

# The Statistical Mechanics of Associative Memories

Mathematics for Artificial Intelligence, Instituto Superior Técnico

Paulo Duarte Mourão

Department of Mathematics - Sapienza University of Rome

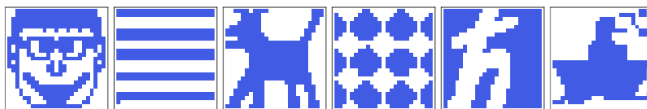
March 2026



# The quintessential example - the Hopfield model

The Hopfield model brings together different major scientific areas. It can be described as

- a recurrent neural network;
- a spin system;
- an associative memory.



**Figure 3.2:** Set of  $P = 6$  patterns stored in a Hopfield network of  $N = 625$  spins. Patterns are black and white images: the network is dealing with digital storage of information [Coolen et al. \(2005\)](#).

# Part I

# Outline

- 1 The Statistical Mechanics framework
- 2 Spin Systems
- 3 Recurrent Neural Networks
- 4 From magnets to memories

# Outline

- 1 The Statistical Mechanics framework
- 2 Spin Systems
- 3 Recurrent Neural Networks
- 4 From magnets to memories

# Motivation

Why do we need Statistical Mechanics in the first place?

- Classical equations are generally not easily treatable for a high number of interacting bodies  $N$  (the system with  $N = 3$  and gravitational interaction is already chaotic).

# Motivation

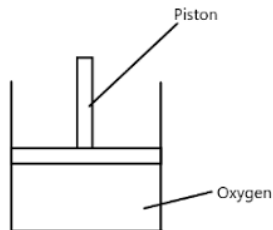
Why do we need Statistical Mechanics in the first place?

- Classical equations are generally not easily treatable for a high number of interacting bodies  $N$  (the system with  $N = 3$  and gravitational interaction is already chaotic).
- Statistical Mechanics proposes leveraging statistics to derive macroscopic behaviors from microscopic properties of very large systems. This is done by studying instead probability distributions over a given set of possible microstates.

# Thermodynamic systems in equilibrium

We will need:

- A set of possible microstates  $\Omega = \{\omega_i\}_{i \in I}$ .
- A function  $S: \Omega \rightarrow \mathbb{R}$  to be maximized (i.e. the entropy).
- A (set of) constraint(s)  $V: P(\Omega) \rightarrow \mathbb{R}$ .



# A natural choice of entropy

A natural choice of entropy is the so-called Shannon entropy

$$S = -k_B \sum_{i \in I} p_i \ln p_i \quad (1)$$

where  $p_i := P(\omega_i)$ , for each  $i \in I$ . We shall also take  $k_B = 1$ .

## Two possibilities for constraints

Let us consider an energy function

$$\begin{aligned}\Omega &\rightarrow \mathbb{R} \\ \omega_i &\mapsto E_i\end{aligned}$$

In the **microcanonical ensemble**, the energy  $E$  is kept fixed, and we have

$$p_i = \begin{cases} \frac{1}{N_i} & \text{if } E_i = E \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

with  $N_i$  denoting the number of states available at energy  $E_i$ .

## Two possibilities for constraints

Let us consider an energy function

$$\Omega \rightarrow \mathbb{R}$$

$$\omega_i \mapsto E_i$$

In the **canonical ensemble**, instead the average energy  $\bar{E} = \sum_i p_i E_i$  is kept fixed and we have

$$p_i = \frac{e^{-\beta E_i}}{Z_N(\beta)}, \quad Z_N(\beta) = \sum_{i \in I} e^{-\beta E_i}, \quad (3)$$

for some constant  $\beta$ . The above is known as the **Gibbs measure**, and  $Z$  as the **partition function**.

Lower values of the energy are more likely. How much depends on how low  $\beta$  is.

# Gibbs measure

The previous maximization problem may be restated (via Lagrange multipliers) as minimizing

$$F_N = \bar{E} - \beta^{-1}S \quad (4)$$

with  $\beta = 1/T$  acting as the multiplier.

## Gibbs measure

The previous maximization problem may be restated (via Lagrange multipliers) as minimizing

$$F_N = \bar{E} - \beta^{-1}S \quad (4)$$

with  $\beta = 1/T$  acting as the multiplier.

Indeed, adding the normalization constraint

$$\mathcal{L}[p] = \sum_i p_i E_i + \beta^{-1} \sum_i p_i \ln p_i - \lambda \sum_i p_i. \quad (5)$$

## Gibbs measure

The previous maximization problem may be restated (via Lagrange multipliers) as minimizing

$$F_N = \bar{E} - \beta^{-1}S \quad (4)$$

with  $\beta = 1/T$  acting as the multiplier.

Indeed, adding the normalization constraint

$$\mathcal{L}[p] = \sum_i p_i E_i + \beta^{-1} \sum_i p_i \ln p_i - \lambda \sum_i p_i. \quad (5)$$

Differentiating gives

$$\frac{\partial \mathcal{L}}{\partial p_i} = E_i + \beta^{-1}(\ln p_i + 1) - \lambda. \quad (6)$$

and so

$$\frac{\partial \mathcal{L}}{\partial p_i} = 0 \iff \ln p_i = -\beta E_i + \beta\lambda - 1 \iff p_i = e^{\beta\lambda - 1} e^{-\beta E_i} \quad (7)$$

# Gibbs measure

Furthermore, we have

$$\begin{aligned}F_N(\beta) &= \sum_i p_i (E_i + \beta^{-1} \ln p_i) \\&= \sum_i p_i (E_i - E_i - \beta^{-1} \ln Z_N(\beta)) \\&= -\beta^{-1} \ln Z_N(\beta).\end{aligned}$$

Which should look familiar from probabilities if we look carefully.

# The free energy is all you need

Both the partition function and the free energy encode the knowledge we have about the system. For instance

$$\beta \frac{\partial F_N(\beta)}{\partial \beta} = \frac{\sum_{i \in I} e^{-\beta E_i} E_i}{Z_N(\beta)} = \sum_i p_i E_i = \bar{E} \quad (8)$$

They are nothing other than the moment generating function and the cumulant generating function, respectively. Furthermore, the free energy density  $f_N(\beta) := F_N(\beta)/N$  should have a limit as  $N \rightarrow \infty$ .

$$f(\beta) := \lim_{N \rightarrow \infty} f_N(\beta) \quad (9)$$

# The strategy

If we can write the partition function, and consequently the free energy, constrained to an appropriate order parameter  $\mathbf{m}$

$$Z_N(\beta) = \int d\mathbf{m} Z_N^c(\mathbf{m}; \beta) = \int d\mathbf{m} e^{-Nf^c(\mathbf{m}; \beta) + \mathcal{O}(1)} \quad (10)$$

Then saddle-point approximation gives

$$f(\beta) = - \lim_{N \rightarrow \infty} \frac{\ln Z_N(\beta)}{\beta N} = \min_{\mathbf{m}} f^c(\mathbf{m}). \quad (11)$$

# Outline

- 1 The Statistical Mechanics framework
- 2 Spin Systems**
- 3 Recurrent Neural Networks
- 4 From magnets to memories

# A quick detour: spins

- Spin is an intrinsic form of angular momentum carried by elementary particles (although without actually "spinning" in the usual sense).
- Ampère's Law: an electric current creates a magnetic field around it.

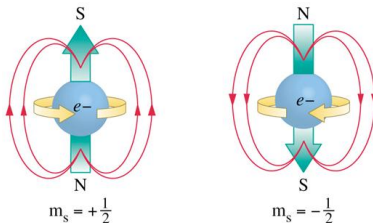


Figure: Magnetic fields resulting from electron spins.<sup>1</sup>

<sup>1</sup>Image from General Chemistry, 3rd edition, by Hill and Petrucci

# The Curie-Weiss model

We now consider a state space of  $N$  spins that can take the values  $\pm 1$ , meaning we have

$$\Omega = \{\pm 1\}^N \quad (12)$$

The Curie-Weiss model is a **mean-field approximation** model defined by the Hamiltonian

$$H_N^{\text{CW}}(\boldsymbol{\sigma}) = -\frac{1}{2N} \sum_{i,j=1}^N \sigma_i \sigma_j - h \sum_{i=1}^N \sigma_i \quad (13)$$

for each state  $\boldsymbol{\sigma} \in \Omega$ , and  $h \in \mathbb{R}$  the intensity of the external field.

# The Curie-Weiss model

Starting from the partition function

$$Z_N^{\text{CW}}(\beta) = \sum_{\sigma \in \{-1, +1\}^N} e^{-\beta H_N^{\text{CW}}(\sigma)}, \quad (14)$$

we insert the order parameter

$$m_N(\sigma) = \frac{1}{N} \sum_{i=1}^N \sigma_i \quad (15)$$

via the Fourier representation of the delta distribution

$$\delta \left( m - \frac{1}{N} \sum_{i=1}^N \sigma_i \right) = \frac{N}{2\pi} \int_{-\infty}^{+\infty} d\hat{m} \exp \left[ i\hat{m} \left( m - \frac{1}{N} \sum_{i=1}^N \sigma_i \right) \right]. \quad (16)$$

# The Curie-Weiss model

This yields

$$Z_N^{\text{CW}}(\beta) = \sum_{\sigma} \int dm \delta \left( m - \frac{1}{N} \sum_{i=1}^N \sigma_i \right) \exp \left[ \frac{N}{2} \beta m^2 + N \beta h m \right]$$

# The Curie-Weiss model

This yields

$$\begin{aligned}
 Z_N^{\text{CW}}(\beta) &= \sum_{\sigma} \int dm \delta \left( m - \frac{1}{N} \sum_{i=1}^N \sigma_i \right) \exp \left[ \frac{N}{2} \beta m^2 + N \beta h m \right] \\
 &= \frac{N}{2\pi} \sum_{\sigma} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m \right) - i \hat{m} \sum_{i=1}^N \sigma_i \right]
 \end{aligned}$$

# The Curie-Weiss model

This yields

$$\begin{aligned}
 Z_N^{\text{CW}}(\beta) &= \sum_{\sigma} \int dm \delta \left( m - \frac{1}{N} \sum_{i=1}^N \sigma_i \right) \exp \left[ \frac{N}{2} \beta m^2 + N \beta h m \right] \\
 &= \frac{N}{2\pi} \sum_{\sigma} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m \right) - i \hat{m} \sum_{i=1}^N \sigma_i \right] \\
 &= \frac{N}{2\pi} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m \right) \right] \sum_{\sigma} \prod_{i=1}^N e^{i \hat{m} \sigma_i}
 \end{aligned}$$

# The Curie-Weiss model

This yields

$$\begin{aligned}
 Z_N^{\text{CW}}(\beta) &= \sum_{\sigma} \int dm \delta \left( m - \frac{1}{N} \sum_{i=1}^N \sigma_i \right) \exp \left[ \frac{N}{2} \beta m^2 + N \beta h m \right] \\
 &= \frac{N}{2\pi} \sum_{\sigma} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m \right) - i \hat{m} \sum_{i=1}^N \sigma_i \right] \\
 &= \frac{N}{2\pi} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m \right) \right] \sum_{\sigma} \prod_{i=1}^N e^{i \hat{m} \sigma_i} \\
 &= \frac{N}{2\pi} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m \right) \right] \prod_{i=1}^N \sum_{\sigma_i = \pm 1} e^{i \hat{m} \sigma_i}
 \end{aligned}$$

# The Curie-Weiss model

This yields

$$\begin{aligned}
 Z_N^{\text{CW}}(\beta) &= \sum_{\sigma} \int dm \delta \left( m - \frac{1}{N} \sum_{i=1}^N \sigma_i \right) \exp \left[ \frac{N}{2} \beta m^2 + N \beta h m \right] \\
 &= \frac{N}{2\pi} \sum_{\sigma} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m \right) - i \hat{m} \sum_{i=1}^N \sigma_i \right] \\
 &= \frac{N}{2\pi} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m \right) \right] \sum_{\sigma} \prod_{i=1}^N e^{i \hat{m} \sigma_i} \\
 &= \frac{N}{2\pi} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m \right) \right] \prod_{i=1}^N \sum_{\sigma_i = \pm 1} e^{i \hat{m} \sigma_i} \\
 &= \frac{N}{2\pi} \int dmd\hat{m} \exp \left[ N \left( \frac{\beta}{2} m^2 + \beta h m + i \hat{m} m + \ln 2 \cos \hat{m} \right) \right].
 \end{aligned}$$

# The Curie-Weiss model

Therefore, we get

$$\begin{aligned} f^{\text{CW}}(\beta) &= - \lim_{N \rightarrow \infty} \frac{1}{\beta N} \ln Z_N^{\text{CW}}(\beta) \\ &= \min_m \left[ -\frac{m^2}{2} - hm - i\beta^{-1} \hat{m} m - \beta^{-1} \ln 2 \cos \hat{m} \right] \end{aligned}$$

# The Curie-Weiss model

Therefore, we get

$$\begin{aligned}
 f^{\text{CW}}(\beta) &= - \lim_{N \rightarrow \infty} \frac{1}{\beta N} \ln Z_N^{\text{CW}}(\beta) \\
 &= \min_m \left[ -\frac{m^2}{2} - hm - i\beta^{-1} \hat{m} m - \beta^{-1} \ln 2 \cos \hat{m} \right]
 \end{aligned}$$

The saddle-point equations are then

$$\left. \frac{\partial f^{\text{CW}}(\beta)}{\partial m} \right|_{\substack{m=m^* \\ \hat{m}=\hat{m}^*}} = 0 \iff m^* + h = -i\beta^{-1} \hat{m}^* \quad (17)$$

$$\left. \frac{\partial f^{\text{CW}}(\beta)}{\partial \hat{m}} \right|_{\substack{m=m^* \\ \hat{m}=\hat{m}^*}} = 0 \iff m^* = -i \tan \hat{m}^* \quad (18)$$

# The Curie-Weiss model

Therefore, we get

$$\begin{aligned}
 f^{\text{CW}}(\beta) &= - \lim_{N \rightarrow \infty} \frac{1}{\beta N} \ln Z_N^{\text{CW}}(\beta) \\
 &= \min_m \left[ -\frac{m^2}{2} - hm - i\beta^{-1} \hat{m} m - \beta^{-1} \ln 2 \cos \hat{m} \right]
 \end{aligned}$$

The saddle-point equations are then

$$\left. \frac{\partial f^{\text{CW}}(\beta)}{\partial m} \right|_{\substack{m=m^* \\ \hat{m}=\hat{m}^*}} = 0 \iff m^* + h = -i\beta^{-1} \hat{m}^* \quad (17)$$

$$\left. \frac{\partial f^{\text{CW}}(\beta)}{\partial \hat{m}} \right|_{\substack{m=m^* \\ \hat{m}=\hat{m}^*}} = 0 \iff m^* = -i \tan \hat{m}^* \quad (18)$$

Combining both gives

$$m^* = \tanh(\beta(m^* + h)) \quad (19)$$

# The Curie-Weiss model

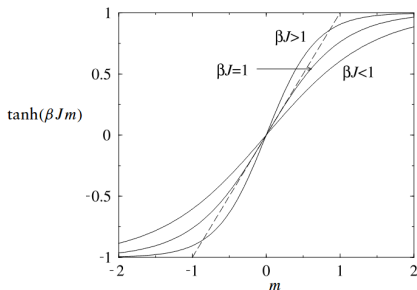
Given the simplicity of the saddle point equation in  $m$ , it can also be used to replace  $\hat{m} = i\beta(m + h)$  into constrained the free energy, yielding the sometimes called pseudo (constrained) free energy

$$\tilde{f}^{\text{CW},c}(m; \beta) = \frac{m^2}{2} - \beta^{-1} \ln 2 \cosh \beta (m + h) \quad (20)$$

which has the same extrema as the true constrained free energy.

# The Curie-Weiss model

The next picture[6] shows clearly possible solutions for the minima of the free energy when  $h = 0$  ( $\beta = 1$  is a phase transition!)



**Figure 20.2** Graphical solution of the Curie-Weiss equation (20.43), viz.  $m = \tanh(\beta J m)$ . Solid lines show the function  $\tanh(\beta J m)$  for different values of  $\beta J$ . Dashed: the diagonal which this function must intersect. For  $\beta J < 1$ ,  $m = 0$  is the only solution; for larger values of  $\beta J$ , there are two additional solutions  $m_\beta$  and  $-m_\beta$ .

# Outline

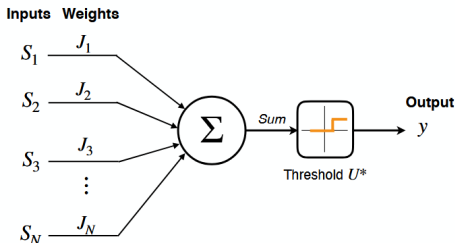
- 1 The Statistical Mechanics framework
- 2 Spin Systems
- 3 Recurrent Neural Networks**
- 4 From magnets to memories

# The McCulloch-Pitts neuron

The MP neuron is a model of an artificial neuron that takes in binary inputs  $\{S_i = 0, 1\}_{i=1}^N$  and outputs a binary response  $y = 0, 1$  given by

$$y = \theta \left( \sum_{i=1}^N J_i S_i - U^* \right) \quad (21)$$

with  $\theta$  being the Heaviside step-function.



**Figure 1.9:** Schematic representation of a MP neuron.

# Recurrent Neural Networks (RNNs)

Fully-connected recurrent neural networks are described by a family of  $N$  MP neurons  $\{S_i = 0, 1\}_{i=1}^N$  and an update rule given by

$$S_i(t+1) = \theta \left( \sum_{j=1}^N J_{ij} S_j(t) - U_i^*(t) \right) \quad (22)$$

And to simulate the effect of noise in the system, we consider the  $U_i^*$  to be a random variable

$$U_i^*(t) = U_i - \frac{1}{2} T z_i(t), \quad (23)$$

with  $\mathbb{E}[z_i(t)] = 0$ ,  $\mathbb{E}[z_i(t)^2] = 1$  and  $T > 0$ .

# From RNNs to Spin Systems

Let us apply the transformation

$$\sigma_i(t) = 2S_i(t) - 1, \quad (24)$$

then the update rule becomes (notice that  $\text{sgn}(x) = 2\theta(x) - 1$ )

$$\sigma_i(t+1) = \text{sgn} \left( \sum_{j=1}^N J_{ij} \sigma_j(t) + h_i + Tz_i(t) \right), \quad (25)$$

with  $h_i = \sum_j J_{ij} - 2U_i$ . But how far does the analogy *really* go?

# Emergence of the Gibbs measure

## Theorem 1

*There exists a choice of noises  $z_i$  such that, for any symmetric interaction ( $J_{ij} = J_{ji}$ ) without self-interactions ( $J_{ii} = 0$ ), the dynamics (25) updated sequentially has the equilibrium distribution given by*

$$p(\boldsymbol{\sigma}) = Z^{-1} e^{-\frac{H(\boldsymbol{\sigma})}{T}} \quad (26)$$

$$Z = \sum_{\boldsymbol{\sigma} \in \{\pm 1\}^N} e^{-\frac{H(\boldsymbol{\sigma})}{T}} \quad (27)$$

with

$$H(\boldsymbol{\sigma}) = -\frac{1}{2} \sum_{i,j=1}^N J_{ij} \sigma_i \sigma_j - \sum_{i=1}^N h_i \sigma_i. \quad (28)$$

## Emergence of the Gibbs measure

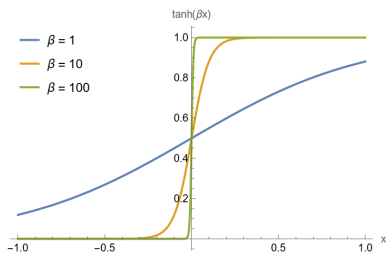
If the  $z_i(t)$ 's are defined according probability distribution function given by

$$f_{z_i(t)}(x) = 1 - \tanh^2(x). \quad (29)$$

It can be shown that this is equivalent to writing the dynamics (25) as

$$\mathbb{P} [\sigma_i(t+1) = \pm \sigma_i(t)] = \frac{1}{2} \left[ 1 \pm \tanh \left( \frac{E_i(\boldsymbol{\sigma}(t))}{T} \right) \right] \quad (30)$$

with  $E_i(\boldsymbol{\sigma}) = \sum_j \sigma_i(t) J_{ij} \sigma_j(t) + h_i \sigma_i(t)$ .



# Emergence of the Gibbs measure

The update rule (30) can also be written as

$$\mathbb{P}[\sigma_i(t+1) \neq \sigma_i(t)] = \frac{1}{1 + e^{\frac{\Delta E_i}{T}}} \quad (31)$$

which corresponds to the Glauber Monte-Carlo dynamics for the canonical ensemble defined by the Hamiltonian

$$H(\sigma) = \sum_{ij} \sigma_i(t) J_{ij} \sigma_j(t) + \sum_i h_i \sigma_i(t), \quad (32)$$

where  $\Delta E_i$  is the energy difference resulting from swapping spin  $i$ .

# Emergence of the Gibbs measure

And so for this choice of noises or system converges to the Boltzmann-Gibbs distribution

$$p(\mathbf{J}, \boldsymbol{\sigma}) = \frac{e^{-\beta H_{N,\mathbf{h},H}(\mathbf{J},\boldsymbol{\sigma})}}{Z_{N,\mathbf{h},H}(\beta, \mathbf{J})} \quad (33)$$

in equilibrium with

$$Z_{N,\mathbf{h},H}(\beta, \mathbf{J}) = \sum_{\boldsymbol{\sigma} \in \{\pm 1\}^N} e^{-\beta H_{N,\mathbf{h},H}(\mathbf{J},\boldsymbol{\sigma})}. \quad (34)$$

# Outline

- 1 The Statistical Mechanics framework
- 2 Spin Systems
- 3 Recurrent Neural Networks
- 4 From magnets to memories**

# Back to the Curie-Weiss model

The Curie-Weiss model is a **mean-field approximation** model defined, with no external field, by the Hamiltonian

$$H_N^{\text{CW}}(\boldsymbol{\sigma}) = -\frac{1}{2N} \sum_{i,j=1}^N \sigma_i \sigma_j \quad (35)$$

# Generalizing Curie-Weiss

One way to think of the Hamiltonian for the Curie-Weiss model (35) is that it "incentivizes" the spins to align themselves with the configuration  $\eta = (1, \dots, 1)$ . In fact, the overall magnetization  $m_N$  may be written as

$$m_N(\sigma; \eta) = \frac{1}{N} \sum_{i=1}^N \eta_i \sigma_i \quad (36)$$

# Generalizing Curie-Weiss

What if we picked a different configuration, say  $\xi \in \{-1, 1\}^N$ ? Then we could define a new "generalized" Curie-Weiss Hamiltonian

$$H_N^{\text{CW}}(\sigma, \xi) = -\frac{N}{2} m_N^2(\sigma; \xi) - h N m_N(\sigma; \xi), \quad (37)$$

or, expanding,

$$H_N^{\text{CW}}(\sigma; \xi) = -\frac{1}{2N} \sum_{i,j=1}^N \xi_i \xi_j \sigma_i \sigma_j - h \sum_{i=1}^N \xi_i \sigma_i, \quad (38)$$

# Generalizing Curie-Weiss

We can see that this "generalization" of the Curie-Weiss model is actually equivalent, if we do a very simple change of variables  $\tilde{\sigma}_i = \xi_i \sigma_i$ . Then

$$\begin{aligned} H_N^{\text{CW}}(\boldsymbol{\sigma}; \boldsymbol{\xi}) &= -\frac{1}{2N} \sum_{i,j=1}^N \xi_i \xi_j \sigma_i \sigma_j - h \sum_{i=1}^N \xi_i \sigma_i \\ &= -\frac{1}{N} \sum_{i,j=1}^N \tilde{\sigma}_i \tilde{\sigma}_j - h \sum_{i=1}^N \tilde{\sigma}_i \end{aligned}$$

Since the  $\tilde{\sigma}$ 's also live in  $\{-1, 1\}^N$ , the systems are equivalent. The only difference is the spins now align themselves with  $\boldsymbol{\xi}$ .

# From Curie-Weiss to Hopfield

Now what if we add more of these patterns? For instance, we can consider  $K$  patterns  $\xi^\mu$ , with  $1 \leq \mu \leq K$  and write

$$H_{N,K}(\sigma; \xi) = -\frac{N}{2} \sum_{\mu=1}^K m_N^2(\sigma; \xi^\mu) - hNm_N(\sigma; \xi^1) \quad (39)$$

or, once again expanded,

$$H_{N,K}(\sigma; \xi) = -\frac{1}{2N} \sum_{\mu=1}^K \sum_{i,j=1}^N \xi_i^\mu \xi_j^\mu \sigma_i \sigma_j - h \sum_{i=1}^N \xi_i^1 \sigma_i \quad (40)$$

This is the Hamiltonian of the Hopfield model! [1, 2]

## Part II

# Outline

- 5 Disordered systems
- 6 The Hopfield model
- 7 Generalizations
- 8 Coming next
- 9 Sources and references

# Outline

- 5 Disordered systems
- 6 The Hopfield model
- 7 Generalizations
- 8 Coming next
- 9 Sources and references

# From Curie-Weiss to Hopfield

Now what if we add more of these patterns? For instance, we can consider  $K$  patterns  $\xi^\mu$ , with  $1 \leq \mu \leq K$  and write

$$H_{N,K}(\sigma; \xi) = -\frac{N}{2} \sum_{\mu=1}^K m_N^2(\sigma; \xi^\mu) - hNm_N(\sigma; \xi^1) \quad (41)$$

or, once again expanded,

$$H_{N,K}(\sigma; \xi) = -\frac{1}{2N} \sum_{\mu=1}^K \sum_{i,j=1}^N \xi_i^\mu \xi_j^\mu \sigma_i \sigma_j - h \sum_{i=1}^N \xi_i^1 \sigma_i \quad (42)$$

This is the Hamiltonian of the Hopfield model! [1, 2]

# Why is this harder?

- There are two main reasons why the Hopfield model is generally (at least for big enough  $K$ ) more complicated than the Curie-Weiss model: disorder and frustration.
- These two properties are characteristic of a wide array of models, known as **disordered systems**.

# Frustration

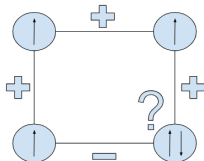
Looking at the interaction matrix for the Hopfield Model

$$J_{ij}^H = \frac{1}{N} \sum_{\mu=1}^K \xi_i^{\mu} \xi_j^{\mu} \quad (43)$$

we see that there is a non-zero probability of having connections that do not agree with a certain pattern.

# Frustration

- The matrix  $\mathbf{J}$  tells us whether the system is trying to align ( $J_{ij} > 0$ ) or oppose ( $J_{ij} < 0$ ) the spins  $i$  and  $j$ .
- Crucially, in the Hopfield Model, contrary to the Curie-Weiss model, not all interactions can be satisfied.



**Figure 2.9: A very simple example of a frustrated system.** The spins tend to be parallel when they interact with a positive coupling and anti-parallel when the interaction is negative. Obviously, not all the conditions can be met simultaneously, meaning that interaction is frustrated.

# Disorder

The partition function of the Hopfield model is given by

$$Z_{N,K}(\beta; \xi) = \sum_{\sigma \in \{\pm 1\}^N} e^{-\beta H_{N,K}(\sigma; \xi)} \quad (44)$$

which depends on the choice of patterns, and so does the free energy.

Calculating these quantities would generally be very difficult....

# Disorder

The solution is then to consider the  $\xi$ 's as being drawn with respect to a certain probability distribution, known as the **disorder**, for instance

$$\mathbb{P}(\xi_i^\mu = \pm 1) = \frac{1}{2} \quad (45)$$

and compute the so called **quenched** free energy density with respect to this distribution

$$\bar{f}(\beta) = \mathbb{E}_\xi f(\beta; \xi) \quad (46)$$

# Minimization of the quenched free energy

The quantity we are interested in is therefore

$$\bar{f}(\beta) = - \lim_{N \rightarrow \infty} \frac{1}{\beta N} \mathbb{E}_{\xi} \ln Z_{N,K}(\beta, \xi) \quad (47)$$

Our goal (now more complicated) is to find order-parameters  $\mathbf{m}$  that self-average in the thermodynamic limit such that (if (47) exists)

$$\bar{f}(\beta) = \min_{\mathbf{m}} \lim_{N \rightarrow \infty} \bar{f}_{N,K}^c(\mathbf{m}) \quad (48)$$

still holds.

# Guerra's interpolation

Given an interpolating Hamiltonian  $H_N^t$ , with  $H_N^1 = H_N$  and  $H_N^0$  an easier model, we compute the free energy via the Fundamental Theorem of Calculus

$$f^1 = f^0 + \int_0^1 dt \frac{d}{dt} f^t \quad (49)$$

# The replica trick

# The replica trick

- Step 1: Calculate  $\langle Z_N^n \rangle$  for  $n \in \mathbb{N}$ .

# The replica trick

- Step 1: Calculate  $\langle Z_N^n \rangle$  for  $n \in \mathbb{N}$ .
- Step 2: Conveniently forget that  $n$  is an integer and compute

$$\langle \ln Z_N \rangle = \lim_{n \rightarrow 0} \frac{1}{n} \ln \langle Z_N^n \rangle. \quad (50)$$

# The replica trick

- Step 1: Calculate  $\langle Z_N^n \rangle$  for  $n \in \mathbb{N}$ .
- Step 2: Conveniently forget that  $n$  is an integer and compute

$$\langle \ln Z_N \rangle = \lim_{n \rightarrow 0} \frac{1}{n} \ln \langle Z_N^n \rangle. \quad (50)$$

- Step 3: Conveniently forget...

# The replica trick

- Step 1: Calculate  $\langle Z_N^n \rangle$  for  $n \in \mathbb{N}$ .
- Step 2: Conveniently forget that  $n$  is an integer and compute

$$\langle \ln Z_N \rangle = \lim_{n \rightarrow 0} \frac{1}{n} \ln \langle Z_N^n \rangle. \quad (50)$$

- Step 3: Conveniently forget...the order of the limits:

$$\lim_{N \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{nN} \ln \langle Z_N^n \rangle = \lim_{n \rightarrow 0} \lim_{N \rightarrow \infty} \frac{1}{nN} \ln \langle Z_N^n \rangle \quad (51)$$

# Outline

- 5 Disordered systems
- 6 The Hopfield model**
- 7 Generalizations
- 8 Coming next
- 9 Sources and references

# The Hopfield model

We recall the Hamiltonian for the Hopfield model

$$H_{N,K}(\sigma; \xi) = -\frac{1}{2N} \sum_{\mu=1}^K \sum_{i,j=1}^N \xi_i^\mu \xi_j^\mu \sigma_i \sigma_j \quad (52)$$



**Figure 3.2:** Set of  $P = 6$  patterns stored in a Hopfield network of  $N = 625$  spins. Patterns are black and white images: the network is dealing with digital storage of information [Coolen et al. \(2005\)](#).

# The free energy

For finite  $K$ , the pseudo constrained free energy of the Hopfield model is

$$\tilde{f}^c(\beta) = \frac{1}{2} \sum_{\mu=1}^K (m_{\mu}^*)^2 - \beta^{-1} \mathbb{E}_{\xi} \ln 2 \cosh \left( \beta \sum_{\mu=1}^K m_{\mu}^* \xi^{\mu} \right) \quad (53)$$

which reduces to that of the Curie-Weiss model when retrieving only a single pattern.

$$\tilde{f}^{\text{CW},c}(m; \beta) = \frac{m^2}{2} - \beta^{-1} \ln 2 \cosh(\beta m) \quad (54)$$

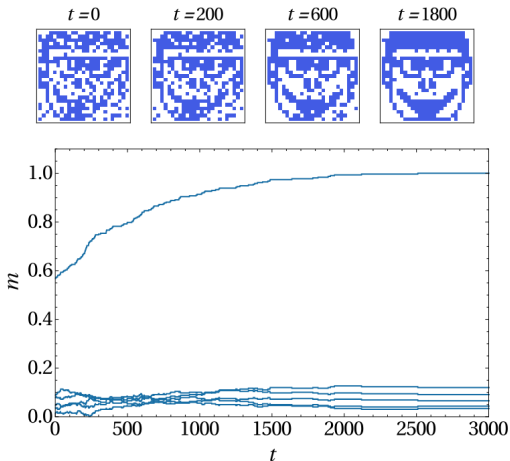
# Memory retrieval in the Hopfield model

If the model is within the retrieval region, and the initial state has enough overlap with one of the patterns, the model successfully retrieves the stored pattern.



**Figure 3.2: Set of  $P = 6$  patterns stored in a Hopfield network of  $N = 625$  spins.** Patterns are black and white images: the network is dealing with digital storage of information [Coolen et al. \(2005\)](#).

# Memory retrieval in the Hopfield model



**Figure 3.3: Example of a successful pattern reconstruction.** The Hopfield network is made of  $N = 625$  neurons and  $P = 6$  patterns are allocated in the synaptic matrix. Starting from a corrupted version of one of the patterns, the network is able to retrieve the associated pattern. We observe that, among the six Mattis magnetizations that quantify the retrieval of the six stored patterns, just one out of them grows up to one and its corresponding pattern is indeed retrieved by the network.

# Spurious states in the Hopfield model

Frustration in the Hopfield model manifests itself in the form of so-called **spurious states**, i.e. equilibrium states of the model that are not original patterns, but instead mixtures, for instance

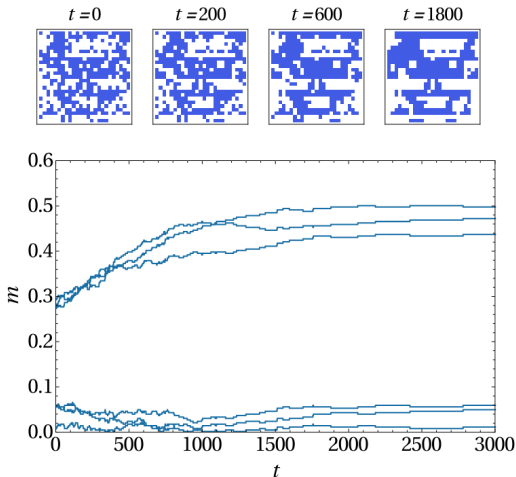
$$\tilde{\xi} = \text{sgn} (\xi^1 + \xi^2 + \xi^3). \quad (55)$$

The amount of spurious states grows exponentially in  $K$  and so eventually takes over the dynamics of the system. The maximum capacity of the Hopfield model turns out to be

$$K \lesssim \alpha_c N, \quad (56)$$

with  $\alpha_c \sim 0.138$ , in the thermodynamic limit.

# Spurious states in the Hopfield model



**Figure 3.4: Example of dynamics reconstruction ending in a spurious state.** The Hopfield network is the same as in Fig. 3.3. In this example, it is possible to observe that several (three) Mattis magnetizations raise sensibly over the noise due to the finite size effects and, correspondingly, the network has not been able to properly retrieve a single pattern, rather obtaining a useless mixture of the stored patterns.

# The free energy

If we instead take

$$\lim_{N \rightarrow \infty} \frac{K}{N} = \alpha \quad (57)$$

# The free energy

If we instead take

$$\lim_{N \rightarrow \infty} \frac{K}{N} = \alpha \quad (57)$$

then the replica-symmetric pseudo free energy for the recovery of  $\ell$  patterns is given by

$$\begin{aligned} \tilde{f}_{\text{RS}}^c(\beta, \alpha) &= \frac{1}{2} \sum_{\mu=1}^{\ell} (m_{\mu}^*)^2 + \frac{\alpha\beta}{2} p(1-q) \\ &+ \frac{\alpha}{2\beta} \left( \beta + \ln(1 - \beta(1-q)) - \frac{q\beta}{1 - \beta(1-q)} \right) \\ &- \beta^{-1} \int \frac{dz e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} \mathbb{E}_{\xi} \ln 2 \cosh \left( \beta \sum_{\mu=1}^{\ell} m_{\mu}^* \xi^{\mu} + \beta z \sqrt{\alpha p} \right) \quad (58) \end{aligned}$$

# Order parameters

To obtain the high load free energy, we need, apart from magnetizations, the replica overlaps

$$m_{\mu,N}(\boldsymbol{\xi}, \boldsymbol{\sigma}) = \frac{1}{N} \sum_{i=1}^N \xi_i^\mu \sigma_i, \quad \mu = 1, \dots, K \quad (59)$$

$$q_N(\boldsymbol{\sigma}^{(1)}, \boldsymbol{\sigma}^{(2)}) = \frac{1}{N} \sum_{i=1}^N \sigma_i^{(1)} \sigma_i^{(2)} \quad (60)$$

## Phase diagram

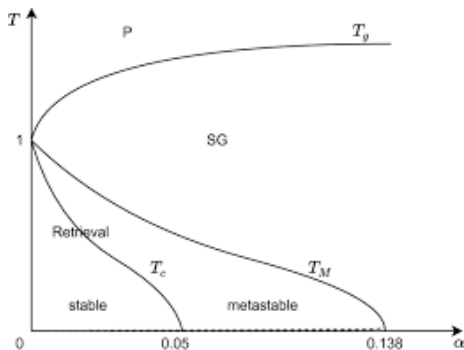


Figure: Phase diagram of the Hopfield model, first obtained by D. Amit, H. Gutfreund and H. Sompolinsky [3].

# Outline

- 5 Disordered systems
- 6 The Hopfield model
- 7 Generalizations**
- 8 Coming next
- 9 Sources and references

# Dataset noise

The first generalization we consider is dataset noise. Here, the patterns are replaced by a noisy dataset (the elements of which we call examples)

$$\{\xi^\mu\}_{\mu=1,\dots,K} \longrightarrow \{\eta_a^\mu\}_{a=1,\dots,M}^{\mu=1,\dots,K}, \quad (61)$$

with  $M$  denoting the size of the dataset and, for each  $\mu = 1, \dots, K$ ,  $a = 1, \dots, M$  and  $i = 1, \dots, N$ , we have

$$\eta_{ia}^\mu = \xi_i^\mu \chi_{ia}^\mu \quad (62)$$

$$P(\chi_{ia}^\mu = \pm 1) = \frac{1 \pm r}{2}, \quad (63)$$

with  $r = (0, 1]$  denoting the quality of the dataset.

## Two kinds of learning

In supervised learning, we assume to have access to the pattern labels corresponding to each example and so we can store the averages

$$J_{ij}^{\text{sup}} \propto \sum_{\mu=1}^K \sum_{a,b=1}^M \eta_{ia}^{\mu} \eta_{jb}^{\mu}, \quad i, j = 1, \dots, N. \quad (64)$$

On the other hand, in unsupervised learning, the best we can do is

$$J_{ij}^{\text{unsup}} \propto \sum_{\mu=1}^K \sum_{a=1}^M \eta_{ia}^{\mu} \eta_{ja}^{\mu}, \quad i, j = 1, \dots, N \quad (65)$$

In the case of supervised learning, it can be shown that the influence of the noise depends only on  $r$  and  $M$  through what we call the entropy:

$$\rho := \frac{1 - r^2}{Mr^2}. \quad (66)$$

# The $L$ -directional associative memory

An  $L$ -directional associative memory is composed of  $L$  independent Hopfield networks (known as modules), with the corresponding state denoted by  $(\sigma_i^k)_{i=1,\dots,N}^{k=1,\dots,L}$ , and the Hamiltonian is

$$J_{ij}^{\ell k} = \frac{1}{\mathcal{R}NM^2} \sum_{\mu=1}^K \sum_{a,b=1}^M g_{\ell k} \eta_{ia}^{\mu} \eta_{jb}^{\mu}, \quad i, j = 1, \dots, N, \quad (67)$$

where  $\mathcal{R} = r^2 + \frac{1-r^2}{M}$  is the second raw moment of the noise.

# Tackling disentanglement

The goal with disentanglement is, in the most simple case, to input a mixture of  $L$  patterns in every module and get back one of its constituents per module

$$(\sigma^{\text{mix}}, \dots, \sigma^{\text{mix}}) \longrightarrow (\xi^1, \dots, \xi^L), \quad (68)$$

for a mixture  $\sigma^{\text{mix}} = \text{sgn}(\xi^1 + \dots + \xi^L)$ . In order to achieve this, we take  $g_{\ell k} = \delta_{\ell k} - \lambda(1 - \delta_{\ell k})$ , with  $0 \leq \lambda < 0.5$ .

A detailed study of disentanglement in this system for  $L = 3$  (known as three-directional associative memory or TAM) was published in *Physica A* [4].

# Order parameters

The order parameters must now also be generalized from the Hopfield model. For each  $\mu = 1, \dots, K$  and  $\ell, k = 1, \dots, L$ , we have

$$m_{\mu}^{\ell}(\boldsymbol{\eta}, \boldsymbol{\sigma}) = \frac{1}{N} \sum_{i=1}^N \xi_i^{\mu} \sigma_i^{\ell} \quad (69)$$

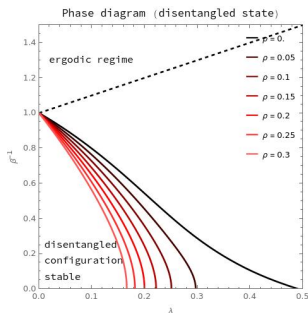
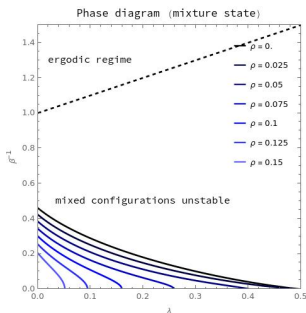
$$q_{12}^{\ell k}(\boldsymbol{\sigma}^{(1)}, \boldsymbol{\sigma}^{(2)}) = \frac{1}{N} \sum_{i=1}^N \sigma_i^{(1),\ell} \sigma_i^{(2),k} \quad (70)$$

$$q_{11}^{\ell k}(\boldsymbol{\sigma}) = \frac{1}{N} \sum_{i=1}^N \sigma_i^{\ell} \sigma_i^k \quad (71)$$

And the entropy now becomes the entropy per layer.

# Self-consistency equations

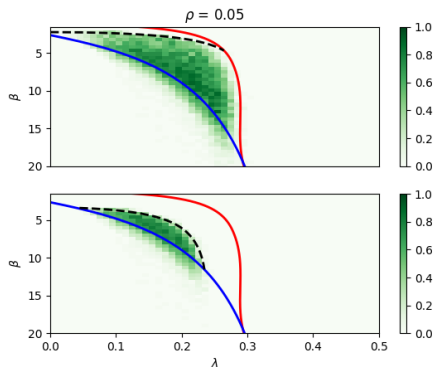
The Statistical Mechanics framework presented before allowed us to get self-consistency equations for the order parameters of this system, which we solved in the low-load (i.e. finite number of patterns  $K$  in the thermodynamic limit  $N \rightarrow \infty$ ).



# Prediction of disentanglement

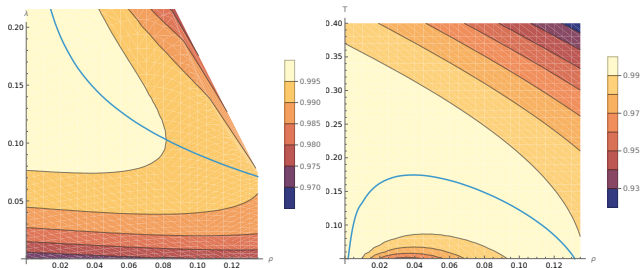
We observe, in the phase diagrams, that there's a region of state space where mixtures are unstable but disentangled states are stable.

Our hypothesis, shown to be valid through Monte-Carlo experiments, is that disentanglement is possible in this region.

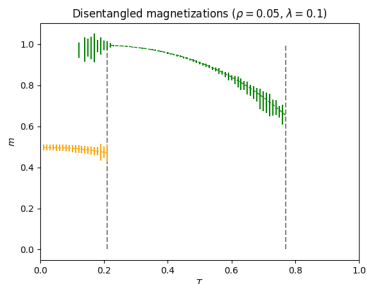
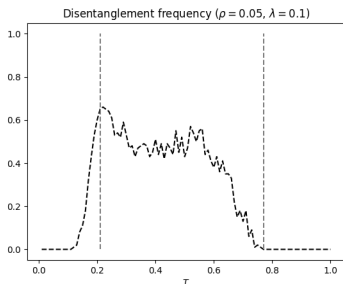


# Optimal parameters

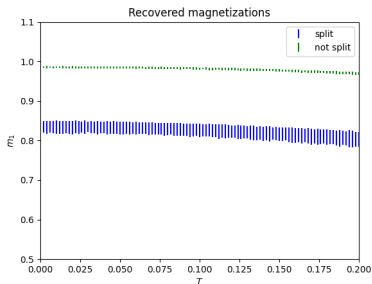
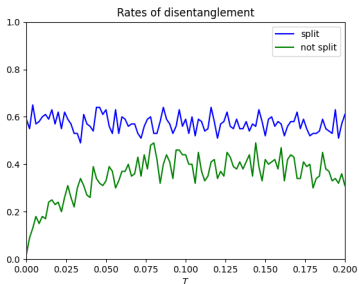
By maximizing retrieval magnetizations (i.e. the magnetizations of the output) within the disentanglement region, we also obtained optimal control parameters for each value of the dataset entropy per module  $\rho$ .



# More experimental evidence



# Splitting



# Outline

- 5 Disordered systems
- 6 The Hopfield model
- 7 Generalizations
- 8 Coming next**
- 9 Sources and references

# Prof. Elena Agliari's talk (April 24-th)

Using the dreaming model

$$J_{ij}^D = \frac{1}{N} \sum_{\mu, \nu=1}^K \xi_i^\mu \left( \frac{1+t}{\mathbb{1} + t\mathbf{C}} \right)_{\mu\nu} \xi_j^\nu, \quad (72)$$

with

$$C_{ab}^{\mu\nu} = \frac{1}{N} \sum_{i=1}^N \xi_i^\mu \xi_i^\nu. \quad (73)$$

to boost learning from unsupervised datasets.

# Outline

- 5 Disordered systems
- 6 The Hopfield model
- 7 Generalizations
- 8 Coming next
- 9 Sources and references**

# Bibliography

- [1] W. Little. “The existence of persistent states in the brain”. In: *Mathematical Biosciences* 19.1 (1974), pp. 101–120.
- [2] J. J. Hopfield. “Neural networks and physical systems with emergent collective computational abilities”. In: *Proceedings of the National Academy of Sciences* 79.8 (1982), pp. 2554–2558.
- [3] D. J. Amit, H. Gutfreund, and H. Sompolinsky. “Storing Infinite Numbers of Patterns in a Spin-Glass Model of Neural Networks”. In: *Phys. Rev. Lett.* 55 (14 Sept. 1985), pp. 1530–1533.
- [4] E. Agliari, A. Fachechi, and P. D. Mourão. “The beneficial role of noises for disentanglement tasks in modular Hebbian networks”. In: *Physica A: Statistical Mechanics and its Applications* (2025), p. 131134.
- [6] C. A.C.C., R. Kuehn, and P. Sollich. *Theory of Neural Information Processing Systems*. Oxford University Press, 2005.

## Other resources

- M. Mezard, G. Parisi, and M. A. Virasoro. *Spin Glass Theory And Beyond: An Introduction To The Replica Method And Its Applications*. World Scientific Publishing Company, 1987.
- M. Talagrand. *Mean Field Models for Spin Glasses*. Springer, 2010.
- A. Bovier. *Statistical Mechanics of Disordered Systems: A Mathematical Perspective*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2006.

My e-mail: paulo.duartemourao@uniroma1.it

# Thank you to my advisors and collaborators!



# Replicated partition function

We recall that our goal is to compute

$$\bar{f}(\beta, \alpha) = \lim_{\substack{N \rightarrow \infty \\ K/N \rightarrow \alpha}} -\frac{1}{\beta N} \mathbb{E}_{\xi} \ln Z_{N,K}(\beta; \xi) \quad (74)$$

with

$$Z_{N,K}(\beta; \xi) = \sum_{\sigma \in \Omega} e^{-\beta H_{N,K}(\sigma; \xi)} = \sum_{\sigma \in \Omega} e^{\frac{\beta}{2N} \sum_{\mu=1}^K \sum_{i,j=1}^N \xi_i^{\mu} \xi_j^{\mu} \sigma_i \sigma_j} \quad (75)$$

Then

$$Z_{N,K}(\beta; \xi)^n = \sum_{\sigma^1, \dots, \sigma^n \in \Omega} e^{\frac{\beta}{2N} \sum_{r=1}^n \sum_{\mu=1}^K \sum_{i,j=1}^N \xi_i^{\mu} \xi_j^{\mu} \sigma_i^{(r)} \sigma_j^{(r)}} \quad (76)$$

# Signal vs noise

Taking  $\ell$  to be the number of patterns recovered, we have

$$Z_{N,K}(\beta; \xi)^n = \sum_{\sigma^1, \dots, \sigma^n \in \Omega} z_{\text{signal}}(\{\sigma\}) \times z_{\text{noise}}(\{\sigma\}) \quad (77)$$

with

$$z_{\text{signal}}(\{\sigma\}) = e^{\frac{\beta}{2N} \sum_{a=1}^n \sum_{\mu=1}^{\ell} \sum_{i,j=1}^N \xi_i^\mu \xi_j^\mu \sigma_i^a \sigma_j^a} \quad (78)$$

$$z_{\text{noise}}(\{\sigma\}) = e^{\frac{\beta}{2N} \sum_{a=1}^n \sum_{\mu=\ell+1}^K \sum_{i,j=1}^N \xi_i^\mu \xi_j^\mu \sigma_i^a \sigma_j^a}. \quad (79)$$

Furthermore, the disorder is independent between both, and the goal is to apply the replica trick to the disorder average of the non-recovered patterns, which we denote  $\mathbb{E}'_{\xi}$

$$\mathbb{E}'_{\xi} Z_{N,K}(\beta; \xi)^n = \sum_{\sigma^1, \dots, \sigma^n \in \Omega} z_{\text{signal}}(\{\sigma\}) \times \mathbb{E}_{\xi} z_{\text{noise}}(\{\sigma\}) \quad (80)$$

## The signal part

For the signal part, we insert the order parameters corresponding to the magnetizations, as we did for the Curie-Weiss model

$$\begin{aligned} Z_{\text{signal}}(\{\sigma\}) &= \exp \left[ \frac{\beta}{2N} \sum_{a=1}^n \sum_{\mu=1}^{\ell} \sum_{i,j=1}^N \xi_i^{\mu} \xi_j^{\mu} \sigma_i^a \sigma_j^a \right] \\ &= \left( \frac{N}{2\pi} \right)^{n\ell} \int \left[ \prod_{a,\mu=1}^{n,\ell} dm_{\mu}^a \delta \left( m_{\mu}^a - \frac{1}{N} \sum_{i=1}^N \xi_i^{\mu} \sigma_i^a \right) \right] \\ &\quad \times \exp \left[ \frac{\beta N}{2} \sum_{a=1}^n \sum_{\mu=1}^{\ell} (m_{\mu}^a)^2 \right] \\ &= \left( \frac{N}{2\pi} \right)^{n\ell} \int \left[ \prod_{a,\mu=1}^{n,\ell} dm_{\mu}^a d\hat{m}_{\mu}^a \exp \left[ iNm_{\mu}^a \hat{m}_{\mu}^a - i\hat{m}_{\mu}^a \sum_{i=1}^N \xi_i^{\mu} \sigma_i^a \right] \right] \\ &\quad \times \exp \left[ \frac{\beta N}{2} \sum_{a=1}^n \sum_{\mu=1}^{\ell} (m_{\mu}^a)^2 \right] \end{aligned}$$

# The Hubbard-Stratonovich Transform

To treat the noise part, we make use of the Hubbard-Stratonovich transform

$$e^{\frac{x^2}{2}} = \sqrt{\frac{1}{2\pi}} \int_{-\infty}^{\infty} dz e^{-\frac{z^2}{2} + zx} \quad (81)$$

to get

$$\begin{aligned} z_{\text{noise}}(\{\sigma\}) &= \exp \left[ \frac{\beta}{2N} \sum_{a=1}^n \sum_{\mu=\ell+1}^K \sum_{i,j=1}^N \xi_i^\mu \xi_j^\mu \sigma_i^a \sigma_j^a \right] \\ &= \int \prod_{\substack{a=1, \dots, n \\ \mu=\ell+1, \dots, K}} \frac{dz_\mu^a}{\sqrt{2\pi}} \exp \left[ -\frac{1}{2} (z_\mu^a)^2 + \sqrt{\frac{\beta}{N}} \sum_{i=1}^N \xi_i^\mu \sigma_i^a z_\mu^a \right] \end{aligned}$$

## Average over the disorder

Averaging over the disorder now gives, with  $\mathcal{D}\mathbf{z} = \prod_{a,\mu} dz_{\mu}^a e^{-(z_{\mu}^a)^2/2}/\sqrt{2\pi}$ ,

$$\begin{aligned}\mathbb{E}_{\xi} \mathbb{Z}_{\text{noise}}(\{\sigma\}) &= \int \mathcal{D}\mathbf{z} \prod_{i,\mu} \mathbb{E}_{\xi} \exp \left[ \sqrt{\frac{\beta}{N}} \sum_{a=1}^n \xi_i^{\mu} \sigma_i^a z_{\mu}^a \right] \\ &= \int \mathcal{D}\mathbf{z} \prod_{i,\mu} \exp \left[ \ln \cosh \left( \sqrt{\frac{\beta}{N}} \sum_{a=1}^n \sigma_i^a z_{\mu}^a \right) \right] \\ &= \left( \int \mathcal{D}\tilde{\mathbf{z}} \exp \left[ \frac{\beta}{2N} \sum_{a,b=1}^n \sum_{i=1}^N z^a z^b \sigma_i^a \sigma_i^b + \mathcal{O}\left(\frac{1}{N}\right) \right] \right)^{K-\ell}\end{aligned}$$

and  $\mathcal{D}\tilde{\mathbf{z}} = \prod_a dz^a e^{-(z^a)^2/2}/\sqrt{2\pi}$ .

## Inserting the overlaps

We now insert the overlaps, again via delta distributions. Denoting  $\tilde{K} := K - \ell$ , we have

$$\begin{aligned}\mathbb{E}_{\xi, Z^{\text{noise}}}(\{\sigma\}) &\approx \left(\frac{N}{2\pi}\right)^{n^2-n} \prod_{\substack{a,b=1 \\ a \neq b}}^n \int dQ_{ab} \delta\left(Q_{ab} - \frac{1}{N} \sum_{i=1}^N \sigma_i^a \sigma_i^b\right) \\ &\times \int \mathcal{D}\tilde{\mathbf{z}} \exp\left[\frac{\beta}{2} \sum_{a,b=1}^n z^a z^b Q_{ab}\right] \\ &= \left(\frac{N}{2\pi}\right)^{n^2-n} (2\pi)^{\frac{n^2}{2}} \int d\mathbf{Q} \exp\left[-\frac{\tilde{K}}{2} \ln \det(\mathbb{1} - \beta\mathbf{Q})\right] \\ &\times \int d\hat{\mathbf{Q}} \exp\left[i \sum_{\substack{a,b=1 \\ a \neq b}}^n \left(NQ_{ab}\hat{Q}_{ab} - \sum_{i=1}^N \hat{Q}_{ab}\sigma_i^a\sigma_i^b\right)\right]\end{aligned}$$

# Replicated partition function

Putting everything together and reparameterizing  $\hat{Q}_{ab} \rightarrow i\frac{\alpha\beta^2}{2}P_{ab}$ , we get

$$\begin{aligned} \mathbb{E}'_{\xi} Z_{N,K}(\beta; \xi)^n &\propto \int d\mathbf{m} d\hat{\mathbf{m}} d\mathbf{Q} d\mathbf{P} \exp \left[ N \sum_{a,\mu=1}^{n,\ell} \left( \frac{\beta}{2} (m_{\mu}^a)^2 + i m_{\mu}^a \hat{m}_{\mu}^a \right) \right] \\ &\times \exp \left[ -\frac{\tilde{K}}{2} \ln \det (\mathbf{1} - \beta \mathbf{Q}) - \frac{\alpha\beta^2}{2} N \sum_{\substack{a,b=1 \\ a \neq b}}^n Q_{ab} P_{ab} \right] \\ &\times \sum_{\sigma^1, \dots, \sigma^n \in \Omega} \prod_{i=1}^N \exp \left[ -\sum_{a=1}^n \sum_{\mu=1}^{\ell} i \hat{m}_{\mu}^a \xi_{\mu}^{\sigma_i} \sigma_i^a + \frac{\alpha\beta^2}{2} \sum_{\substack{a,b=1 \\ a \neq b}}^n P_{ab} \sigma_i^a \sigma_i^b \right] \end{aligned}$$

# Replicated partition function

Furthermore, assuming the magnetizations self-average, the last product in the  $i$ -index factorizes, meaning that we get

$$\begin{aligned} \mathbb{E}'_{\xi} Z_{N,K}(\beta; \xi)^n &\propto \int d\mathbf{m} d\hat{\mathbf{m}} d\mathbf{Q} d\mathbf{P} \exp \left[ N \sum_{a,\mu=1}^{n,\ell} \left( \frac{\beta}{2} (m_{\mu}^a)^2 + i m_{\mu}^a \hat{m}_{\mu}^a \right) \right] \\ &\times \exp \left[ N \left( -\frac{\tilde{K}}{2N} \ln \det (\mathbf{1} - \beta \mathbf{Q}) - \frac{\alpha \beta^2}{2} \sum_{\substack{a,b=1 \\ a \neq b}}^n Q_{ab} P_{ab} \right) \right] \\ &\times \left( \sum_{\sigma^1, \dots, \sigma^n = \pm 1} \exp \left[ -\sum_{a=1}^n \sum_{\mu=1}^{\ell} i \hat{m}_{\mu}^a \xi^{\mu} \sigma^a + \frac{\alpha \beta^2}{2} \sum_{\substack{a,b=1 \\ a \neq b}}^n P_{ab} \sigma^a \sigma^b \right] \right)^N \end{aligned}$$

# Taking the thermodynamic limit

Furthermore, assuming the magnetizations self-average, the last product in the  $i$ -index factorizes, meaning that we get

$$\lim_{N \rightarrow \infty} -\frac{1}{\beta N} \mathbb{E}'_{\xi} Z_{N,K}(\beta; \xi)^n = \min_{\mathbf{m}, \hat{\mathbf{m}}, \mathbf{Q}, \mathbf{P}} f(\mathbf{m}, \hat{\mathbf{m}}, \mathbf{Q}, \mathbf{P}) \quad (82)$$

with

$$f = -\frac{1}{2} \sum_{a,\mu=1}^{n,\ell} (m_{\mu}^a)^2 - i\beta^{-1} \sum_{a,\mu=1}^{n,\ell} m_{\mu}^a \hat{m}_{\mu}^a + \frac{\alpha\beta}{2} \sum_{\substack{a,b=1 \\ a \neq b}}^n Q_{ab} P_{ab} + \frac{\alpha}{2\beta} \ln \det(\mathbf{1} - \beta \mathbf{Q})$$
$$- \beta^{-1} \ln \sum_{\sigma} \exp \left[ - \sum_{a=1}^n \sum_{\mu=1}^{\ell} i \hat{m}_{\mu}^a \xi^{\mu} \sigma^a + \frac{\alpha\beta^2}{2} \sum_{\substack{a,b=1 \\ a \neq b}}^n P_{ab} \sigma^a \sigma^b \right]$$

# Taking the thermodynamic limit

By minimizing already with respect to  $\mathbf{m}$ , we get  $\hat{m}_\mu^a = i\beta m_\mu^a$ , which gives us

$$\lim_{N \rightarrow \infty} -\frac{1}{\beta N} \mathbb{E}'_{\xi} Z_{N,K}(\beta; \xi)^n = \min_{\mathbf{m}, \mathbf{Q}, \mathbf{P}} \tilde{f}(\mathbf{m}, \mathbf{Q}, \mathbf{P}) \quad (83)$$

with

$$\begin{aligned} \tilde{f} = & \frac{1}{2} \sum_{a,\mu=1}^{n,\ell} (m_\mu^a)^2 + \frac{\alpha\beta}{2} \sum_{\substack{a,b=1 \\ a \neq b}}^n Q_{ab} P_{ab} + \frac{\alpha}{2\beta} \ln \det(\mathbb{1} - \beta \mathbf{Q}) \\ & - \beta^{-1} \ln \sum_{\sigma} \exp \left[ \beta \sum_{a=1}^n \sum_{\mu=1}^{\ell} m_\mu^a \xi^\mu \sigma^a + \frac{\alpha\beta^2}{2} \sum_{\substack{a,b=1 \\ a \neq b}}^n P_{ab} \sigma^a \sigma^b \right] \end{aligned}$$

# Replica symmetry ansatz

To be able to proceed with the computations and take the limit  $n \rightarrow 0$ , we must make assumptions on the structure of the matrices  $\mathbf{m}$ ,  $\hat{\mathbf{m}}$ ,  $\mathbf{P}$  and  $\mathbf{Q}$ . The so-called replica-symmetry ansatz assumes complete independence on the choice of replicas, meaning that

$$m_{\mu}^a = m_{\mu} \quad (84)$$

$$Q_{ab} = \delta_{ab} + (1 - q)\delta_{ab} \quad (85)$$

$$P_{ab} = p \quad (86)$$

# Replica symmetry ansatz

From this ansatz, we get, for the first and second terms

$$\lim_{n \rightarrow 0} \frac{1}{n} \sum_{a, \mu=1}^{n, \ell} (m_{\mu}^a)^2 = \sum_{\mu=1}^{\ell} (m_{\mu})^2 \quad (87)$$

$$\lim_{n \rightarrow 0} \frac{1}{n} \sum_{\substack{a, b=1 \\ a \neq b}}^n Q_{ab} P_{ab} = -pq \quad (88)$$

# Replica symmetry ansatz

For the third term, we have that

$$(\mathbb{1} - \beta \mathbf{Q})_{ab} = (1 - \beta)\delta_{ab} + (1 - \beta q)(1 - \delta_{ab}) \quad (89)$$

whose eigenvalues are

$$\lambda_1 = 1 - \beta(1 - q) - \beta qn \quad (90)$$

$$\lambda_2 = 1 - \beta(1 - q) \quad (91)$$

with multiplicities 1 and  $n - 1$ , respectively.

# Replica symmetry ansatz

We therefore have

$$\begin{aligned}\log \det (\mathbb{1} - \beta \mathbf{Q}) &= \ln(1 - \beta(1 - q) - \beta q n) + (n - 1) \ln(1 - \beta(1 - q)) \\ &= n \ln(1 - \beta(1 - q)) - n \frac{\beta q}{1 - \beta(1 - q)} + \mathcal{O}(n^2)\end{aligned}$$

and so

$$\lim_{n \rightarrow 0} \log \det (\mathbb{1} - \beta \mathbf{Q}) = \ln(1 - \beta(1 - q)) - \frac{\beta q}{1 - \beta(1 - q)} \quad (92)$$