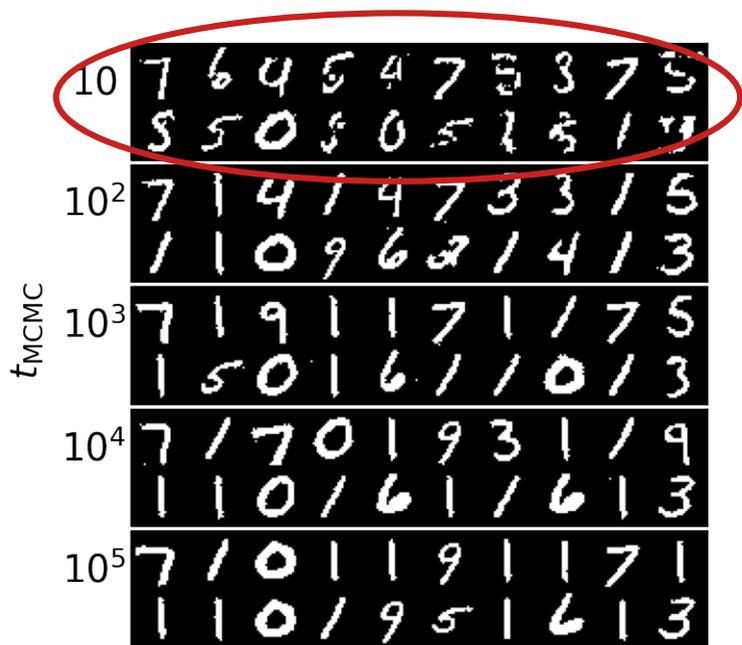


- “Learns the **dynamics**”

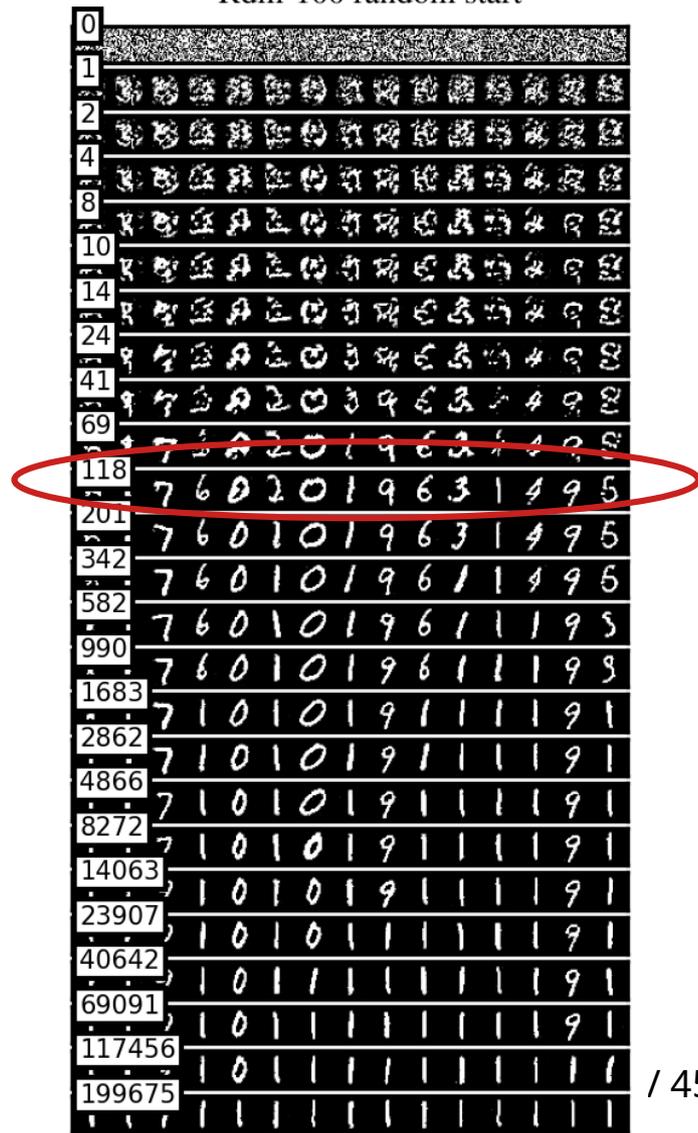
- Learns a good **model** for the data

Non-equilibrium regime

Rdm - 10



Rdm-100 random start



Non-equilibrium regime

Training using CD: chain initialisation at the dataset

If we sample the RBM from:

→ random **NOTHING**

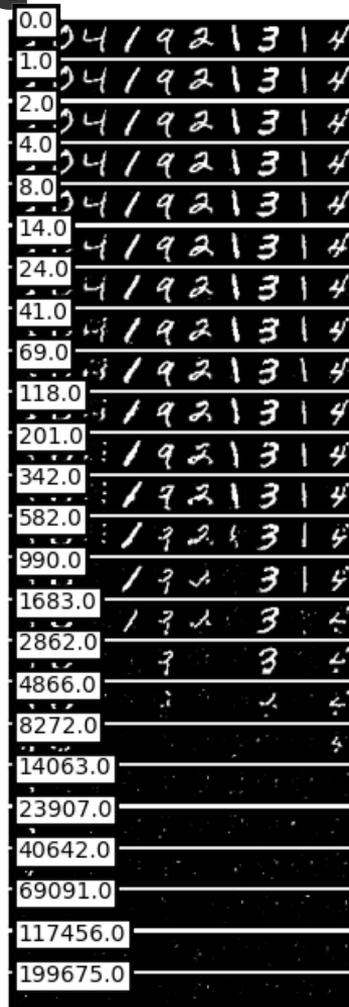


Non-equilibrium regime

Training using CD: chain initialisation
at the dataset

If we sample the RBM from:

- random **NOTHING**
- the dataset we **do not get anything new**



CD-100 random start



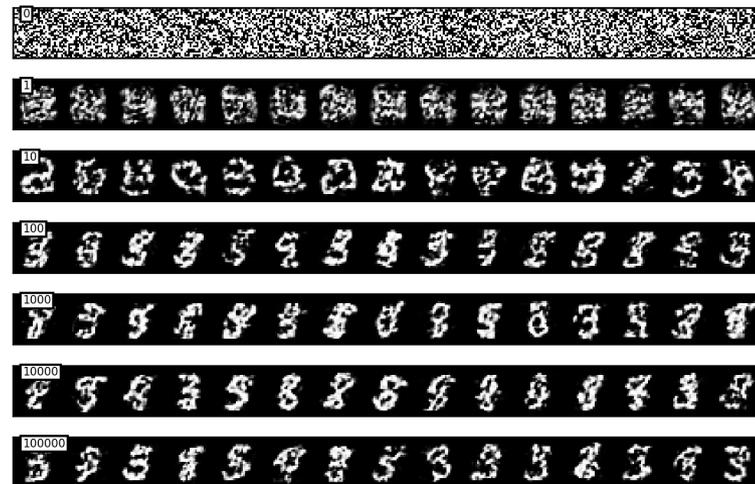
Non-equilibrium regime

Training using Mean Field dynamics

If we sample the RBM :

- Heat bath dynamics
NOTHING
- With MF **good data at t~k**
and it does not change
anymore

With HB MCMC dynamics

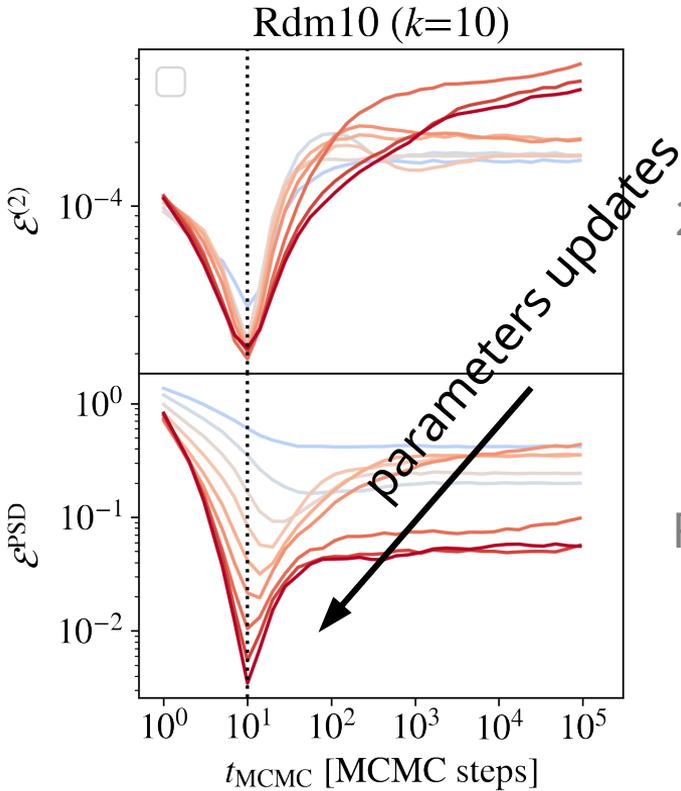


With MF dynamics



Non-equilibrium regime

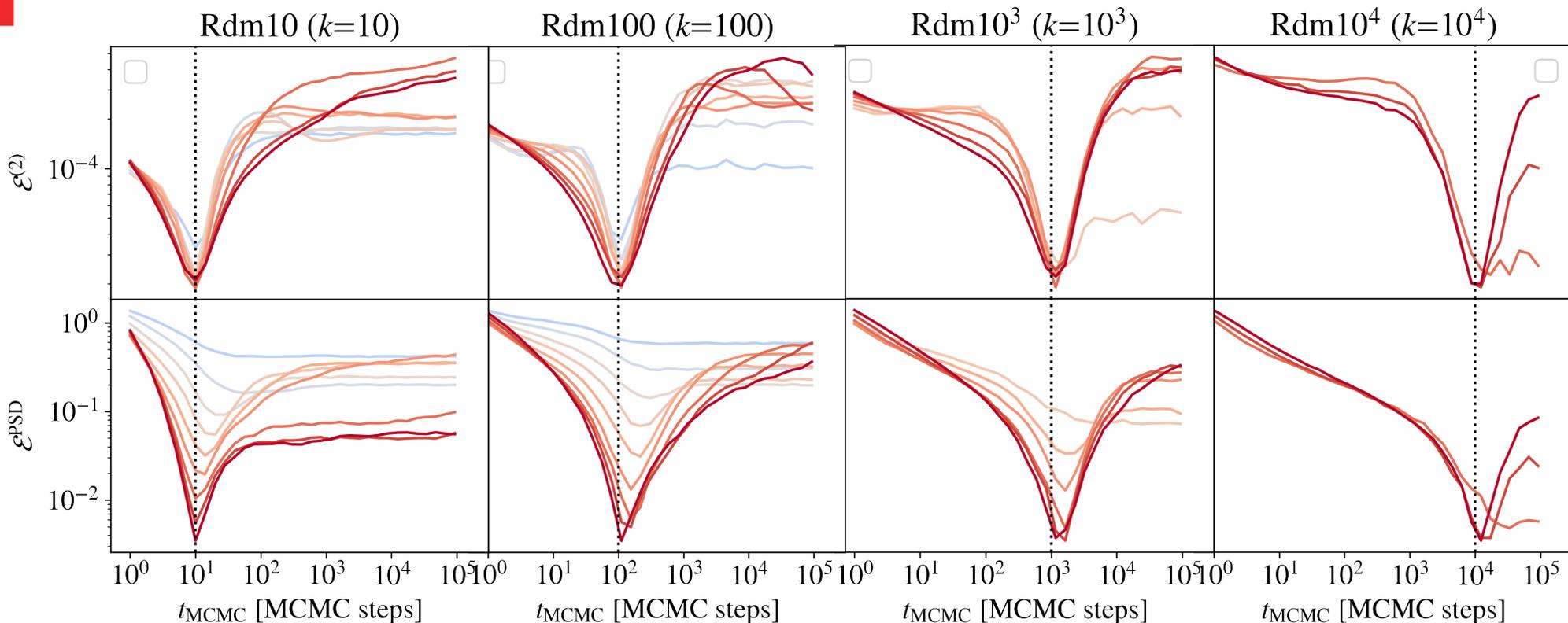
Quality of the generated samples



t_{age} [parameter updates]		
2863	14064	69092
4867	23908	117457
8273	40643	199676

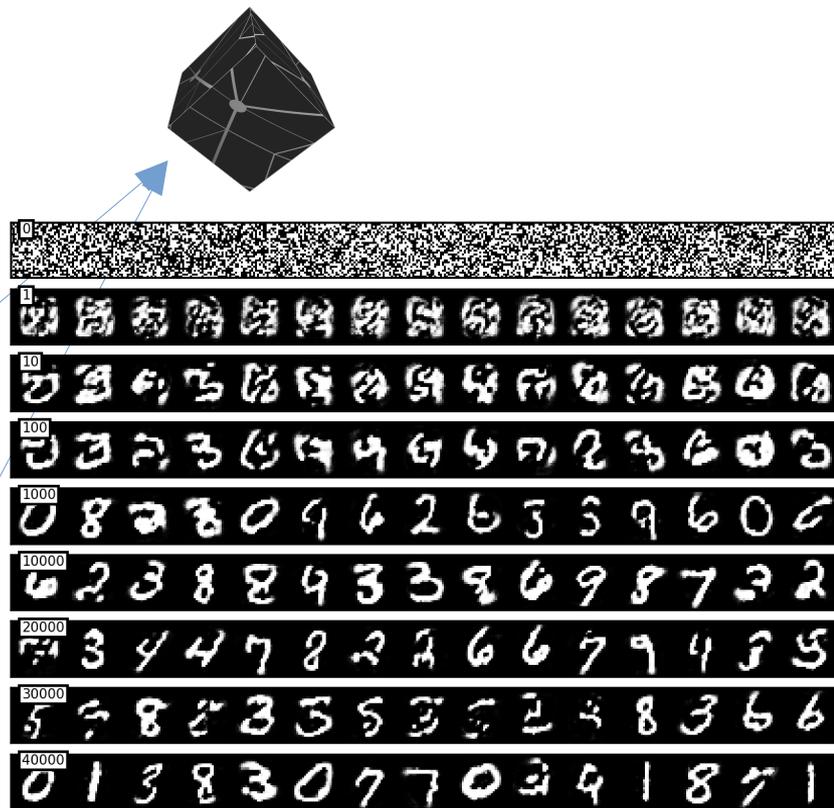
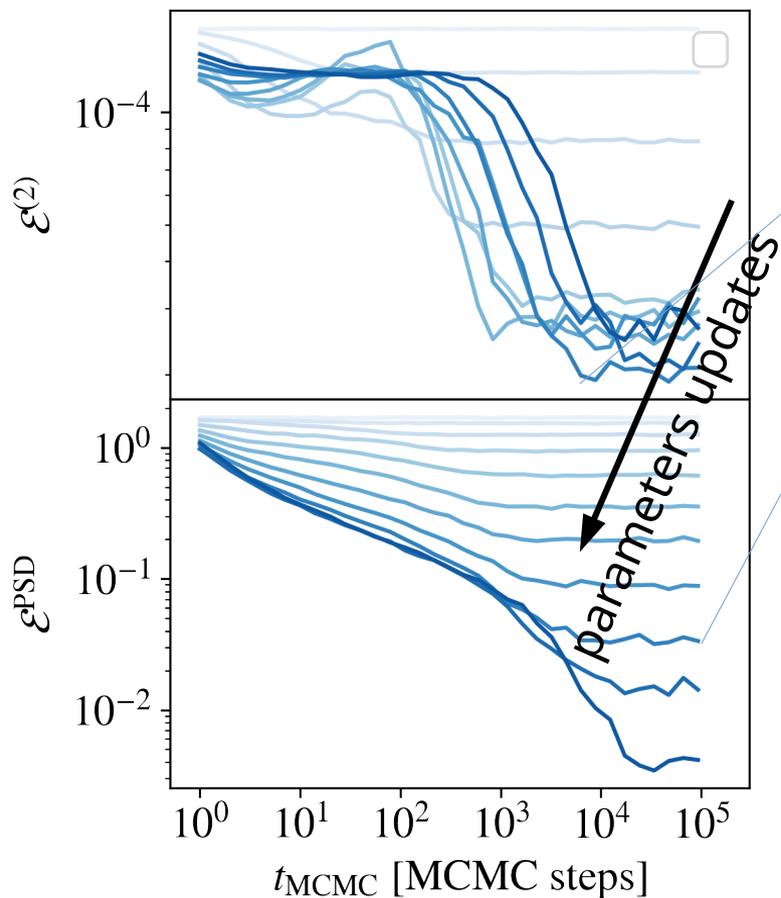
Non-equilibrium regime : generation

Quality of the generated samples



Best quality samples are obtained at $t_{\text{MCMC}} \sim k$

Equilibrium regime



Dynamics are much faster

Equilibrium vs. Non-eq. regimes

Non-equilibrium

$$k < t_{\text{therm}}$$

- “Learns the **dynamics**”
- **Advantage:** Optimal for data generation

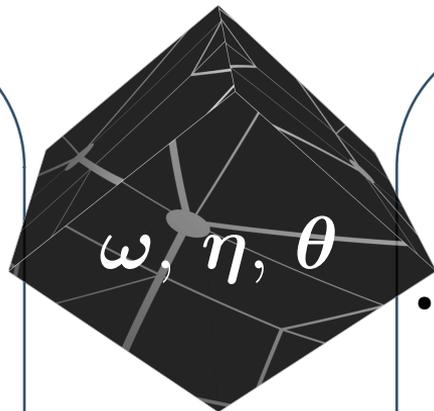
→ random noise initialization

Nijkamp, Hill, Han, Wu, Zhu.
NIPS 2019, AAI 2020.

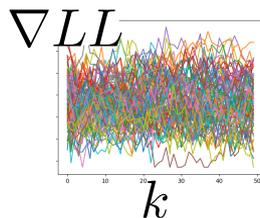


BM: Muntoni, Pagnani,
Martin Weigt, Zamponi (2021)

- **Drawbacks:**
 - Unpredictable if not controlled
 - not a good model for the data
 - Extremely slow dynamics



Training



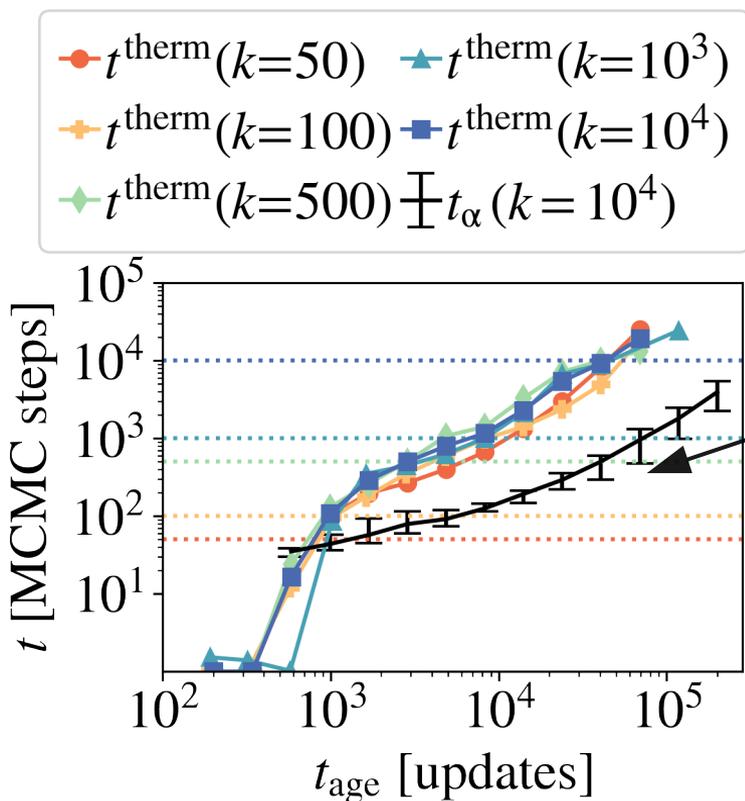
Equilibrium

$$k > t_{\text{therm}}$$

- Learns the (unnormalized) **prob. Distribution** of the data
- **Advantage:**
 - Fits a good model for the data
 - Sampling is stable
- **Drawbacks:**
 - Very slow training: **need large k**
 - Generating new configurations can become prohibitive

Equilibration: how long?

Easy : **MNIST**



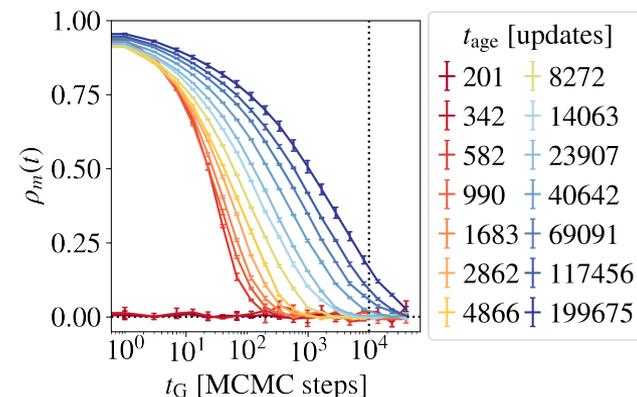
Mixing time

- Grows with the learning updates
- Always over 50

$t^{\text{therm}} \sim n\tau_{\text{mixing}}$

$$C_m^{\text{eq}}(t) = \frac{1}{N_v} \sum_i (m_i(t) - m)(m_i(0) - m)$$

$$\rho_m(t) = \frac{C_m(t)}{C_m(0)} \sim A \exp(-t/t_{\alpha})$$

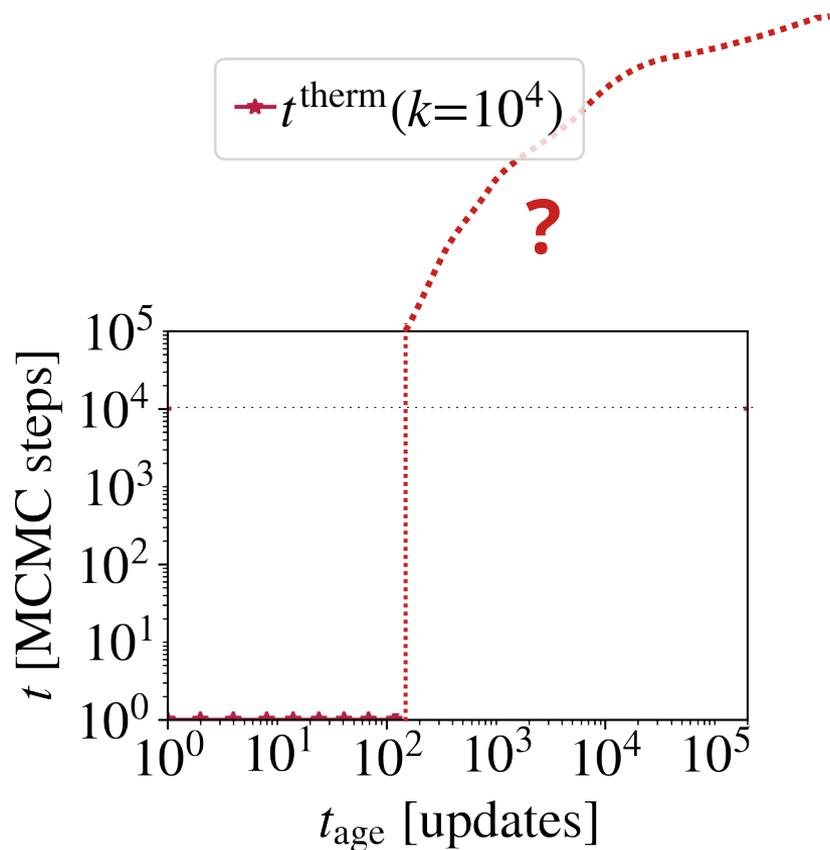


Literature $k \sim 10$

Always out-of-equilibrium regime !

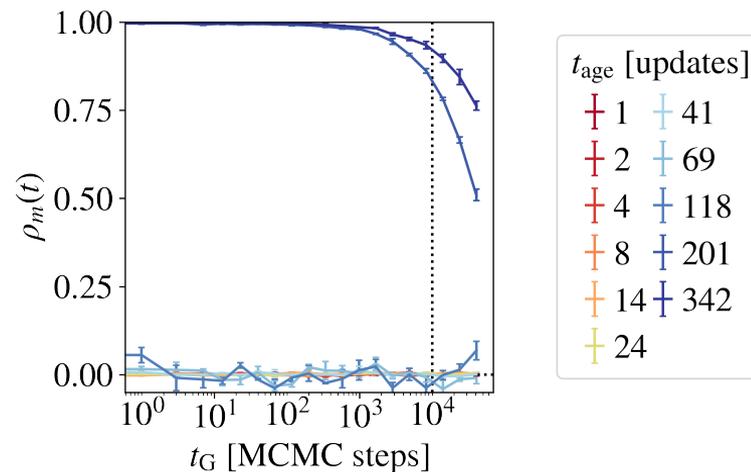
Equilibration: how long?

Hard : **GENE**



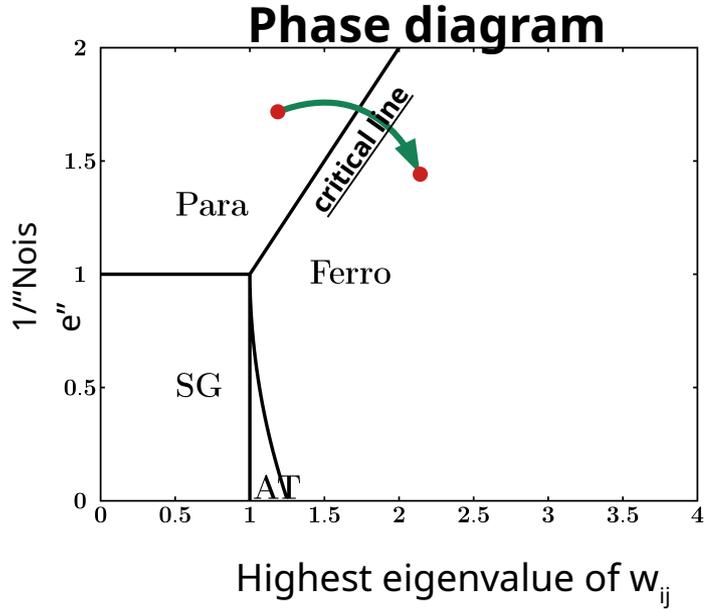
The thermalisation time jumps suddenly beyond 10^5 MCMC steps.

The equilibrium regime is beyond our reach...



What does it happen?

$$E[v, h] = - \sum_{ia} v_i \mathbf{w}_{ia} h_a - \sum_i \eta_i v_i - \sum_a \theta_a h_a$$



ω_i : i -th eigenvector of the W matrix

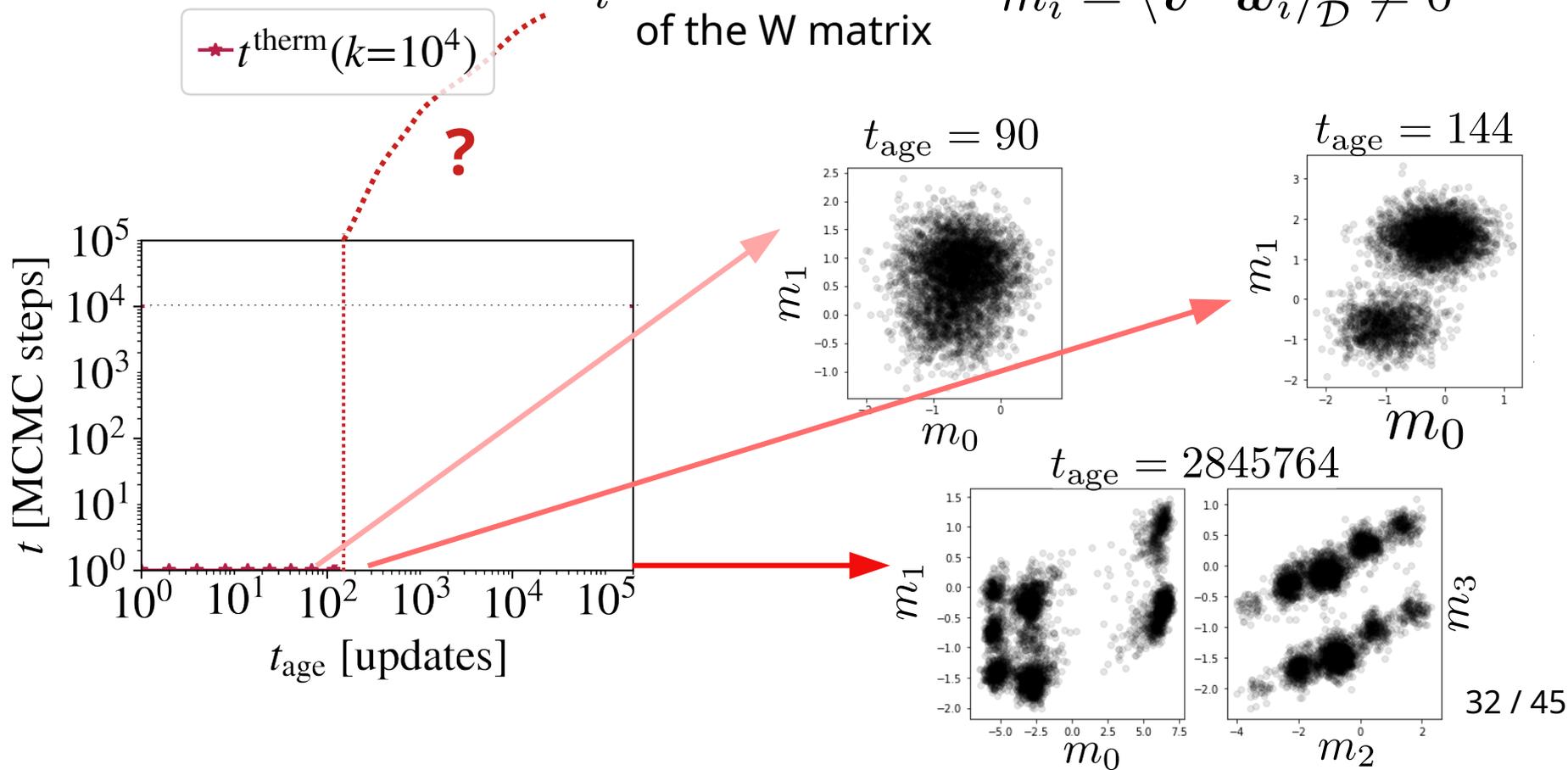
$$m_i = \langle \mathbf{v} \cdot \boldsymbol{\omega}_i \rangle_{\mathcal{D}} \neq 0$$

What does it happen?

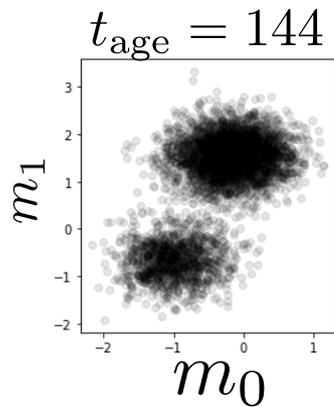
$$E[v, h] = - \sum_{ia} v_i \mathbf{w}_{ia} h_a - \sum_i \eta_i v_i - \sum_a \theta_a h_a$$

ω_i : i -th eigenvector of the W matrix

$$m_i = \langle \mathbf{v} \cdot \omega_i \rangle_{\mathcal{D}} \neq 0$$



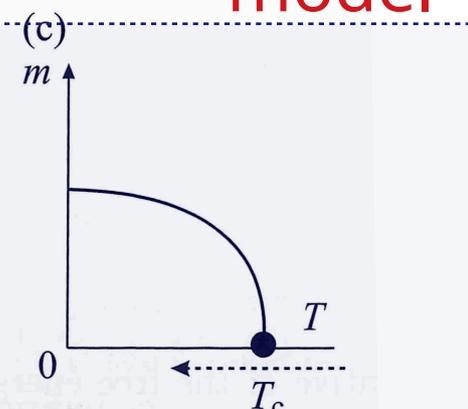
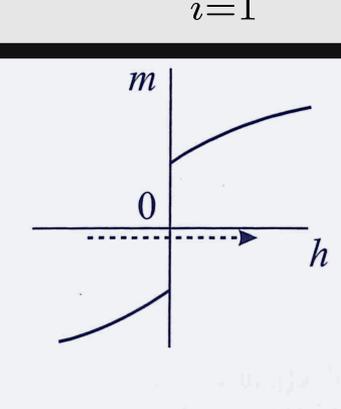
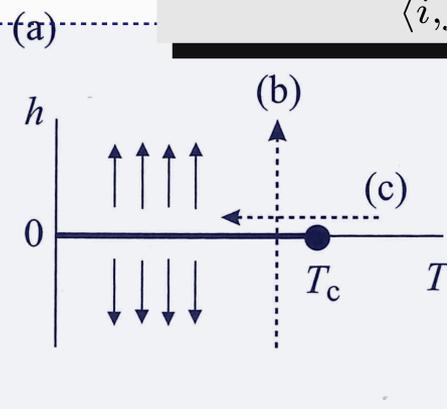
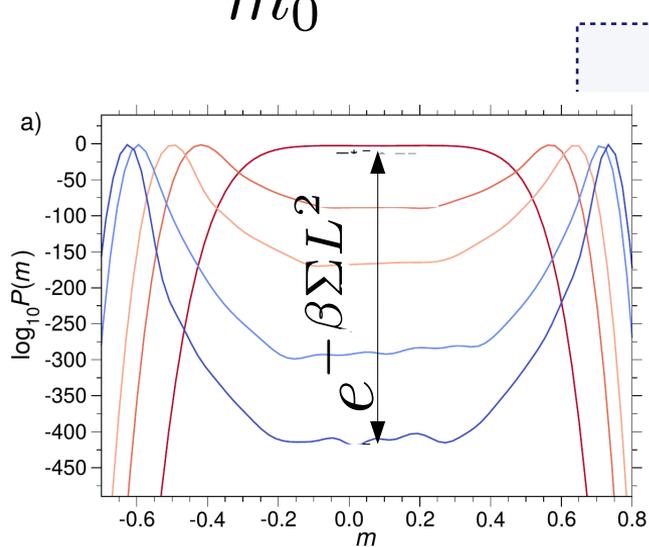
What does it happen?



As learning advances we start to have **metastable states**

$$\mathcal{H} = -J \sum_{\langle i,j \rangle} \sigma_i \sigma_j - h \sum_{i=1}^N \sigma_i,$$

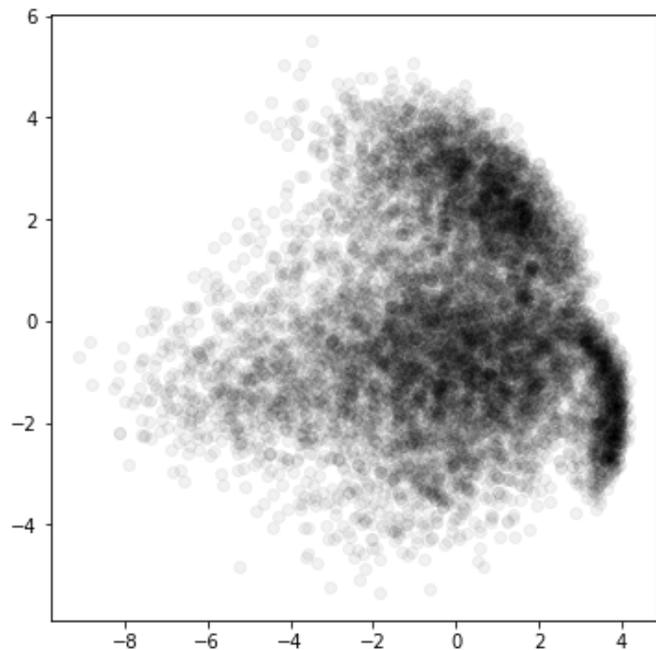
Ising
model



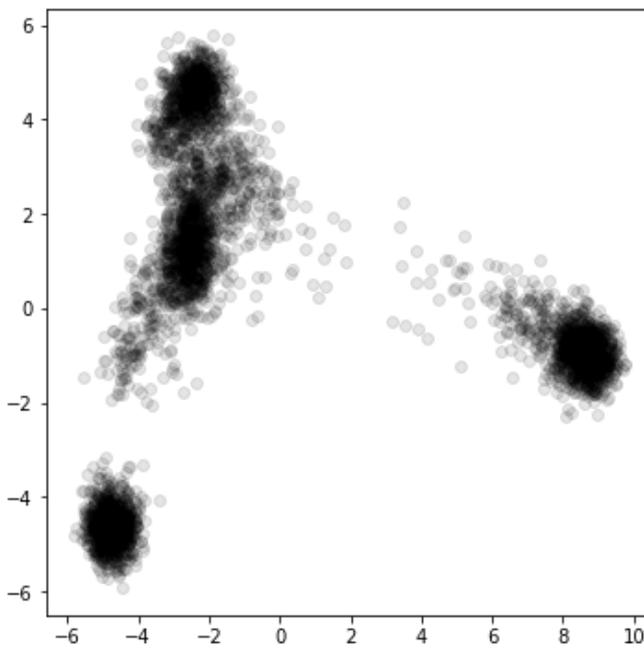
Structured datasets

We do the PCA of the data and project the data along the first 2 eigenvectors

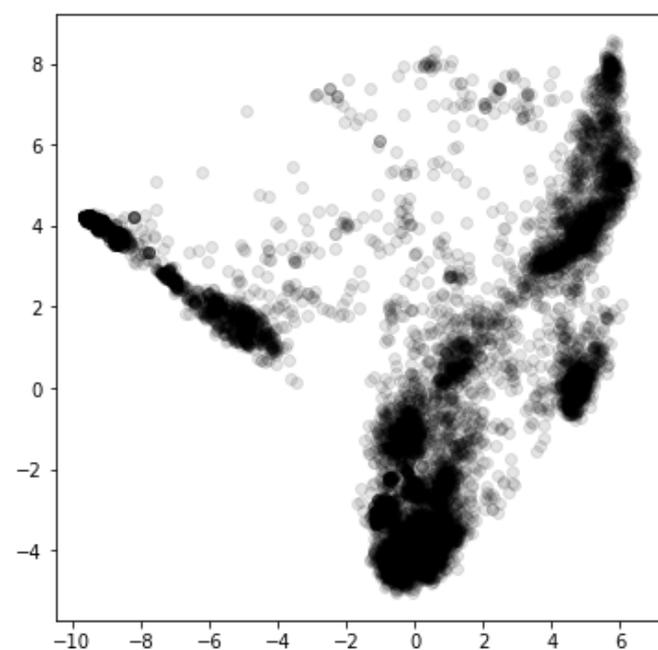
MNIST



GENE

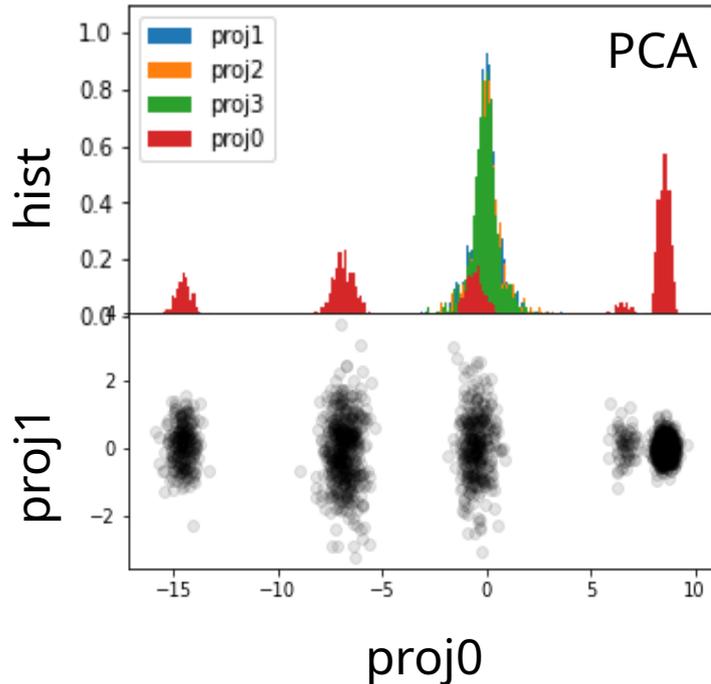


PROTEIN FAMILY



Step back : high dimensional clusters along 1D

We feed the RBM with points belonging to different clusters (in high dimensions) but Separated only in one

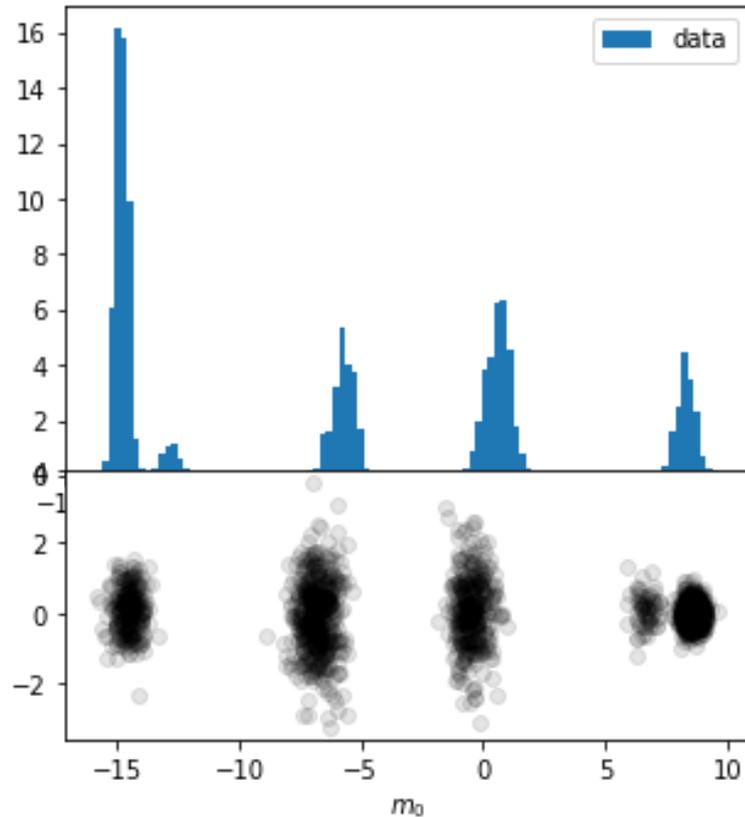


- Standard RBM training procedure fails completely to fit such dataset
- Yet, this simple low dimensional dataset can be trained analytically [Decelle, Furtlehner, PRL 2021]

→ We can have a perfect model ω, θ, η to test the biased sampling

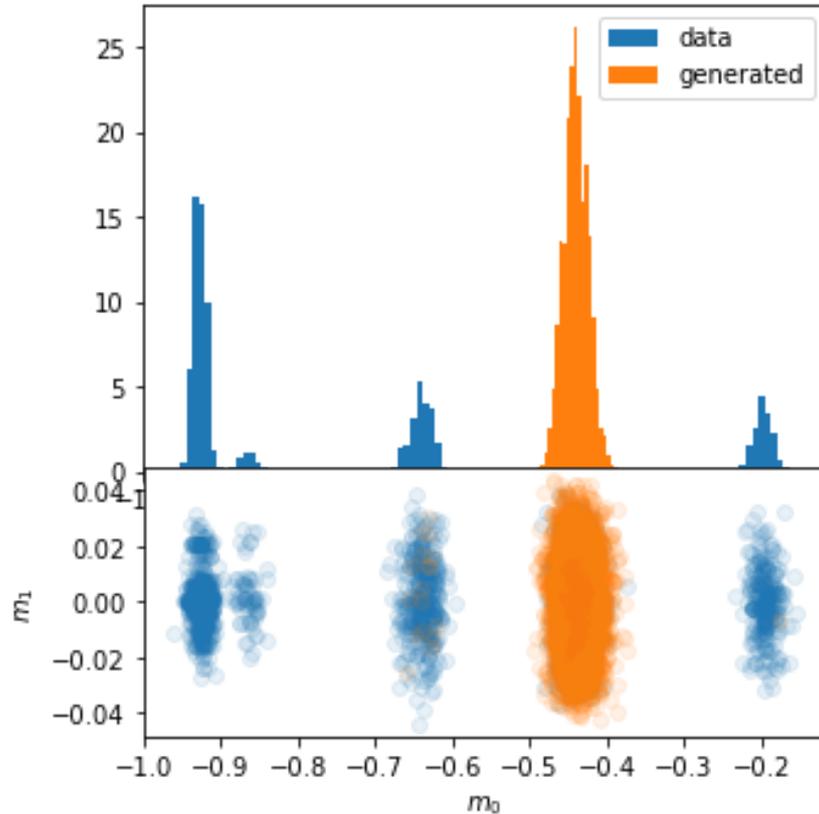
[Bereux, Decelle, Furtlehner, Seoane, *In preparation*]

Problems of the standard MCMC sampling



Projection of the first eigenvector
of the W matrix (normalized)

Problems of the standard MCMC sampling



Projection of the first eigenvector
of the W matrix (normalized)

The Tethered Monte Carlo approach (I)

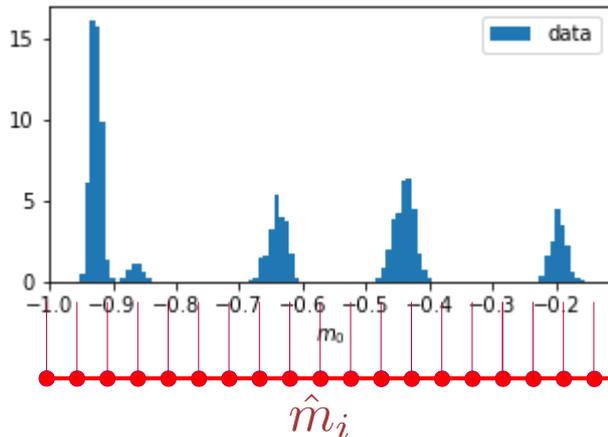
[Fernandez, Martin-Mayor, Yllanes - Nuclear physics (2009),
Martin-Mayor, Seoane, Yllanes, Journal of Statistical Physics (2011),
Fernández, Martín-Mayor, Seoane, Verrocchio, PRL (2012)]

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} = \sqrt{\frac{N}{2\pi}} \int_{-\infty}^{\infty} d\hat{m} \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} e^{-N(\hat{m} - m_0(\mathbf{v}))^2/2} = \int_{-\infty}^{\infty} d\hat{m} e^{-N\Omega(\hat{m})}$$

$$\langle O(\mathbf{v}, \mathbf{h}) \rangle = \frac{\sum_{\mathbf{v}, \mathbf{h}} O e^{-E(\mathbf{v}, \mathbf{h})}}{\sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}} = \int_{-\infty}^{\infty} d\hat{m} \langle O \rangle_{\hat{m}} e^{-N\Omega(\hat{m})}$$

$$\langle O(\mathbf{v}, \mathbf{h}) \rangle_{\hat{m}} = \frac{\sum_{\mathbf{v}, \mathbf{h}} O \omega(\hat{m}, \mathbf{v}, \mathbf{h})}{\sum_{\mathbf{v}, \mathbf{h}} \omega(\hat{m}, \mathbf{v}, \mathbf{h})}$$

$$\omega(\hat{m}, \mathbf{v}, \mathbf{h}) = e^{-E(\mathbf{v}, \mathbf{h})} e^{-N(\hat{m} - m_0(\hat{v}))^2/2}$$



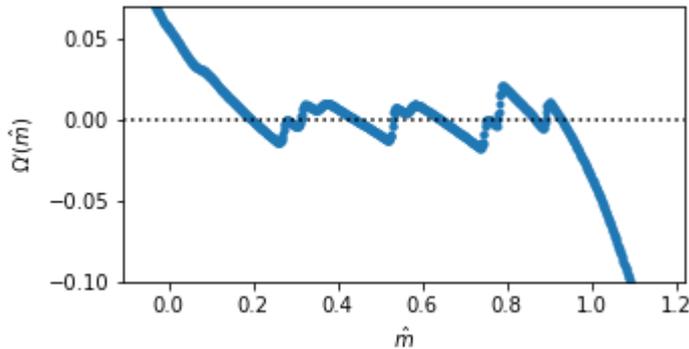
- Run K simulations at \hat{m}_i , with $i=1, \dots, K$ fixed
- We break the metastability: fast thermalisation

The Tethered Monte Carlo approach (II)

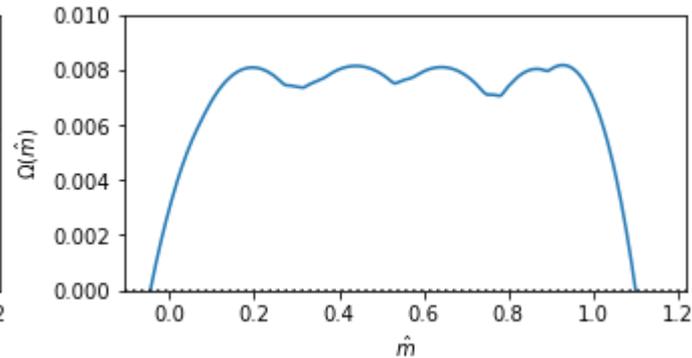
[Fernández, Martín-Mayor, Yllanes - Nuclear physics (2009),
Martín-Mayor, Seoane, Yllanes, Journal of Statistical Physics (2011),
Fernández, Martín-Mayor, Seoane, Verrocchio, PRL (2012)]

$$\frac{d\Omega}{d\hat{m}} = \langle \hat{m} - m_0(\mathbf{v}) \rangle_{\hat{m}}$$

1) Compute Ω'

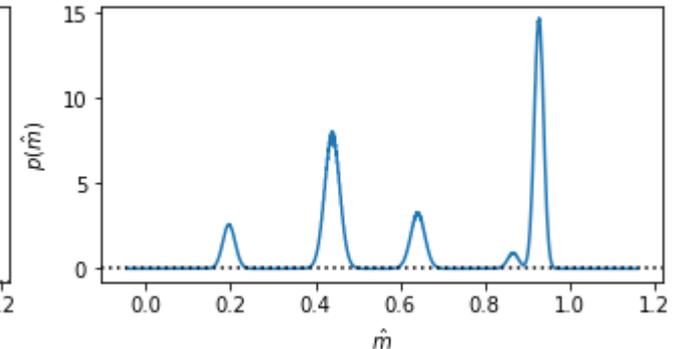


2) Integrate it numerically



3) Extract

$$p(\hat{m}) = \frac{e^{-N\Omega(\hat{m})}}{\int_{-\infty}^{\infty} e^{-N\Omega(\hat{m})}}$$

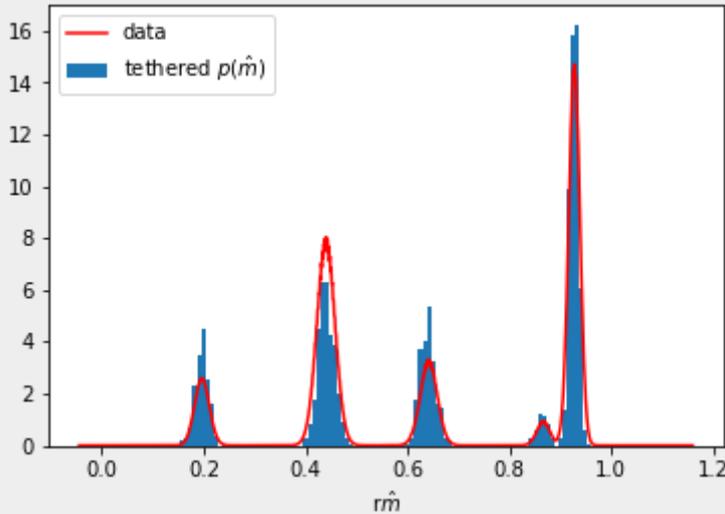
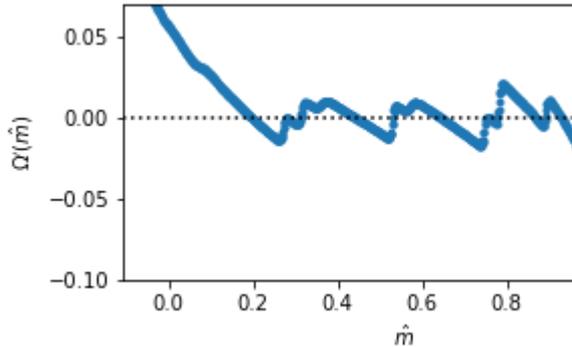


The Tethered Monte Carlo approach (II)

[Fernandez, Martin-Mayor, Yllanes - Nuclear physics (2009),
 Martin-Mayor, Seoane, Yllanes, Journal of Statistical Physics (2011),
 Fernández, Martín-Mayor, Seoane, Verrocchio, PRL (2012)]

$$\frac{d\Omega}{d\hat{m}} = \langle \hat{m} - n \rangle$$

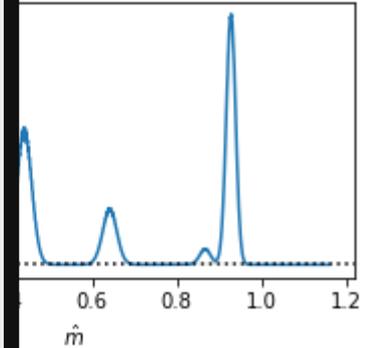
1) Compute Ω'



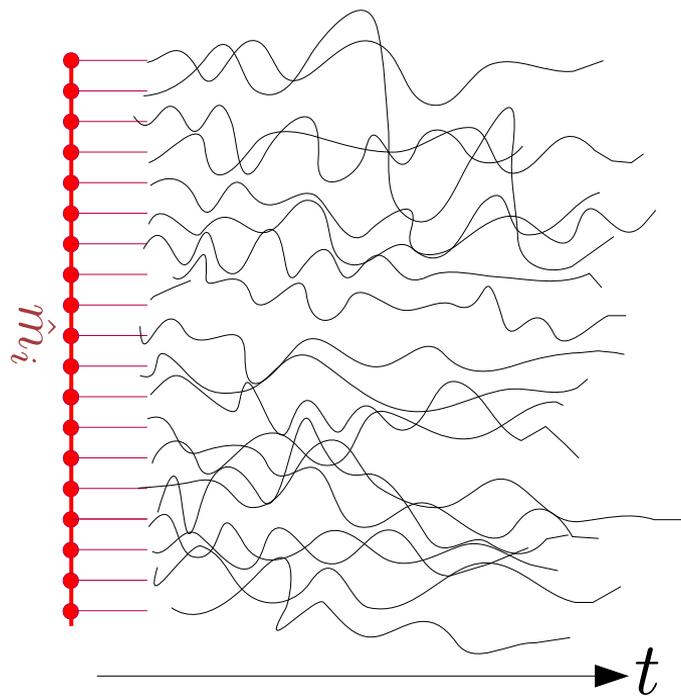
Once having $p(m)$, we can use **inverse sampling** for sample generation!

extract

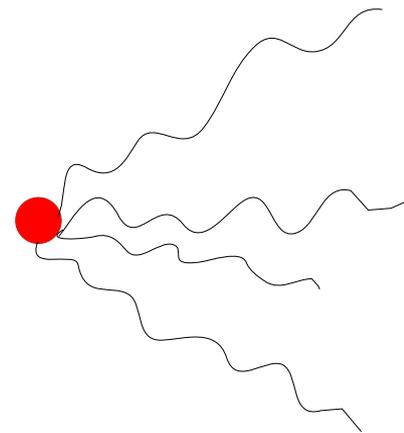
$$\frac{e^{-N\Omega(\hat{m})}}{\int_{-\infty}^{\infty} e^{-N\Omega(\hat{m})}$$



Learning with TMCMC



$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial w_{ia}} &= \langle v_i h_a \rangle_{\mathcal{D}} - \langle v_i h_a \rangle_{\mathcal{H}} \\
 \frac{\partial \mathcal{L}}{\partial \eta_i} &= \langle v_i \rangle_{\mathcal{D}} - \langle v_i \rangle_{\mathcal{H}} \\
 \frac{\partial \mathcal{L}}{\partial \theta_a} &= \langle h_a \rangle_{\mathcal{D}} - \langle h_a \rangle_{\mathcal{H}}
 \end{aligned}$$



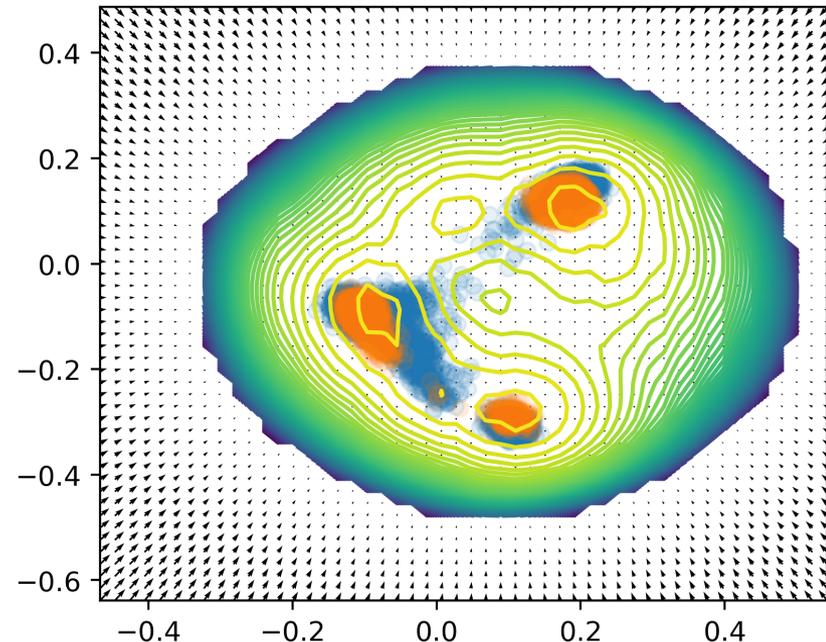
$$\langle O(v, h) \rangle_{\mathcal{H}} = \int_{-\infty}^{\infty} d\hat{m} \langle O \rangle_{\hat{m}} e^{-N\Omega(\hat{m})}$$

Learning with TMCMC

Generalization to higher number of conserved observables is straightforward...

$$\omega(\hat{m}_1, \hat{m}_2, \mathbf{v}, \mathbf{h}) = e^{-E(\mathbf{v}, \mathbf{h})} e^{-N(\hat{m}_2 - m_0(\hat{v}))^2 / 2} e^{-N(\hat{m}_2 - m_1(\hat{v}))^2 / 2}$$

$$\Omega'(\hat{m}) \rightarrow \nabla \Omega = \left(\langle \hat{m}_0 - m_0 \rangle_{\hat{m}_0, \hat{m}_1} \langle \hat{m}_1 - m_1 \rangle_{\hat{m}_0, \hat{m}_1} \right)$$



Conclusions

Decelle, Furtlehner, Seoane
[ArXiv:2105.13889](https://arxiv.org/abs/2105.13889)

- RBM have a major advantage in terms of interpretability of the extracted patterns, but training is very unstable following the standard recipes.
- **Instability** is a consequence of the **nonequilibrium** sampling during the sampling and can be controlled and taken in **advance to generate good samples with short trainings**.
- In order to fit a **good model for the data**, the sampling during the learning **must equilibrate**:
 - Datasets without structure : mixing time grows with Nb. Updates
 - Structured datasets: thermalisation is hampered by coexistence of states → biased sampling

Parameters MNIST

- Number of hidden nodes: $N_h = 500$
- Learning rate: $\alpha = 0.01$
- Minibatch size: $n_{mb} = 500$
- no ℓ_2 regularization of momentum.
- The gradient is centered according to [1]
- The visible biases are initialized to match the empirical frequency of the training dataset:

$$\eta_i = \log \left(\frac{\bar{m}_i}{1 - \bar{m}_i} \right) \text{ where } \bar{m}_i = \frac{1}{M} \sum_m s_i^{(m)} \quad (1)$$

- The number of MC chains used for the negative term was always equal to n_{mb}
- The number of MC steps for the negative chains is indicated by the variables t_{GL} and can vary.

RBM: learning and phase transition

We can confirm experimentally that the divergence of the mixing time correspond to the 2nd order phase transition

