The early bird regularizes better

Elena Agliari

Dipartimento di Matematica, Sapienza Università di Roma Istituto Nazionale di Alta Matematica, Roma





Joint work with: Adriano Barra, Salento Alberto Fachechi, Sapienza Linda Albanese, Salento Miriam Aquaro, Sapienza Francesco Alemanno, Bologna



iNδAM

A Day on Statistical Physics for Machine Learning Dipartimento di Matematica — Università Tor Vergata 15 September 2023

Hard sciences for machine learning



Mathematical control

- design optimal architecture a priori
- estimate hyperparameters
- check dataset size

. . . .

 \rightarrow reduce AI carbon footprint!



*transformer (213M parameters) w/ neural architecture search

Hard sciences for machine learning



Statistical physics perspective

Information-processing capability as emerging collective behavior





temperature/depth

Mimics retrieval capabilities

Pattern recognition Pattern reconstruction Denoising

Content addressable memory



CAPTCHA









Mimics retrieval capabilities

Pattern recognition Pattern reconstruction Denoising

Content addressable memory



CAPTCHA













[Hopfield - PNAS '82; Pastur, Figotin - Theor. Math. Phys. '78]

0

 σ_2

Set of N binary neurons $\boldsymbol{\sigma} \in \{-1, +1\}^N$ with pair-wise interaction strengths $\boldsymbol{J} \in \mathbb{R}^{N \times N}$

Neural dynamics
$$\sigma_i^{(n+1)} = \operatorname{sign}\left[\sum_j J_{ij}\sigma_j^{(n)} + T\zeta_i\right]$$
, with $T \in \mathbb{R}^+$, ζ_i iid r.v.

Task: reconstruct *K* binary vectors (patterns) of length *N* $\boldsymbol{\xi} := \{\boldsymbol{\xi}^{\mu}\}_{\mu=1,...,K} \in \{-1, +1\}^{N \times K}$

⇒ choose $J = J(\xi)$ s.t. each pattern is (δ, ϵ) -stable for $\delta, \epsilon > 0$



Set of N binary neurons $\boldsymbol{\sigma} \in \{-1, +1\}^N$ with pair-wise interaction strengths $\boldsymbol{J} \in \mathbb{R}^{N \times N}$

Neural dynamics
$$\sigma_i^{(n+1)} = \operatorname{sign}\left[\sum_j J_{ij}\sigma_j^{(n)} + T\zeta_i\right]$$
, with $T \in \mathbb{R}^+$, ζ_i iid r.v.

Task: reconstruct *K* binary vectors (patterns) of length *N* $\boldsymbol{\xi} := \{\boldsymbol{\xi}^{\mu}\}_{\mu=1,...,K} \in \{-1, +1\}^{N \times K}$

⇒ choose $J = J(\xi)$ s.t. each pattern is (δ, ϵ) -stable for $\delta, \epsilon > 0$



Gibbs measure is reversible w.r.t. neural dynamics

$$\rho_{N,\boldsymbol{J},\boldsymbol{\beta}}(\boldsymbol{\sigma}) := \frac{e^{-\boldsymbol{\beta}\mathcal{H}_{N,\boldsymbol{J}}(\boldsymbol{\sigma})}}{\mathcal{Z}_{N,\boldsymbol{J},\boldsymbol{\beta}}}, \text{ with } \mathcal{H}_{N,\boldsymbol{J}}(\boldsymbol{\sigma}) = -\sum_{i,j=1}^{N} \sum_{\mu=1}^{K} \sigma_{i} J_{ij} \sigma_{j}$$

$$\Rightarrow \text{ choose } \boldsymbol{J} = \boldsymbol{J}(\boldsymbol{\xi}) \text{ s.t. } \boldsymbol{\xi}^{\mu} = \operatorname*{argmin}_{\boldsymbol{\sigma} \in \{-1,+1\}^{N}} \mathcal{H}_{N,\boldsymbol{J}}(\boldsymbol{\sigma}) \text{ for } \mu = 1, ..., K$$

Hebb's rule
$$J_{ij}(\xi) := (1 - \delta_{ij}) \sum_{\mu=1}^{N} \xi_i^{\mu} \xi_j^{\mu}$$

AGS solution

[Amit, Gutfreund, Sompolinsky - Phys. Rev. A '87]



Noise

Load
$$\alpha := \lim_{N \to \infty} \frac{K}{N}$$

AGS solution

$$\mathcal{H}_{N,K,\boldsymbol{\xi}}(\boldsymbol{\sigma}) = -\frac{1}{2N} \sum_{\substack{i, j = 1 \\ i \neq j}}^{N} \sum_{\substack{\mu=1 \\ i \neq j}}^{K} \sigma_i \ \xi_i^{\mu} \xi_j^{\mu} \ \sigma_j$$

Patterns
$$\xi_{i}^{\mu} \sim_{iid} \mathcal{R}, \forall i, \mu$$

Magnetization $m_{\mu} := \frac{1}{N} \sum_{i=1}^{N} \xi_{i}^{\mu} \sigma_{i},$
Overlap $q_{ab} := \frac{1}{N} \sum_{i=1}^{N} \sigma_{i}^{(a)} \sigma_{i}^{(b)}$

T

Noise

Load
$$\alpha := \lim_{N \to \infty} \frac{K}{N}$$

Revise Hebb's rule to enhance storing (e.g., Kohonen '72, Personnaz '85, Kanter&Sompolinsky '86, Opper al. '22)

possibly including iterative protocols

(e.g., Hopfield '83, Hopfield '84, Parisi '86, Sherrington '95, van Hemmen '97, Marinari '18)

 $1/\beta$

[Amit, Gutfreund, Sompolinsky - Phys. Rev. A '87]



 α



[Hopfield et al. - Nature Lett. '83]



[Hopfield et al. - Nature Lett. '83]

Analytical unlearning rule

$$J_{ij}^{(n)} \leftarrow J_{ij}^{(n-1)} - \varphi_i(\sigma^*) \varphi_j(\sigma^*)$$

$$J_{ij}(t) = \frac{1}{N} \sum_{\mu,\nu} \xi_i^{\mu} (1 + tC)_{\mu\nu}^{-1} \xi_j^{\nu}$$

where $C_{\mu\nu} := \frac{1}{N} \sum_{i=1}^N \xi_i^{\mu} \xi_j^{\mu}$ pattern correlation matrix

[Dotsenko et al. - J. Phys. A '91, Plakhov - IEEE Patt. Recogn. '94]

1.0
0.8

$$t = 0.1$$

 $t = 0.2$
 0.4
 0.2
 0.0
 0.2
 0.4
 0.5
 $t = 1.0$
 α

Remotion&Reinforcement

X X /

$$\begin{split} J_{ij}^{(n)} &\leftarrow J_{ij}^{(n-1)} + \frac{1}{1+n} [J^{(n-1)} - (J^{(n-1)})^2] \\ J_{ij}(t) &= \frac{1}{N} \sum_{\mu,\nu} \xi_i^{\mu} \Big(\frac{1+t}{1+tC} \Big)_{\mu\nu} \xi_j^{\nu} \\ \text{where } C_{\mu\nu} &:= \frac{1}{N} \sum_{i=1}^N \xi_i^{\mu} \xi_j^{\mu} \text{ pattern correlation matrix} \end{split}$$

[Fachechi, Agliari, Barra - Neur. Net. '19]

$$(1+t) \quad \cdot \quad (\boldsymbol{I}+t\boldsymbol{C})^{-1}$$

reinforcement (SW) remotion (REM)

t "sleeping time"

[Crick, Mitchinson - Nature '83; Stickgold - Nature '05; Diekelmann, Born - Nature Rev. Neurosc. '10]



Remotion&Reinforcement

X X /

$$\begin{split} J_{ij}^{(n)} &\leftarrow J_{ij}^{(n-1)} + \frac{1}{1+n} [J^{(n-1)} - (J^{(n-1)})^2] \\ J_{ij}(t) &= \frac{1}{N} \sum_{\mu,\nu} \xi_i^{\mu} \Big(\frac{1+t}{1+tC} \Big)_{\mu\nu} \xi_j^{\nu} \\ \text{where } C_{\mu\nu} &:= \frac{1}{N} \sum_{i=1}^N \xi_i^{\mu} \xi_j^{\mu} \text{ pattern correlation matrix} \end{split}$$

[Fachechi, Agliari, Barra - Neur. Net. '19]

$$J_{ij}(t) = \frac{1}{N} \sum_{\mu,\nu} \xi_i^{\mu} \left(\frac{1+t}{1+tC} \right)_{\mu\nu} \xi_j^{\nu}$$

 $t \rightarrow 0$: Hebb's rule

 $t \rightarrow \infty$: Kohonen rule '72, Kanter-Sompolinsky '86

"Dreaming" Hopfield network

$$\mathcal{H}_{N,K,\boldsymbol{\xi},t}^{(\mathrm{DHN})}(\boldsymbol{\sigma}) = -\frac{1}{2N} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{K,K} \sigma_i \xi_i^{\mu} \left(\frac{1+t}{\boldsymbol{I}+t\boldsymbol{C}}\right)_{\mu,\nu} \xi_j^{\nu} \sigma_j$$



Interpolating technique

$$\mathscr{H}_{N}(\boldsymbol{\sigma};\boldsymbol{J}) \to \tilde{\mathscr{H}}_{N}(\boldsymbol{\sigma};\boldsymbol{J},\boldsymbol{J}') = s \ \mathscr{H}_{N}(\boldsymbol{\sigma};\boldsymbol{J}) + (1-s) \ \mathscr{H}_{N}^{easy}(\boldsymbol{\sigma};\boldsymbol{J}'), \ s \in [0,1]$$

 $ilde{\mathscr{H}}_N(\pmb{\sigma};\pmb{J},\pmb{J}') ext{ yields } ilde{\mathscr{Z}}_{N,\beta}(\pmb{J},\pmb{J}'), extit{ } ilde{
ho}_{N,\beta}(\pmb{\sigma};\pmb{J},\pmb{J}')$

$$\mathscr{A}_{N} = \tilde{\mathscr{A}}_{N}(1) = \tilde{\mathscr{A}}_{N}(0) + \int_{0}^{1} \frac{d\tilde{\mathscr{A}}_{N}(s)}{ds} \bigg|_{s=s'} ds'$$



[Guerra - Fields Inst. Comm. '01]

Hopfield w/reinforcement&remotion

$$\mathscr{H}_{N,K,\boldsymbol{\xi},t}^{(\mathrm{DHM})}(\boldsymbol{\sigma}) = -\frac{1}{2N} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{K,K} \sigma_i \xi_i^{\mu} \left(\frac{1+t}{\boldsymbol{I}+t\boldsymbol{C}}\right)_{\mu,\nu} \xi_j^{\nu} \sigma_j$$



Hopfield w/reinforcement&remotion

$$\mathscr{H}_{N,K,\boldsymbol{\xi},t}^{(\mathrm{DHM})}(\boldsymbol{\sigma}) = -\frac{1}{2N} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{K,K} \sigma_i \xi_i^{\mu} \left(\frac{1+t}{\boldsymbol{I}+t\boldsymbol{C}}\right)_{\mu,\nu} \xi_j^{\nu} \sigma_j$$

Goal: explicit expression for quenched pressure





Retrieval region gets wider

Spin-glass collapses

Ergodic region gets wider

Allocate more information with the same resources

[Agliari, Alemanno, Barra, Fachechi - Neur. Net. '19, J. Stat. '19]

Replace *archetypes* by a sample of *examples*

Replace *archetypes* by a sample of *examples*



[Agliari, Alemanno, Barra, De Marzo - Neur. Net. '22]

Replace *archetypes* by a sample of *examples*

$$K \text{ "archetypes" } \{ \boldsymbol{\xi}^{\mu} \}_{\mu=1,\dots,K} \text{ 123} \quad K \times M \text{ "blurred" examples } \{ \boldsymbol{\eta}^{\mu,a} \}_{\mu=1,\dots,K}^{a=1,\dots,M} \text{ } \{$$

$$J_{ij}^{(sup)}(\eta) \propto \sum_{\mu,\nu=1}^{K} \sum_{a,b=1}^{M} \eta_{i}^{\mu,a} \left(\frac{1+t}{I+tC(\bar{\eta})}\right)_{\mu\nu} \eta_{j}^{\nu,b} = \frac{1}{N} \sum_{\mu,\nu=1}^{K} \bar{\eta}_{i}^{\mu} \left(\frac{1+t}{I+tC(\bar{\eta})}\right)_{\mu\nu} \bar{\eta}_{j}^{\nu}$$

•.

$$J_{ij}(\boldsymbol{\xi}) \propto \sum_{\mu,\nu=1}^{K} \xi_i^{\mu} \left(\frac{1+t}{\boldsymbol{I}+t\boldsymbol{C}(\boldsymbol{\xi})} \right)_{\mu\nu} \xi_j^{\nu}$$

$$J_{ij}^{(unsup)}(\eta) \propto \sum_{(\mu,a)=(1,1)}^{(K,M)} \sum_{(\nu,b)=(1,1)}^{(K,M)} \eta_i^{\mu,a} \left(\frac{1+t}{I+tC(\eta)}\right)_{(\mu,a),(\nu,b)} \eta_j^{\nu,b}$$

Replace *archetypes* by a sample of *examples*

$$K \text{ "archetypes" } \{\boldsymbol{\xi}^{\mu}\}_{\mu=1,\dots,K} \text{ 123} \quad K \times M \text{ "blurred" examples } \{\boldsymbol{\eta}^{\mu,a}\}_{\mu=1,\dots,K}^{a=1,\dots,M} \text{ } \{\boldsymbol{\eta}^{\mu,a}\}_{\mu=1,\dots,K}^{$$

$$J_{ij}^{(sup)}(\eta) \propto \sum_{\mu,\nu=1}^{K} \sum_{a,b=1}^{M} \eta_{i}^{\mu,a} \Big(\frac{1+t}{I+tC(\bar{\eta})} \Big)_{\mu\nu} \eta_{j}^{\nu,b} = \frac{1}{N} \sum_{\mu,\nu=1}^{K} \bar{\eta}_{i}^{\mu} \Big(\frac{1+t}{I+tC(\bar{\eta})} \Big)_{\mu\nu} \bar{\eta}_{j}^{\nu}$$

•.

$$J_{ij}(\boldsymbol{\xi}) \propto \sum_{\mu,\nu=1}^{K} \xi_i^{\mu} \left(\frac{1+t}{\boldsymbol{I}+t\boldsymbol{C}(\boldsymbol{\xi})} \right)_{\mu\nu} \xi_j^{\nu}$$

$$\mathscr{H}_{N,K,\boldsymbol{\eta},t}^{(\text{DHN})}(\boldsymbol{\sigma}) = -\frac{1}{2NM^2} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{K,K} \sum_{a,b=1}^{M} \sigma_i \eta_i^{\mu,a} \Big(\frac{1+t}{\boldsymbol{I}+t\boldsymbol{C}(\boldsymbol{\bar{\eta}})}\Big) \eta_j^{\nu,b} \sigma_j$$

$$T, \alpha, t$$

 $\rho(M, r) := \frac{1 - r^2}{Mr^2}$ Dataset "entropy"

 T, α, t $\rho(M, r) := \frac{1 - r^2}{Mr^2}$ Dataset "entropy"



 T, α, t $\rho(M, r) := \frac{1 - r^2}{Mr^2}$ Dataset "entropy"





 $M_c(m_{\rm X},t)$ lowest number of examples for success ($m_{\mu}>m_{\rm X}$)

$$S(m_{\times}, t) = 1 - \frac{M_c}{\max_{t \in \mathbb{R}^+} M_c}$$
 data "saving"

[Aquaro, Alemanno, Kanter, Durante, Barra, Agliari, submitted]







[Aquaro, Alemanno, Kanter, Durante, Barra, Agliari, submitted]





Why should we ever wake up?

[Kleinfeld, Pendergraft - Biophys. J. '87][van Hemmen, loffe, Kühn, Vaas - Physica A '90][Benedetti, Ventura, Marinari, Ruocco, Zamponi - J. Chem. Phys. '22]

A qualitative picture

 $K \times M$ "blurred" examples $\{\eta^{\mu,a}\}_{\mu=1,\dots,K}^{a=1,\dots,M}$ generated from a set of K (unknown) archetypes $\{\xi^{\mu}\}_{\mu=1,\dots,K}$



An analytical formulation

Look for a set of couplings $\{J_{ij}\}_{i,j=1,...,N}$ s.t. any pattern ξ^{μ} is a fixed point, for $\mu = 1,...,K$ $\xi_i^{\mu} = \operatorname{sign}\left(\sum_j J_{ij}\xi_j^{\mu} + h_i\right)$

with T = 0 and $h_i \in \mathbb{R}$ are external fields.

Stronger condition to guarantee a finite basin of attraction

$$\xi_i^{\mu} \left(\sum_j J_{ij} \xi_j^{\mu} + h_i \right) > \kappa \to \sum_j J_{ij} \xi_j^{\mu} + h_i = \gamma \xi_i^{\mu}, \ \gamma \ge \kappa,$$

with κ , γ finite constants.

[Gardner - J Phys A '86; Gardner, Derrida J Phys A '86] [Personnaz, Guyon, Dreyfus - J Phys Lett '85]

$$\mathscr{L}_{\boldsymbol{\xi}}(\boldsymbol{J},\boldsymbol{h}) = \frac{1}{2N} \sum_{i,\mu} \left(\sum_{j} \xi_{j}^{\mu} J_{ij} + h_{i} - \gamma \xi_{i}^{\mu} \right)^{2} + \underbrace{\frac{1}{2N} \sum_{j,\mu} \left(\sum_{i} \xi_{i}^{\mu} J_{ij} + h_{j} - \gamma \xi_{j}^{\mu} \right)^{2}}_{\text{symmetrize}} + \underbrace{\epsilon_{J} \sum_{i,j} J_{ij}^{2}}_{\text{regularize}} + \underbrace{\epsilon_{h} \sum_{i} h_{i}^{2}}_{\text{regularize}} \underbrace{\epsilon_{h} \sum_{i} h_{i}^{2}}_{\text{regularize}}$$

$$\begin{aligned} \boldsymbol{J}_{\epsilon_{J}=t^{-1}}(\tau \to \infty) &\to \boldsymbol{J}_{t}^{(DHN)} \\ \tau^{*}(\epsilon_{J}) &= \operatorname*{argmin}_{\tau} \| \boldsymbol{J}_{\epsilon_{J}=0,\gamma}(\tau) - \boldsymbol{J}_{\epsilon_{J},\gamma}(\tau \to \infty) \| \\ &\to \tau^{*} = A \log(1 + tB) + C \end{aligned}$$

[Girosi, Jones, Poggio - Neur. Com. '95] [Zhang, Yu - Ann. Stat. '05]



Hopfield network / Restricted Boltzmann machines equivalence



- [Tubiana, Monasson Rhys. Rev. Lett. '17]
- [Barra, Genovese, Sollich, Tantari Rhys. Rev. E '18]
- [Decelle, Fissore, Furtlehner J. Stat. Phys. '18]
- [Smart, Zilman ICLR '21]

learning with labels







$$\mathcal{S} = \{ \boldsymbol{\eta}^{\mu,a} \}_{a=1,\dots,M}^{\mu=1,\dots,K}, \, \boldsymbol{\eta}^{\mu,a} \to z_{\nu}^{(\mu)} = \delta_{\mu,\nu}$$

$$q(\boldsymbol{\sigma}, \boldsymbol{x}) = \sum_{\mu,a} \delta_{\boldsymbol{\eta}^{\mu,a},\boldsymbol{\sigma}} \delta_{\boldsymbol{z}^{(\mu)},\boldsymbol{z}}$$

$$\rho(\boldsymbol{\sigma}, \boldsymbol{z}; \boldsymbol{W}) \propto \exp\left[-\frac{\beta z_{\mu}^{2}}{2} + \frac{\beta}{\sqrt{N}} \sum_{i,\mu} \sigma_{i} W_{i}^{\mu} z_{\mu} \right]$$

$$\Delta W^{\mu} \propto \left(\langle \boldsymbol{\sigma}, \boldsymbol{z} \rangle \right) = \langle \boldsymbol{\sigma}, \boldsymbol{z} \rangle = \langle \boldsymbol{\sigma}, \boldsymbol{z} \rangle$$

$$\Delta W_i^{\mu} \propto \left(\langle \sigma_i z_{\mu} \rangle_{\text{clamped}} - \langle \sigma_i z_{\mu} \rangle_{\text{free}} \right), \quad \forall (i, \mu) \in (N \times K)$$

 $W = \bar{\eta}$ For structureless datasets: fixed point For structured datasets: effective initialisations

[Agliari, Leonelli, Marullo - Appl. Math. Comput. '22]

Hopfield network / Restricted Boltzmann machines equivalence







[Agliari, Alemanno, Aquaro, Barra, Kanter - Europhys. Lett. Perspective '23]



34/35

"Dreaming" Boltzmann machine

BM with $\sigma_i \in \{-1, +1\}, z_{\mu} \sim \mathcal{N}[0, T(1+t)]$

$$\mathscr{H}_{N,K,W,t}^{(\mathrm{DBM})}(\boldsymbol{\sigma},\boldsymbol{z}) = -\frac{1}{\sqrt{N}} \sum_{i,\mu}^{N,K} W_{i\mu} \sigma_i z_{\mu} + \frac{t}{1+t} \sum_{\mu<\nu}^{K,K} C_{\mu\nu} z_{\mu} z_{\nu},$$

equivalent to "Dreaming" Hopfield network as $W_{i\mu} = \xi_i^{\mu}$ for any i, μ

$$\mathscr{Z}_{\beta,N,K,\xi,t}^{(\text{DBM})} = \sum_{\sigma} \int \prod_{\mu} dz_{\mu} e^{-\frac{\beta z_{\mu}^{2}}{2(1+t)}} e^{-\beta \mathscr{H}_{N,K,\xi,t}^{(\text{DBM})}(\sigma,z)} = \sum_{\sigma} e^{-\beta \mathscr{H}_{N,K,\xi,t}^{(\text{DHN})}(\sigma)} = \mathscr{Z}_{\beta,N,K,\xi,t}^{(\text{DHN})}$$

Gradient descent over KL divergence

$$\Delta W_{j\rho} = \epsilon \beta \Big[\langle \sigma_j z_\rho \rangle_+ - \langle \sigma_j z_\rho \rangle_- - \frac{N}{2} \frac{t}{1+t} \sum_{\mu\nu} \frac{\partial C_{\mu\nu}}{\partial W_{j\rho}} \big(\langle z_\mu z_\nu \rangle_+ - \langle z_\mu z_\nu \rangle_- \big) \Big]$$

$$\Delta t = \epsilon \beta \frac{N}{2(1+t)^2} \sum_{\mu\nu} \big(1 - C)_{\mu\nu} (\langle z_\mu z_\nu \rangle_+ - \langle z_\mu z_\nu \rangle_- \big)$$

 $\langle \cdot \rangle_+$ clamped average $\langle \cdot \rangle_-$ free average

[Agliari et al. Neur. Netw. '21, IEEE Trans Neural Netw Learn Syst. '22]






Regularisation and Early-Stopping

$$\frac{dJ}{d\tau} = -\nabla_J \mathscr{L}_{\xi}(J, h)$$
$$\frac{dh}{d\tau} = -\nabla_h \mathscr{L}_{\xi}(J, h)$$

Natural basis: eigenvectors of $\hat{C} = \hat{\xi}\hat{\xi}^T \in \mathbb{R}^{N \times N}$

$$J_{ab}(\tau) \sim_{\tau \gg 1} \begin{cases} 0 & \text{if} \quad a \neq b \\ \frac{\lambda_a(1+\epsilon_J)}{\lambda_a+\epsilon_J} & \text{if} \quad a = b \end{cases} \qquad \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \qquad \begin{pmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \lambda_3 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$
$$\epsilon_J \to \infty \ (t \to 0)$$
$$J^{(KS)} = \xi^T C^{-1} \xi \qquad J^{(Hebb)} = \xi^T \xi$$

$$J_{aa}(\tau) = 1 - (1 - \lambda_a)e^{-2\tau\lambda_a}$$

Finite stopping-time τ^* corresponding to λ^*



$$\lambda_1 \leq \lambda_2 \ldots \leq \lambda_{N-\ell+1} \leq \ldots \leq \lambda_N$$

 ℓ fastest

$$\lambda_{N-K+1} = \frac{MK}{N} [r^2 (1 - \alpha^{-1/2})^2 + (1 - r^2)]$$

with K, M, r determinable by \hat{C} 's moments

Numerical experiments



Numerical experiments

Dataset $K \times M$ "blurred" examples $\{\eta^{\mu,a}\}_{\mu=1,...,K}^{a=1,...,M}$

Success: if $m_{\mu,a} > r$ and $m_{\mu} > n_{\mu,a}$

Overfit: elseif $n_{\mu,a} > r^2$ and $n_{\mu,a} > m_{\mu}$

Failure: else



40/35













Learning Hebb's rule

Hopfield network and blurred dataset

Does Hebbian learning work also for imperfect (unlabelled) training data?

K "archetypes" $\{\boldsymbol{\xi}^{\mu}\}_{\mu=1,...,K}$ *K* × *M* "blurred" examples $\{\boldsymbol{\eta}^{\mu,a}\}_{\mu=1,...,K}^{a=1,...,M}$

$$\eta_i^{\mu,a} = \xi_i^{\mu} \chi_i^{\mu,a} \qquad \text{with } \mathcal{P}(\chi_i^{\mu,a} = +1) = 1 - \mathcal{P}(\chi_i^{\mu,a} = -1) = p \in [1/2, 1]$$

М

Hebb's rule
$$J_{ij} = \frac{1}{N} \sum_{\mu=1}^{K} \sum_{a=1}^{M} \eta_i^{\mu,a} \eta_j^{\mu,a}$$

 $\mathscr{H}_{N,K,M}(\boldsymbol{\sigma}; \boldsymbol{\chi}, \boldsymbol{\xi}) = -\frac{1}{2N} \sum_{a=1}^{M} \sum_{\mu=1}^{K} \left(\sum_{i=1}^{N} \xi_i^{\mu} \chi_i^{\mu,a} \sigma_i\right)^2.$

Additional parameters

$$r := \langle \chi \rangle = (2p - 1)$$
 sample quality

sample size

Stat-mech solution

$$\mathcal{H}_{N,K,M}(\boldsymbol{\sigma};\boldsymbol{\chi},\boldsymbol{\xi}) = -\frac{1}{2N} \sum_{a=1}^{M} \sum_{\mu=1}^{K} \left(\sum_{i=1}^{N} \xi_{i}^{\mu} \chi_{i}^{\mu,a} \sigma_{i} \right)^{2}.$$

<u>Def.</u> Intensive quenched free-energy

$$f_{N,K,M}(\beta) := -\frac{1}{N\beta} \mathbb{E} \ln Z_{N,K,M}(\beta; \boldsymbol{\chi}, \boldsymbol{\xi}),$$

where
$$\mathbb{E} := \mathbb{E}_{\chi} \mathbb{E}_{\xi}$$
 and $Z_{N,K,M}(\beta; \chi, \xi) = \sum_{\sigma}^{2^{N}} \exp\left[\frac{\beta}{2N} \sum_{a=1}^{M} \sum_{\mu=1}^{K} \left(\sum_{i=1}^{N} \xi_{i}^{\mu} \chi_{i}^{\mu,a} \sigma_{i}\right)^{2}\right].$

<u>Def.</u> Mattis magnetizations

$$n_{\mu,a} := \frac{1}{N} \sum_{i=1}^{N} \xi_i^{\mu} \chi_i^{\mu,a} \sigma_i, \quad \mu = 1, \dots, K, \ a = 1, \dots, M; \qquad \qquad m_{\mu} := \frac{1}{N} \sum_{i=1}^{N} \xi_i^{\mu} \sigma_i, \quad \mu = 1, \dots, K.$$

<u>Def.</u> Two-replica overlap

$$q_{12} := \frac{1}{N} \sum_{i=1}^{N} \sigma_i^{(1)} \sigma_i^{(2)} \qquad \qquad p_{12} := \frac{1}{K} \sum_{\mu=1}^{K} z_{\mu}^{(1)} z_{\mu}^{(2)}$$

<u>Hp.</u> Replica symmetry & Thermodynamic limit RS $n_{1,a}$, q_{12} , p_{12} self-average at \bar{n} , \bar{q} , \bar{p} , respectively

M fixed, finite; $K, N \to \infty$ s.t. $\alpha := \lim_{N \to \infty} K/N$ finite

As $M \gg 1$, recalling r = (2p - 1)

$$\bar{n} = \frac{\bar{m}r}{1 - \beta(1 - \bar{q})(1 - r^2)}$$

Being $Z \sim \mathcal{N}(0,1)$, under noise rescaling $\beta \rightarrow \frac{\beta}{r^2 + \beta(1-\bar{q})(1-r^2)}$

$$\bar{m} = \mathbb{E}_{Z} \tanh\left[\beta\bar{m}M + Z\beta\sqrt{M\frac{1-r^{2}}{r^{2}}\bar{m}^{2} + \frac{\alpha\bar{p}}{r^{4}\beta}M}\right]$$
$$\bar{q} = \mathbb{E}_{Z} \tanh^{2}\left[\beta\bar{m}M + Z\beta\sqrt{M\frac{1-r^{2}}{r^{2}}\bar{m}^{2} + \frac{\alpha\bar{p}}{r^{4}\beta}M}\right]$$

$$\bar{p} = \frac{\beta \bar{q}}{[1 - \beta(1 - \bar{q})]^2}$$



0.2 0.4 0.6 0.8

$$\bar{n} = \frac{\bar{m}r}{1 - \beta(1 - \bar{q})(1 - r^2)}$$
$$\bar{m} = \mathbb{E}_Z \tanh\left[\frac{\beta \bar{m}M + Z\beta \sqrt{M \frac{1 - r^2}{r^2} \bar{m}^2 + \frac{\alpha \bar{p}}{r^4 \beta}M}}{r^4 \beta}\right]$$

• Archetype stability

signal carried by \bar{m} and *two* sources of (slow) noise

i. pattern interferenceii. example interference

$$\rightarrow M > M_c \sim r^{-4}$$

• Archetype vs Example retrieval

$$\bar{m} > \bar{n} \rightarrow M > M_{\times} \sim \frac{1+r}{r}$$



$$S = \{ \boldsymbol{\eta}^{\mu,a} \}_{\mu=1,...,K}^{a=1,...,K} \quad \text{Training set}$$

$$\mathscr{Z} = \{ \boldsymbol{z}^{\nu} \}_{\mu=1,...,K} \quad \text{Set of one-hot vectors of length } K(z_{\mu}^{\nu} = \delta_{\mu,\nu})$$

$$\mathscr{H}_{N,K}(\boldsymbol{\sigma}, \boldsymbol{z}; \boldsymbol{w}) = -\frac{1}{\sqrt{N}} \sum_{\mu=1}^{K} \sum_{i=1}^{N} \sigma_{i} w_{i}^{\mu} z_{\mu}$$

$$\mathscr{P}_{N,K}(\boldsymbol{\sigma}, \boldsymbol{z}; \boldsymbol{w}) = \frac{1}{Z_{N,K}(\boldsymbol{w})} e^{-\beta \mathscr{H}_{N,K}(\boldsymbol{\sigma}, \boldsymbol{z}; \boldsymbol{w})}$$

 $\Delta w_i^{\mu} \propto \left(\langle \sigma_i z_{\mu} \rangle_{\text{clamped}} - \langle \sigma_i z_{\mu} \rangle_{\text{free}} \right), \quad \forall (i, \mu) \in (N \times K) \qquad \text{Supervised learning rule}$

``clamped" average is evaluated over the pairs $(\sigma_E, z_E) \in S \times \mathscr{Z}$, e.g. $\sigma_E = \eta^{\nu,1}, z_E = z^{\nu}$

``free" average is sampled via a single step of Gibbs dynamics: select (σ_E, z_E) $\in S \times \mathscr{Z}$, then using the Gibbs-chain $z_E \to \sigma_{\text{free}} \to z_{\text{free}}$ obtain the state ($\sigma_{\text{free}}, z_{\text{free}}$) <u>Classifier</u> $\boldsymbol{\sigma} = \boldsymbol{\xi}^{\nu}, \boldsymbol{\eta}^{\nu,a} \mapsto \boldsymbol{z}^{\nu}$

Performance measure: log $\left[\frac{\mathscr{P}(z_E | \boldsymbol{\sigma} = \boldsymbol{\xi}_E)}{\mathscr{P}(z_E | \boldsymbol{\sigma} = \boldsymbol{\eta}_E)}\right]$

Generative model
$$z^{\nu} \mapsto \sigma = \xi^{\nu}, \eta^{\nu,a}$$

Performance measure:
$$\log \frac{\langle m_{\nu} \rangle_{z^{\nu}}}{\langle n_{\nu} \rangle_{z^{\nu}}}$$

As training is running, as long as $M < M_{\times}$, the saturation values for $\mathscr{P}(z_E | \boldsymbol{\sigma} = \boldsymbol{\eta}_E)$ and $\langle n \rangle$ are larger than those obtained for $\mathscr{P}(z_E | \boldsymbol{\sigma} = \boldsymbol{\xi}_E)$ and $\langle m \rangle$; the opposite holds when $M > M_{\times}$



EA, F. Alemanno, A. Barra, G. Di Marzo, The emergence of a concept in shallow neural networks (2021)

Restricted Boltzmann machines



 $h \in \{-1, +1\}^{N_h}$

hidden layer

Two layers composed of $N = N_v + N_h$ binary neurons

 $s = (v, h) \in \{-1, +1\}^{N_v + N_h}$

 $\{J_{ij}\}$, with $i, j \in \{1, ..., N\}$ interaction matrix (symmetric, zero eye) $\{\vartheta_i\}$, with $i \in \{1,...,N\}$ bias vector

<u>Theor.</u> The partition function of the RR algorithm, given by

$$Z_{\beta,N,P,t}(\boldsymbol{\xi}) = \sum_{\{\boldsymbol{\sigma}\}} e^{-\beta H_{N,P,t}^{(RR)}(\boldsymbol{\sigma}|\boldsymbol{\xi})} = \sum_{\{\boldsymbol{\sigma}\}} \exp\left[\frac{\beta}{2N} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{P,P} \xi_i^{\mu} \xi_j^{\nu} \left(\frac{1+t}{\mathbb{I}+tC}\right)_{\mu,\nu} \sigma_i \sigma_j\right]$$

where $C_{\mu,\nu} := \frac{1}{N} \sum_{i=1}^{N} \xi_i^{\mu} \xi_i^{\nu}$

can be represented in Gaussian integral form as

$$Z_{\beta,N,P,t}(\boldsymbol{\xi}) = \sum_{\{\sigma\}} \int \Big(\prod_{\mu=1}^{P} d\mu(z_{\mu})\Big) \Big(\prod_{i=1}^{N} d\mu(\phi_{i})\Big) \exp\left(\sqrt{\frac{\beta}{N}(t+1)} \sum_{\mu,i}^{P,N} z_{\mu}\xi_{i}^{\mu}\sigma_{i} + i\sqrt{\frac{t}{N}} \sum_{\mu,i}^{P,N} z_{\mu}\xi_{i}^{\mu}\phi_{i}\right)$$

that is equivalent to the partition function of a tripartite spin-glass with intermediate layer made of a set of real neurons $\{z_{\mu}\}_{\mu=1,...,P}$ with $z_{\mu} \sim \mathcal{N}[0,1]$ external layers made, respectively, of a set of Boolean neurons $\{\sigma_i\}_{i=1,...,N}$ and of a set of imaginary neurons with magnitude $\{\phi_i\}_{i=1,...,N}$, being $\phi_i \sim \mathcal{N}[0,1]$.

Dreaming Hopfield network

Hopfield network with "wide" retrieval region

Maximise the retrieval region to enhance storing

For symmetric neural networks capacity upper bounded by $a_G = 1$ E. Gardner - J. Phys. A (1988)





Hopfield network with "wide" retrieval region

Maximise the retrieval region to enhance storing

For symmetric neural networks capacity upper bounded by $a_G = 1$ E. Gardner - J. Phys. A (1988)



Projection rule (decorrelation prescription) Kohonen - IEEE Trans. Comp. '72 Personnaz - J. Phys. Lett. '85 Kanter&Sompolinsky - Phys. Rev. A '86

$$\begin{split} J_{ij} &= \frac{1}{N} \sum_{\mu=1}^{K} \xi_{i}^{\mu} \xi_{j}^{\mu} \to \sum_{\mu,\nu=1}^{K} \xi_{i}^{\mu} C_{\mu\nu}^{-1} \xi_{j}^{\nu} \\ C_{\mu\nu} &:= N^{-1} \sum_{i=1}^{N} \xi_{i}^{\mu} \xi_{i}^{\nu} \end{split}$$

Interpreted as an iterative rule accounting for learning (patterns) and unlearning (spurious states) Hopfield et al. - Nature Lett. '83 Opper - Europhys. Lett. '89 Plakhov - IAPR '94 $J_{ij} = \frac{1}{N} \sum_{\nu=1}^{K} \xi_i^{\mu} \xi_j^{\mu} \rightarrow \left[J_{ij} - \epsilon \langle h_i h_j \rangle_{t'} \right]_{t'=0}^{t} - \sum_{\nu=1}^{K} \xi_i^{\mu} (I_K + tC)_{\mu,\nu}^{-1} \xi_j^{\nu}$

Reinforcement&Removal (RR) algorithm

$$\frac{dJ}{dt} = \frac{1}{1+t} (J - J^2)$$
$$J_{ij}(t) = \frac{1}{N} \sum_{\mu,\nu=1}^{K,K} \xi_i^{\mu} \xi_j^{\nu} \left(\frac{1+t}{I_K + tC}\right)_{\mu\nu}$$

with $C_{\mu\nu} = \frac{1}{N} \sum_{i=1}^{N} \xi_i^{\mu} \xi_j^{\nu}$ the pattern correlation matrix

$$\mathscr{H}_{N,K}^{(RR)}(\boldsymbol{\sigma};\boldsymbol{\xi},t) = -\frac{1}{2N} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{K,K} \xi_i^{\mu} \xi_j^{\nu} \left(\frac{1+t}{\boldsymbol{I}_K + t\boldsymbol{C}}\right)_{\mu,\nu} \sigma_i \sigma_j$$

Goal: thermodynamic limit of the quenched-averaged intensive pressure

$$f(\beta, \alpha, t) := -\lim_{N \uparrow \infty} \frac{1}{N\beta} \mathbb{E} \ln Z_{N,K}(\beta, \xi, t)$$

being
$$Z_{N,K}(\boldsymbol{\xi},t) = \sum_{\boldsymbol{\sigma}} e^{-\beta \mathcal{H}_{N,K}^{(RR)}(\boldsymbol{\sigma};\boldsymbol{\xi},t)}$$

A. Fachechi, EA, A. Barra, *Dreaming neural networks: forgetting spurious memories and reinforcing pure ones*, Neur. Netw. (2019)

$$-\beta f^{RS}(\beta, \alpha, t) = \log 2 - \frac{\beta m^2}{2(1+t)} \left(1 + \frac{t}{\Delta}\right) - \frac{(1+t)(\Delta - 1)}{2t} \beta \bar{q} + \mathbb{E}_{\eta} \log \cosh\left[\frac{\beta}{\Delta}(m + \sqrt{\alpha \bar{p}} \eta)\right] - \frac{\log \Delta}{2} - \frac{\alpha \beta \bar{p} t}{2(1+t)\Delta} - \frac{\alpha}{2} \left(\log[1 - \beta(1+t)(\bar{q} - \bar{q}')] + \frac{\bar{q}' \beta^2(1+t)}{1 - \beta(1+t)(\bar{q} - \bar{q}')}\right) - \frac{(1+t)(1-\Delta)\beta}{2t\Delta} - \frac{\alpha \beta^2}{2} \bar{p}(\bar{q} - \bar{q}').$$

$$\begin{split} m &= \frac{1+t}{\Delta+t} \mathbb{E}_{\eta} \tanh\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha\bar{p}}\eta)\right],\\ \bar{p} &= \frac{\bar{q}'(1+t)^2}{[1-\beta(1+t)(\bar{q}-\bar{q}')]^2},\\ \Delta &= 1+\frac{\alpha t}{1-\beta(1+t)(\bar{q}-\bar{q}')},\\ \bar{q}' &= \bar{q}+\frac{t}{\beta(1+t)\Delta} - \frac{1}{\Delta^2} \mathbb{E}_{\eta} \cosh^{-2}\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha\bar{p}}\eta)\right],\\ \bar{q}\Delta^2 &= 1-\frac{t\Delta}{\beta(1+t)} + \frac{\alpha\bar{p}t^2 - m^2t(t+2\Delta)}{(1+t)^2} - \frac{2\alpha\beta\bar{p}t}{(1+t)\Delta} \mathbb{E}_{\eta} \cosh^{-2}\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha\bar{p}}\eta)\right]. \end{split}$$

Critical surface delimiting the ergodic region

$$\beta_c = \frac{1}{1+t} \Big[\frac{\Delta^2}{1+\sqrt{\alpha}} + t\Delta \Big],$$

$$\Delta = 1 + \sqrt{\alpha} (1+\sqrt{\alpha})t.$$

Let the network... dream!

F. Crick, G. Mitchinson, *The function of dream sleep*, Nature (1983).

«We propose that the function of dream sleep (more properly rapid-eye movement or REM sleep) is to remove certain undesirable modes of interaction in networks of cells in the cerebral cortex. We postulate that this is done in REM sleep by a **reverse learning mechanism**, so that the trace in the brain of the unconscious dream is weakened, rather than strengthened, by the dream.»

REM sleeping is associated to the **removal** of unwanted attractors (*unlearning*)

J.J. Hopfield, D.I. Feinstein, R.G. Palmer, Unlearning has a stabilizing effect in collective memories, Nature (1983).

Unlearning rule aiming to increase the energy of non pure attractors (they get less stable)





Hopfield's unlearning succeeds in reducing the accessibility of spurious memories, yet...

Too much unlearning can also destroy pure memories
No analytical treatment

J.L. Van Hemmen, L.B. loffe, R. Kühn, M. Vaas, *Increasing the efficiency of a neural network through unlearning*, Phys. A (1990).

A. Y. Plakhov, S.A. Semenov, *The modified unlearning procedure for enhancing storage capacity in Hopfield network*, IEEE Trans. (1992).



In the limit of large t, the coupling matrix converges to the **projector matrix** $P \rightarrow$ Analytically treatable (but no longer local)

$$P_{ij} = \frac{1}{N} \sum_{\mu,\nu=1}^{P,P} \xi_i^{\mu} (C^{-1})_{\mu\nu} \xi_j^{\mu},$$

where $C_{\mu\nu} = \frac{1}{N} \sum_{i=1}^{N} \xi_i^{\mu} \xi_i^{\nu}$

The unlearning strength ϵ identifies a temporal scale

Taking the continuous time limit $(dt \sim \varepsilon)$, the unlearning procedure is described by a the differential equation

$$\frac{dJ}{dt} = -J^2 \quad \text{with solution} \quad J_{ij}(t) = \frac{1}{N} \sum_{\mu,\nu=1}^{P,P} \xi_i^{\mu} \xi_j^{\nu} (1+tC)_{\mu\nu}^{-1}$$

where $C_{\mu\nu} = \frac{1}{N} \sum_{i=1}^{N} \xi_i^{\mu} \xi_j^{\nu}$ is the pattern correlation matrix

V. Dotsenko, N.D. Yarunin, E.A. Dorotheyev, Statistical mechanics of Hopfield-like neural networks with modified interactions, J. Phys. A (1991).

Analytical solution of the model in the replica-symmetric regime by means of replica trick technique.



The critical capacity at zero temperature grows with *t*, but the area retrieval region vanishes in the large unlearning time limit

Critical load $(t = \infty) \approx 1.07$

V. Dotsenko, N.D. Yarunin, E.A. Dorotheyev, Statistical mechanics of Hopfield-like neural networks with modified interactions, J. Phys. A (1991).



Too much unlearning destroys pure memories With this unlearning rule all couplings tend to zero in the large unlearning-time limit

The area of the retrieval region vanishes in the large unlearning-time limit

Intuitively...

The unlearning algorithm does not distinguish between pure and spurious memories \rightarrow Too much unlearning destroys also the former.

Mathematically...

The generalized kernel in the coupling matrix equals the zero matrix in the large unlearning time limit \rightarrow All synaptic strengths vanish.

Biologically...

The function of sleeping (REM sleep + <u>Slow Wave Sleep</u>) is to weaken fictitious synaptic strengths *and* <u>consolidate relevant ones</u> (Walker 2009, Born 2010...).

A. Fachechi, E. Agliari, A. Barra, *Dreaming neural networks: forgetting spurious memories and reinforcing pure ones*, Neur. Net. (2019).

Introduce a regularizator

$$(1+tC)_{\mu\nu}^{-1} \to \left(\frac{1+t}{1+tC}\right)_{\mu\nu}$$
$$J_{ij}(t) = \frac{1}{N} \sum_{\mu,\nu=1}^{P,P} \xi_i^{\mu} \xi_j^{\nu} \left(\frac{1+t}{1+tC}\right)_{\mu\nu}$$

The resulting model accounts for both **removal and consolidation**!

<u>Argument 1</u>: Look at the evolution process leading to J(t)

 $\frac{dJ}{dt} = \underbrace{1}_{1+t}(J - J^2)$

Or, in discrete sleeping time...

$$\begin{split} dt &\to \epsilon, \\ t &\to k\epsilon, \\ \dot{J} &\to \epsilon^{-1} [J(k+1) - J(k)], \end{split}$$

$$J(k+1) = J(k) + \frac{\epsilon}{1+\epsilon k} \left[J(k) - J(k)^2 \right]$$

effective sleep strength depends on the sleep session k

61 /19

$$J(k+1) = J(k) + \frac{\epsilon}{1+\epsilon k} \left[J(k) - J(k)^2 \right]$$

with this prescription $||\mathbf{J}(k) - \mathbf{P}|| \xrightarrow{k \to \infty} 0$ under operator norm

Sketch of the proof

Introduce a matrix **G** such that $J_{ij}(k) = \frac{1}{N} \sum_{\mu,\nu=1}^{P} \xi_i^{\mu} \xi_j^{\nu} G_{\mu\nu}(k)$

then, prescription recast as $\mathbf{G}(k+1) = \left(1 + \frac{\epsilon}{1+\epsilon k}\right)\mathbf{G}(k) - \frac{\epsilon}{1+\epsilon k}\mathbf{G}(k)\mathbf{CG}(k)$, with initial condition $G(0) = \mathbb{I}$

- ||*C*|| ≥ 1
- G(k) commutes with $C \forall k$
- G(k) are invertible for $k \ge 0$ and $\varepsilon > \varepsilon_c$
- There exists finite, real $c_k \ge ||C G(k)||$. The sequence is therefore upper bounded by $\bar{c}=\max_k c_k$
- The unlearning algorithm converges to the stationary solution G(∞) = C⁻¹ in the sense defined by the operator norm
- Critical value for ε to ensure convergence $\varepsilon_c = (||C||-1)^{-1}$



Projector (or pseudo-inverse) $P_{ij} = \frac{1}{N} \sum_{\mu,\nu=1}^{P,P} \xi_i^{\mu} (C^{-1})_{\mu\nu} \xi_j^{\mu},$ where $C_{\mu\nu} = \frac{1}{N} \sum_{i=1}^{N} \xi_i^{\mu} \xi_i^{\nu}$ Operator norm $||\mathbf{A}|| = \sqrt{\max\{a | a \in \sigma(\mathbf{A}^T \mathbf{A})\}}$



 ε_c higher than Plakhov et al.'s and decreases slower with N

A. Fachechi, E. Agliari, A. Barra, *Dreaming neural networks: forgetting spurious memories and reinforcing pure ones*, Neur. Net. (2019).

Introduce a regularizator

$$(1+tC)_{\mu\nu}^{-1} \to \left(\frac{1+t}{1+tC}\right)_{\mu\nu}$$
$$J_{ij}(t) = \frac{1}{N} \sum_{\mu,\nu=1}^{P,P} \xi_i^{\mu} \xi_j^{\nu} \left(\frac{1+t}{1+tC}\right)_{\mu\nu}$$

The resulting model accounts for both **removal and consolidation**!

<u>Argument 1</u>: Look at the evolution process leading to J(t)

 $\frac{dJ}{dt} = \underbrace{1}_{1+t}(J - J^2)$

Or, in discrete sleeping time...

$$\begin{split} dt &\to \epsilon, \\ t &\to k\epsilon, \\ \dot{J} &\to \epsilon^{-1} [J(k+1) - J(k)], \end{split}$$

$$J(k+1) = J(k) + \frac{\epsilon}{1+\epsilon k} \left[J(k) - J(k)^2 \right]$$

effective sleep strength depends on the sleep session k

63 /19

A. Fachechi, E. Agliari, A. Barra, *Dreaming neural networks: forgetting spurious memories and reinforcing pure ones*, Neur. Net. (2019).

Introduce a regularizator

$$(1+tC)_{\mu\nu}^{-1} \to \left(\frac{1+t}{1+tC}\right)_{\mu\nu}$$
$$J_{ij}(t) = \frac{1}{N} \sum_{\mu,\nu=1}^{P,P} \xi_i^{\mu} \xi_j^{\nu} \left(\frac{1+t}{1+tC}\right)_{\mu\nu}$$

The resulting model is accounts for both **removal and consolidation**!

<u>Argument 2</u>: consider the statistical mechanics of two separated models



The model is analytically treatable, even in the high-load regime!

Self-consistent equations for the order parameters obtained (rigorously) under replicasymmetry hypothesis.

Order parameters

 $m \rightarrow$ assess retrieval of pattern ξ^1 (measure overlap between neural configuration and pattern)

 $q, p \rightarrow$ assess glassiness in the system (measure overlap between neural configurations of two replicas)

Replica symmetry ansatz

 $m_1^{\alpha} = m \quad \forall \alpha,$ $a_{\alpha\beta} = O\delta_{\alpha\beta} + q(1 - \delta_{\alpha\beta}),$

$$q_{\alpha\beta} = Q\delta_{\alpha\beta} + q(1 - \delta_{\alpha\beta}),$$

$$p_{\alpha\beta} = P\delta_{\alpha\beta} + p(1-\delta_{\alpha\beta}),$$

Useful quantity

$$\Delta := 1 + \alpha \beta t (1+t)^{-1} (P-p)$$

The model is analytically treatable, even in the high-load regime!

Self-consistent equations for the order parameters obtained (rigorously) under replicasymmetry hypothesis.

$$\begin{split} m &= \frac{1+t}{\Delta+t} \int Dx \tanh\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha p}x)\right] \\ p &= \frac{q(1+t)^2}{[1-\beta(1+t)(Q-q)]^2} \\ \Delta &= 1 + \frac{\alpha t}{1-\beta(1+t)(Q-q)} \\ q &= Q + \frac{t}{\beta(1+t)\Delta} - \frac{1}{\Delta^2} \int Dx \cosh^{-2}\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha p}x)\right] \\ Q\Delta^2 &= 1 - \frac{t\Delta}{\beta(1+t)} + \frac{\alpha p t^2}{(1+t)^2} - \frac{m^2 t(t+2\Delta)}{(1+t)^2} - \frac{2\alpha\beta p t}{(1+t)\Delta} \int Dx \cosh^{-2}\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha p}x)\right] \end{split}$$

By evaluating these observables as functions of the parameters (T, a) we build the phase diagram, focusing on the retrieval region



Mattis magnetization for t = 1000 as a function of temperature, for various storage capacity values ($\alpha = 0, 0.05, 0.2$ and 0.5)

E. Agliari, F. Alemanno, A. Barra, A. Fachechi, *Dreaming neural networks: rigorous results*, J. Stat. Phys. (2019)

66 /19



Retrieval region gets wider

Spin-glass (where slow noise, i.e., spurious states, prevail) region collapses

Ergodic (where fast noise, i.e., purely random states, prevail) region gets wider

Ergodic line changes concavity

Perfect retrieval is accomplished!

A technical note: the underlying mathematics is rigorous (at replica-symmetry resolution)



Comparison with numerics

Check hp underlying analytics (thermodynamic limit, replica symmetry)



Average values for the Mattis magnetization m_1 corresponding to the retrieved pattern ξ^1 obtained from numerical simulations for fixed a=0.08 and M=10



More techanical details

Definition 1 The Hamiltonian of the reinforcement&removal algorithm reads as:

$$H_{N,P}^{(RR)}(\sigma|\xi,t) = -\frac{1}{2N} \sum_{i=1}^{N} \sum_{j=1}^{N} \sum_{\mu=1}^{P} \sum_{\nu=1}^{P} \xi_{i}^{\mu} \xi_{j}^{\nu} \left(\frac{1+t}{\mathbb{I}+tC}\right)_{\mu,\nu} \sigma_{i}\sigma_{j},\tag{1}$$

where ξ^1 , the pattern candidate to be retrieved, has binary entries $\xi^1_i \in \{-1, +1\}$ drawn from $P(\xi^{\mu}_i = +1) = P(\xi^{\mu}_i = -1) = \frac{1}{2}$, while the remaining P - 1 patterns $\{\xi^{\mu}\}_{\mu=2,...,P}$, have i.i.d. standard Gaussian entries $\xi^{\mu}_i \sim \mathcal{N}[0,1]$, and the correlation matrix C is defined as

$$C_{\mu,\nu} := \frac{1}{N} \sum_{i=1}^{N} \xi_i^{\mu} \xi_i^{\nu}$$

Definition 2 Being $\beta \in \mathbb{R}^+$ a parameter tuning the level of fast noise in the network (with the physical meaning of inverse temperature), the partition function of the model (1) is introduced as

$$Z_{N,P}(\sigma|\xi,t) = \sum_{\{\sigma\}} e^{-\beta H_{N,P}^{(RR)}(\sigma|\xi,t)} = \sum_{\{\sigma\}} \exp\left[\frac{\beta}{2N} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{P,P} \xi_i^{\mu} \xi_j^{\nu} \left(\frac{1+t}{\mathbb{I}+tC}\right)_{\mu,\nu} \sigma_i \sigma_j\right].$$
 (2)

Definition 3 Denoting with \mathbb{E}_{ξ} the average over the quenched patterns, for a generic function $O(\sigma,\xi)$ of the neurons and the couplings, we can define the Boltzmann $\omega(O(\sigma,\xi))$ and the quenched $\langle O(\sigma,\xi) \rangle$ averages as

$$\omega\left(O(\sigma,\xi)\right) = \frac{\sum_{\{\sigma\}} O(\sigma,\xi) e^{-\beta H_{N,P}^{(RR)}(\sigma|\xi,t)}}{Z_{N,P}(\sigma|\xi,t)},$$
$$\left\langle O(\sigma,\xi)\right\rangle = \mathbb{E}_{\xi} \omega\left(O(\sigma,\xi)\right).$$

Definition 4 Once introduced the partition function $Z_{N,P}(\sigma|\xi,t)$, we can define the infinite volume limit of the intensive quenched free-energy $F_N(\alpha,\beta,t)$ and of the intensive quenched pressure $A(\alpha,\beta,t)$ associated to the model (1) as

$$-\beta F(\alpha,\beta,t) = A(\alpha,\beta,t) = \lim_{N \to \infty} \frac{1}{N} \mathbb{E} \ln Z_{N,P}(\sigma|\xi,t).$$

Remark 1 For the sake of mathematical convenience, we take solely the pattern candidate for retrieval (i.e. the signal) to be Boolean, while all the remaining ones (acting as slow noise on the retrieval) are chosen as Gaussian. "Universality" as long as the distribution function of the quenched patterns generating the slow noise satisfies

$$\exists L > 0 \mid p > 2 \implies |\mathbb{E}_{\xi}\xi^p| < L, \mathbb{E}_{\xi}\xi = 0, \mathbb{E}_{\xi}\xi^2 = 1.$$

Lemma 1 The partition function (2) can be represented in Gaussian integral form as

$$Z_{N,P}(\sigma|\xi,t) = \sum_{\{\sigma\}} \int \left(\prod_{\mu=1}^{P} d\mu(z_{\mu})\right) \left(\prod_{i=1}^{N} d\mu(\phi_{i})\right) \exp\left(\sqrt{\frac{\beta}{N}(t+1)} \sum_{\mu,i}^{P,N} z_{\mu}\xi_{i}^{\mu}\sigma_{i} + i\sqrt{\frac{t}{N}} \sum_{\mu,i}^{P,N} z_{\mu}\xi_{i}^{\mu}\phi_{i}\right), \quad (3)$$

that is equivalent to the partition function of a tripartite spin-glass where the intermediate layer is made of a set of real neurons $\{z_{\mu}\}_{\mu=1,...,N}$ with $z_{\mu} \sim \mathcal{N}[0,1], \forall \mu$, while the external layers are made, respectively, of a set of Boolean neurons $\{\sigma_i\}_{i=1,...,N}$ and of a set of imaginary neurons with magnitude $\{\phi\}_{i=1,...,N}$, being $\phi_i \sim \mathcal{N}[0,1], \forall i$.

Definition 5 Once expressed the partition function (2) in its integral representation (3), we can introduce the related tripartite spin-glass Hamiltonian as

$$H_{N,P} = \frac{a}{\sqrt{N}} \sum_{i=1}^{N} \sum_{\mu=1}^{P} z_{\mu} \xi_{i}^{\mu} k_{i}, \qquad (4)$$

where we introduced the "multi-spin" $k_i = \sigma_i + b\phi_i$ and where

$$a = \sqrt{\beta(t+1)}, \quad b = i\sqrt{\frac{t}{\beta(t+1)}}.$$
(5)

The cost function (4) and the one associated to the original model (1) share the same partition function and therefore exhibit the same thermodynamic.

Definition 6 The natural order parameters for the neural network model (1) are the overlaps q_{ab} and p_{ab} between the k's and the z's variables, respectively, as functions of two replicas (a,b) of the system, and the generalized Mattis overlap m_1 , namely

$$q_{ab} := \frac{1}{N} \sum_{i=1}^{N} k_i^{(a)} k_i^{(b)},$$

$$p_{ab} := \frac{1}{P} \sum_{\mu \ge 2} z_{\mu}^{(a)} z_{\mu}^{(b)},$$

$$m_1 := \frac{1}{N} \sum_{i=1}^{N} \xi_i^1 k_i.$$

Remark 2 RS approximation \rightarrow order-parameters do not fluctuate in the thermodynamic limit,

$$\begin{array}{ll} q_{ab} & \stackrel{\text{RS}}{\to} & W\delta_{ab} + q(1 - \delta_{ab}), \\ p_{ab} & \stackrel{\text{RS}}{\to} & X\delta_{ab} + p(1 - \delta_{ab}), \\ m_1 & \stackrel{\text{RS}}{\to} & m. \end{array}$$

71 /19

Theorem 1 In the infinite volume limit, the replica symmetric free energy related to the neural network (1) can be expressed in terms of the natural order parameters of the theory as

$$A_{RS}(\alpha,\beta,t) = \log 2 - \frac{\beta m^2}{2(1+t)} \left(1 + \frac{t}{\Delta}\right) - \frac{(1+t)(\Delta-1)}{2t} \beta W + \mathbb{E}_{\eta} \log \cosh\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha p}\eta)\right] - \frac{\log \Delta}{2} - \frac{\alpha\beta pt}{2(1+t)\Delta} - \frac{\alpha}{2} \left(\log[1-\beta(1+t)(W-q)] + \frac{q\beta^2(1+t)}{1-\beta(1+t)(W-q)}\right) - \frac{(1+t)(1-\Delta)\beta}{2t\Delta} - \frac{\alpha\beta^2}{2}p(W-q).$$
(6)

Proposition 1 Using the standard variational principle $\vec{\nabla}A_{RS} = 0$ on the free energy (6), namely by extremizing the latter over the order parameters, we obtain the following set of self-consistent equations

$$\begin{split} m &= \frac{1+t}{\Delta+t} \mathbb{E}_{\eta} \tanh\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha p}\eta)\right], \\ p &= \frac{q(1+t)^2}{[1-\beta(1+t)(W-q)]^2}, \\ \Delta &= 1+\frac{\alpha t}{1-\beta(1+t)(W-q)}, \\ q &= W+\frac{t}{\beta(1+t)\Delta} - \frac{1}{\Delta^2} \mathbb{E}_{\eta} \cosh^{-2}\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha p}\eta)\right], \\ W\Delta^2 &= 1-\frac{t\Delta}{\beta(1+t)} + \frac{\alpha p t^2 - m^2 t (t+2\Delta)}{(1+t)^2} - \frac{2\alpha\beta p t}{(1+t)\Delta} \mathbb{E}_{\eta} \cosh^{-2}\left[\frac{\beta}{\Delta}(m+\sqrt{\alpha p}\eta)\right]. \end{split}$$

The proof is based on Guerra's interpolating techniques...
Definition 6 Being $s \in [0,1]$ an interpolating parameter, $\{\eta_i\}_{i \in (1,...,N)}$ a set of N i.i.d. Gaussian variables, $\{\lambda_\mu\}_{\mu \in (2,...,P)}$ a set of P-1 i.i.d. Gaussian variables, and the scalars C_1, C_2, C_3, C_4, C_5 to be set a posteriori, we use as interpolating pressure the following quantity

$$\mathcal{A}(s) = \frac{1}{N} \mathbb{E}_{\xi,\eta,\lambda} \ln \sum_{\sigma} \int d\mu(z,\phi) \exp\left[\sqrt{s} \frac{a}{\sqrt{N}} \sum_{i,\mu \ge 2} z_{\mu} \xi_{i}^{\mu} k_{i} + \sqrt{s} \frac{a}{\sqrt{N}} \sum_{i} z_{1} \xi_{i}^{1} k_{i} \right.$$
(7)
+ $\sqrt{1-s} \left(C_{1} \sum_{i}^{N} \eta_{i} k_{i} + C_{2} \sum_{\mu \ge 2} \lambda_{\mu} z_{\mu} \right) + \frac{1-s}{2} \left(C_{3} \sum_{\mu \ge 2} z_{\mu}^{2} + C_{4} \sum_{i} k_{i}^{2} + C_{5} a \sum_{i} \xi_{i}^{1} k_{i} \right) \right].$

When s = 1 we recover the original model, namely $A(\alpha, \beta, t) = \lim_{N \to \infty} \mathcal{A}(s = 1)$, while for $s \to 0$ we are left with a one-body problem, and, consequently, the probabilistic structure of $\mathcal{A}(s = 0)$ is more tractable.

Proposition 2 The infinite volume limit of the quenched free energy related to the model (1) can be obtained by using the Fundamental Theorem of Calculus as

$$A(\alpha, \beta, t) = \lim_{N \to \infty} \mathcal{A}(s=1) = \lim_{N \to \infty} \left(\mathcal{A}(s=0) + \int_0^1 \frac{d\mathcal{A}(s)}{ds} \, ds \right).$$

Performing calculations and setting the free scalars $C_{1,..,5}$ as

$$C_1^2 = a^2 \alpha p, \quad C_2^2 = a^2 q, \quad C_3 = a^2 (W - q), \quad C_4 = a^2 \alpha (X - p), \quad C_5 = 2ma,$$
(8)

one gets that under RS $d_s \mathcal{A}(s)$ is independent of s and calculations straightforwardly yields to we finally get (6).

Theorem 2 The ergodic region of the model defined by the cost function (??) is delimited by the following critical surface in the (α, β, t) space of the tunable parameters

$$\beta_c = \frac{1}{1+t} \Big[\frac{\Delta^2}{1+\sqrt{\alpha}} + t\Delta \Big],$$

$$\Delta = 1 + \sqrt{\alpha} (1+\sqrt{\alpha})t.$$

Remark 3 At t = 0, where the model reduces to Hopfield's scenario, the critical surface correctly collapses over the Amit-Gutfreund-Sompolinsky critical line $\beta_c = (1 + \sqrt{\alpha})^{-1}$, but in the large t limit the ergodic region collapses to the axis T = 0: we speculate that this may have a profound implications, namely that the ergodic region -during the sleep state- invades the spin-glass region, the latter being de facto suppressed.

The proof is again based on Guerra's interpolating techniques...

Definition 7 The centered and rescaled overlap fluctuations θ_{lm} and ρ_{lm} are introduced as

$$\theta_{lm} = \sqrt{N} [q_{lm} - \delta_{lm} W - (1 - \delta_{lm})q]$$

$$\rho_{lm} = \sqrt{P} [p_{lm} - \delta_{lm} X - (1 - \delta_{lm})p]$$

Remark 4 As we will address the problem of the overlap fluctuations in the ergodic region, the signal is absent, thus there is no need to introduce a rescaled Mattis order parameter.

Proposition 2 We introduce the r-replicated interpolating free energy $\mathcal{A}_{J}^{r}(s)$, where we further added a source field J, coupled to an observable O (that is a smooth function of the neurons of the r-replicas) as

$$\begin{aligned} \mathcal{A}_{J}^{r}(s) &= \mathbb{E}_{\xi,\eta,\lambda} \ln \sum_{\sigma_{R}} \int d\mu(z_{R},\phi_{R}) \exp\left[\sqrt{s} \frac{a}{\sqrt{N}} \sum_{l=1}^{r} \sum_{i,\mu} z_{\mu}^{(l)} \xi_{i}^{\mu} k_{i}^{(l)} + J\hat{O} \right. \end{aligned} \tag{9} \\ &+ \sqrt{1-s} \Big(C_{1} \sum_{l=1}^{r} \sum_{i} \eta_{i} k_{i}^{(l)} + C_{2} \sum_{l=1}^{r} \sum_{\mu} \lambda_{\mu} z_{\mu}^{(l)} \Big) + \frac{1-s}{2} \Big(C_{3} \sum_{l=1}^{r} \sum_{\mu} (z_{\mu}^{(l)})^{2} + C_{4} \sum_{l=1}^{r} \sum_{i} (k_{i}^{(l)})^{2} \Big) \Big]. \tag{9}$$

where k_i and the interpolation constants $C_{1,2,3,4}$ are the same given before.

One can verify that

$$\langle O(s) \rangle = \left. \frac{\partial \mathcal{A}_J^r(s)}{\partial J} \right|_{J=0}, \qquad \partial_s \langle O(s) \rangle = \left. \frac{\partial (\partial_s \mathcal{A}_J^r)}{\partial J} \right|_{J=0}.$$
 (11)

Therefore, in order to evaluate the fluctuations of O we need to evaluate first $\partial_s \mathcal{A}_J^r$ and, by a routine calculation, we get

$$\partial_s \mathcal{A}_J^r = \frac{1}{2} \sqrt{\alpha} \beta (1+t) \sum_{l,m=1}^r \left[\langle g_{l,m} \rangle - \langle g_{l,m+r} \rangle \right], \qquad g_{l,m} = \theta_{lm} \rho_{lm}.$$
(12)

Proposition 3 Given O as a smooth function of r replica overlaps (q_1, \ldots, q_r) and (p_1, \ldots, p_r) , the following streaming equation holds:

$$d_{\tau}\langle O\rangle = \frac{1}{2} \sum_{a,b}^{r} \langle O \cdot g_{a,b} \rangle - r \sum_{a=1}^{r} \langle O \cdot g_{a,r+1} \rangle + \frac{r(r+1)}{2} \langle O \cdot g_{r+1,r+2} \rangle - \frac{r}{2} \langle O \cdot g_{r+1,r+1} \rangle, \tag{13}$$

where we used the operator d_{τ} defined as

$$d_{\tau} = \frac{1}{\beta(1+t)\sqrt{\alpha}} \frac{d}{ds},\tag{14}$$

l=1 i

in order to simplify calculations and presentation.

Definition 8 The centered and rescaled overlap fluctuations $\xi_{l,m}$ and $\rho_{l,m}$ are introduced as

$$\xi_{l,m} = \sqrt{N} [q_{l,m} - \delta_{l,m} W - (1 - \delta_{l,m})q] \rho_{l,m} = \sqrt{P} [p_{l,m} - \delta_{l,m} X - (1 - \delta_{l,m})p].$$

Posing $Y(s) = \langle \xi_{12} \rho_{12} \rangle_s + \frac{\langle \xi_{12}^2 \rangle_0}{\langle \rho_{12}^2 \rangle_0} \langle \rho_{12} \rangle_s$ we get $d_\tau Y = Y^2$, whose solution evaluated at $\tau = \beta (1+t) \sqrt{\alpha} s$, s = 1 is given by

$$Y(s=1) = \frac{W}{1-\beta(1+t)W(1+\sqrt{\alpha})},$$

$$W\Delta^{2} = 1 - \frac{t\Delta}{\beta(1+t)},$$

$$\Delta = 1 + \frac{\alpha t}{1-\beta(1+t)W}.$$
(15)

Since we are interested in obtaining the critical temperature for ergodicity breaking, where fluctuations (in this case Y) grow arbitrarily large we check where the denominator at the r.h.s. of eq. (15) vanishes.

Theorem 2 The ergodic region of the model defined by the cost function (1) is delimited by the following critical surface in the (α, β, t) space of the tunable parameters

$$\beta_c = \frac{1}{1+t} \Big[\frac{\Delta^2}{1+\sqrt{\alpha}} + t\Delta \Big],$$
$$\Delta = 1 + \sqrt{\alpha} (1+\sqrt{\alpha})t.$$

Remark 4 At t = 0, where the model reduces to Hopfield's scenario, the critical surface correctly collapses over the Amit-Gutfreund-Sompolinsky critical line $\beta_c = (1 + \sqrt{\alpha})^{-1}$, but -in the deep sleep phase (i.e. in the large t limit)-the ergodic region collapses to the axis T = 0: the ergodic region -during the sleep state- invades the spin-glass region, the latter being de facto suppressed.

Dreaming Boltzmann machine

Boltzmann Machines with "large" hidden-layer size

$$\mathscr{H}_{N,K}^{(RR-HN)}(\boldsymbol{\sigma};\boldsymbol{\xi},t) = -\frac{1}{2N} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{K,K} \xi_i^{\mu} \,\xi_j^{\nu} \left(\frac{1+t}{\boldsymbol{I}_K + t\boldsymbol{C}}\right)_{\mu,\nu} \sigma_i \,\sigma_j$$

HN with Reinforcement&Removal

$$\mathscr{H}_{N,K}^{(RR-BM)}(\boldsymbol{\sigma}, \boldsymbol{z}; \boldsymbol{w}, t) = -\frac{1}{\sqrt{N}} \sum_{i,\mu} w_i^{\mu} \sigma_i z_{\mu} + \frac{t}{2(1+t)} \sum_{\mu,\nu} C_{\mu,\nu} z_{\mu} z_{\nu}$$

BM with Reinforcement&Removal



$$C_{\mu\nu} := \frac{1}{N} \sum_{i=1}^{N} \xi_i^{\mu} \xi_j^{\nu}$$

$$z_{\mu} \sim \mathcal{N}(0, [\beta/(1+t)]^{-1})$$

EA, F. Alemanno, A. Barra, A. Fachechi, Dreaming Neural Networks for learning (2021)

Boltzmann Machines with "large" hidden-layer size

$$\mathscr{H}_{N,K}^{(RR-HN)}(\boldsymbol{\sigma};\boldsymbol{\xi},t) = -\frac{1}{2N} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{K,K} \xi_i^{\mu} \xi_j^{\nu} \left(\frac{1+t}{\boldsymbol{I}_K + t\boldsymbol{C}}\right)_{\mu,\nu} \sigma_i \sigma_j$$

HN with Reinforcement&Removal

$$\mathscr{H}_{N,K}^{(RR-BM)}(\boldsymbol{\sigma}, \boldsymbol{z}; \boldsymbol{w}, t) = -\frac{1}{\sqrt{N}} \sum_{i,\mu} w_i^{\mu} \sigma_i z_{\mu} + \frac{f(t)}{2[1+f(t)]} \sum_{\mu,\nu} P_{\mu,\nu} z_{\mu} z_{\nu}$$

BM with Reinforcement&Removal

P empirical Pearson correlation matrix

$$P_{\mu\nu} = \frac{\frac{1}{N} \sum_{i=1}^{N} (W_{i\mu} - \bar{W}_{i\mu})(W_{i\nu} - \bar{W}_{i\nu})}{\sqrt{(\frac{1}{N} \sum_{i=1}^{N} W_{i\mu}^2 - \bar{W}_{i\mu}^2)(\frac{1}{N} \sum_{i=1}^{N} W_{i\nu}^2 - \bar{W}_{i\nu}^2)}}$$

f(t) s.t. f(0) = 0 and $[I_K + f(t)P]$ is pos. def. (e.g., f(t) = |t|)

 $\boldsymbol{z} \sim \mathcal{N}(0, \boldsymbol{I}_{\boldsymbol{K}}(\beta/[1+f(t)])^{-1})$



Boltzmann Machines with "large" hidden-layer size

$$\mathscr{H}_{N,K}^{(RR-HN)}(\boldsymbol{\sigma};\boldsymbol{\xi},t) = -\frac{1}{2N} \sum_{i,j=1}^{N,N} \sum_{\mu,\nu=1}^{K,K} \xi_i^{\mu} \xi_j^{\nu} \left(\frac{1+t}{\boldsymbol{I}_K + t\boldsymbol{C}}\right)_{\mu,\nu} \sigma_i \sigma_j$$

HN with Reinforcement&Removal

$$\mathscr{H}_{N,K}^{(RR-BM)}(\boldsymbol{\sigma}, \boldsymbol{z}; \boldsymbol{w}, t) = -\frac{1}{\sqrt{N}} \sum_{i,\mu} w_i^{\mu} \sigma_i z_{\mu} + \frac{f(t)}{2[1+f(t)]} \sum_{\mu,\nu} P_{\mu,\nu} z_{\mu} z_{\nu}$$

BM with Reinforcement&Removal

P empirical Pearson correlation matrix

$$P_{\mu\nu} = \frac{\frac{1}{N} \sum_{i=1}^{N} (W_{i\mu} - \bar{W}_{i\mu})(W_{i\nu} - \bar{W}_{i\nu})}{\sqrt{(\frac{1}{N} \sum_{i=1}^{N} W_{i\mu}^2 - \bar{W}_{i\mu}^2)(\frac{1}{N} \sum_{i=1}^{N} W_{i\nu}^2 - \bar{W}_{i\nu}^2)}}$$

$$f(t)$$
 s.t. $f(0) = 0$ and $[I_K + f(t)P]$ is pos. def. (e.g., $f(t) = |t|)$

$$\boldsymbol{z} \sim \mathcal{N}(0, \boldsymbol{I}_{K}(\beta/[1+f(t)])^{-1})$$

By gradient descent over KL divergence



Sub-critical regime (N = 100, K = 5, T = 0.2)



Super-critical regime (N = 100, K = 20, T = 0.2)





N= 100, 200, 400, CD-1, ϵ = 0.0005 *T*, p= 0.05 example noise, M= 600 examples per pattern (M_t = 500 for training, M_v = 100 for validation), M_b = 200 minibatch size, τ = 3000 epochs, network parameter initial condition $w_{i,\mu} \sim \mathcal{N}(0,1)$, supervised.



N= 100, 200, 400, CD-1, ϵ = 0.0005 T, p= 0.05 examples noise, M= 600 examples per pattern (M_t = 500 for training, M_v = 100 for validation), M_b = 200 minibatch size, τ = 3000 epochs, network parameter initial condition $w_{i,\mu} \sim \mathcal{N}(0,1)$.



The graphical representation of the RBM weights is surrounded by a dark halo, whose origin lies in the fact that none of the examples exhibits activated pixels in that area. This means that these sites would contribute positively to the Mattis magnetizations, thus increasing the signal for all the hidden units. This does not happen for the DBM where entries corresponding to the external part of patterns are randomly oriented; hence, their contribution to the signal is lower w.r.t. the RBM case.

Results for the MNIST dataset with supervised training and random initialization of the weights. (a) the pseudo-(log)likelihood, averaged over five different realizations of (b) training procedure, (c) dreaming time f(t) as a function of the learning time τ , and (d) the final accuracy as a function of the temperature, computed after 100 updates of the whole network. (e) and (f) First four digits, respectively, for the RBM and the DBM. The learning parameters are $\epsilon = 0.005T$, the size of the minibatch is fixed to 200 examples for epoch.



Learning results for the MNIST dataset with unsupervised training and random initialization of the weights. The plots show the results for the pseudo-(log)likelihood (a) and (b) averaged over five different realizations of the training procedure, and the dreaming time as a function of the learning time (c) and (d) for each of the five different learning sessions. The initial condition is $W_{i\mu} \sim \mathcal{N}(0,0.1)$. The learning strength is $\epsilon = 0.005T$, the temperature is T = 1, the size of the minibatch is fixed to 200 examples, and the number of epochs is 2000.

t as a free parameter: it smoothly interpolates between the RBM and the DBM models; thus, the minima landscape of the objective function will be modified in a continuous way



