# Exploring the Role of Liquid States in Neural Networks: From Feedforward Networks to Attractor Models

Riccardo Zecchina



Finanziato dall'Unione europea **NextGenerationEU** 







Ministero dell'Università e della Ricerca





and negative margin classifiers

and continuous weights

neural networks



- 1. Accessible liquid states in non-convex feedforward models
- 2. Geometry of perfect learning: differences between discrete

3. Chaotic and stable liquid attractors in recurrent asymmetric

# Basic problems

- Non-convexity and first-order algorithms

Overparametrization & generalization (overfitting under control)

Internal representations if feedforward and attractor NN



Let's consider feedforward architecture first

Training set:

 $\hat{y}^{\mu} \leftrightarrow y^{\mu}$ 

 $\{(x^{\mu}, y^{\mu})\}_{\mu=1,...,M}$ 

 $\Delta^{\mu} \stackrel{\circ}{=} \left( W^{K} \cdot \sigma_{K}(x^{\mu}) \right) y^{\mu}$ 

pre-activations at layer K



- Energy = "0-1 loss": number of errors on the training set (not differentiable)

 $\mathcal{L}_{NE} = \sum \left[ (1 - \delta(\hat{y}^{\mu}, y^{\mu})) \right]$  $\mu$ 

### - Surrogate differentiable losses

Cross-entropy: "loglikelihood", softmax

 $\mathcal{L}_{CE} = -\sum \left( \hat{\Delta}_{y^{\mu}}^{\mu} - \log \sum \exp \gamma \hat{\Delta}_{k}^{\mu} \right)$  $\mu$ k $\Delta^{\mu} \stackrel{\circ}{=} (W^K \cdot \sigma_K(x^{\mu})) y^{\mu}$ 

Mean Square Error, ...

The simplest non convex models: binary perceptron and negative stability spherical perceptrons

**Random Training set:**  $\{(\bar{x}^{\mu}, y^{\mu})\}$   $\mu =$ 

Control parameter:

 $\alpha = \frac{\# \text{ patterns}}{\# \text{ weigts}}$ 

**Learning:** find  $W_i$  such that



1,..., 
$$P = \alpha N$$
  $x_i^{\mu} = \pm 1$  with  $p = \frac{1}{2}$ ,  $i = 1, .$   
 $y^{\mu} = \pm 1$  with  $p = \frac{1}{2}$ 

$$\operatorname{Sign}(\sum_{i=1}^{N} W_{i} x_{i}^{\mu}) = \operatorname{Sign}(W \cdot x^{\mu}) \quad \forall \mu$$

$$= \operatorname{Sign}(\Delta^{\mu}) \text{ or } \Delta^{\mu} > 0$$

$$\mathbf{A}^{\mu} \equiv y^{\mu} \sum_{i=1}^{N} W_{i} x_{i}^{\mu} = y^{\mu} W \cdot x^{\mu} \qquad \text{(stability, margin)}$$



the site of the store of the set i.i.d. random variables, and 2) the generalization (or teacher-student) scenario, in which they are provided.byd. Binary perception: another  $i_{p}$  berceptron with sympler tic weights  $W^{\mathscr{T}}$ . In the classification scenario, the typiinput.

Simplified algorithms.—Only a handful of algorithms: are currently believed to be able to solve the classification large N in a sub-exponential running time: they are all.—

Propagation using a solution from SBPI, wi classificational or estoraged times from eq. (3); (dotted black) upper bound ( random Variables), and and and and and the randomderestimitenthe spectra finneshielthe connected clusters; the BP curve is lower th cal problem is solvable with probability 1 in the lime of teacherse in large N up to  $\alpha = 0.833$  [5], after which the probability we distance, While in the three tit is fixed the With (negative) margin: With (negative) margin: of finding a solution drops to zero.  $\alpha_c$  is called the car is between a typical solution and one found with the probability of finding a solution drops to zero.  $\alpha_c$  is called the car is problem. pacity; we also use this term for the maximum value of problem is solve by with probable  $i \in \mathbb{R}$  $\alpha$  for which a solution can be found by a specific algoes a specific algoes and the specific algoes and the problem of the specific algoes are the specific algoes and the problem of the specific algoes are the specific algoes and the problem of the specific algoes are the specific algoes and the problem of the specific algoes are the specific algoes and the problem of the specific algoes are the specific algoes and the problem of the problem of the specific algoes are the problem of the problem of the specific algoes are the problem of th ble: the teacher itself [2, 6]. And additional quantity of or which a solution can be found interest in this scenario is the generalization error rate  $p_e = \frac{1}{\pi} \arccos\left(\frac{1}{N}W \cdot W^{\mathscr{T}}\right)$ , which is the probability that update schemes), making them appealing  $\tau(W,\xi^*) = \tau(W^{\mathscr{T}},\xi^*)$  when  $\xi^*$  is a previously unseen ltipheplementations, up be out  $\varphi f_{\mathscr{T}}$  achieving  $\xi$ is a first citaler transition and algorith Max-Sum scheme, can be shown to have s theistest of the Australia and the set of the set problem and achieve a non-zero capacity in the limit erestmen which see enanies is early with input. A qualitatively similar scenario ho

![](_page_6_Figure_3.jpeg)

### MLP 3072-1000-1 Cifar10 (bird, deer)

![](_page_7_Figure_1.jpeg)

### Start with the Binary perceptron

- $\forall \alpha > 0$  typical solutions are isolated (Huang, Kabashima (2014));

![](_page_8_Figure_4.jpeg)

• The space of solution splits into separated states of vanishing entropy (Krauth, Mézard (1989));

• Rigorous: Abbe, Li, Sly (2021), Perkins, Xu (2021), Nakajima, Sun (2022), Aubin, Perkins, Zdeborova (2019)

 $\alpha$ 

If this would be the case, locally stable algorithms would fail and learning would be hard: The so called *Overlap Gap Property* would hold

### f У in 🖂 🦲 PERSPECTIVE 🛛 🛇 The overlap gap property: A topological barrier to optimizing over random structures

David Gamarnik 问 🏼 Authors Info & Affiliations

Edited by David L. Donoho, Stanford University, Stanford, CA, and approved August 19, 2021 (received for review May 5, 2021)

**October 1, 2021** 118 (41) e2108492118 <u>https://doi.org/10.1073/pnas.2108492118</u>

If this would be the case, locally stable algorithms would fail and learning would be hard: The so called *Overlap Gap Property* would hold

### f 🎔 in 🖂 🦲 PERSPECTIVE 🛛 🛇 The overlap gap property: A topological barrier to optimizing over random structures

David Gamarnik 问 🏼 Authors Info & Affiliations

Edited by David L. Donoho, Stanford University, Stanford, CA, and approved August 19, 2021 (received for review May 5, 2021)

**October 1, 2021** 118 (41) e2108492118 <u>https://doi.org/10.1073/pnas.2108492118</u>

### This contradicts empirical evidence!

### If this would be the case, locally stable algorithms would fail and learning would be hard: The so called *Overlap Gap Property* would hold

### У in 🖂 🦲 PERSPECTIVE | The overlap gap property: A topological barrier to optimizing over random structures

![](_page_11_Figure_3.jpeg)

![](_page_11_Picture_6.jpeg)

![](_page_12_Figure_0.jpeg)

- 1) Regions observed with a diameter d
- 2) Internal entropy:  $S_I(d, y) = -\langle \mathcal{E}(\tilde{W}) \rangle_{\xi, \tilde{W}} =$
- 3) External entropy (number of regions of diameter d with a given internal entropy):

![](_page_13_Figure_3.jpeg)

$$\frac{1}{N} \langle \log \mathcal{N}(\tilde{W}, d) \rangle_{\xi, \tilde{W}} \quad S_I(d, y) = \partial_y(y \mathcal{F}(d, y))$$

$$\mathcal{S}_E(d, y) = -y[\mathcal{F}(d, y) + \mathcal{S}_I(d)]$$

![](_page_13_Picture_7.jpeg)

(perception) maps vectors qi, ir supurs The vector qits binary outputs as  $\tau(W, \xi) = \text{sign}(W, \xi)$ , where  $W \in [1, 1]^N$  is the vector of synaptic weights. Given  $W \in [1, 1]^N$ 0.005 input patterns  $\xi^{\mu}$  with  $\chi espending N deside de <math>\omega t p u t s$  $\tau(W,\xi^*) = \tau(W^{\mathscr{T}},\xi^*)$  which is solutions upsed by probably the physical set of input. Simplified algorithms of the teacher it selfpic: are currently believed to be able to solve the classification; f(x) = 0.69; Yet another algorithm based on a subject of the solution of the solu

problem and achieve a interest apacithis scenario eresthing this scenario discrimination of the sign of the sign of the station of the second state of the second stat

![](_page_14_Figure_2.jpeg)

### distance from

![](_page_15_Figure_1.jpeg)

Geometrical phase discontinuous transition:  $\alpha_U \simeq 0.75$ 

C. Baldassi, A. Ingrosso, C. Lucibello, L. Saglietti, and R. Zecchina Phys. Rev. Lett. 115, 128101 (2015)

Check the existence of subdominant dense regions of solutions in binary networks

![](_page_15_Figure_5.jpeg)

![](_page_15_Figure_6.jpeg)

![](_page_15_Picture_7.jpeg)

![](_page_15_Picture_8.jpeg)

![](_page_15_Picture_9.jpeg)

![](_page_16_Figure_1.jpeg)

### Generalization

### Student

### Algorithmic follow up: real replicas

Local free entropy: 
$$\phi(W, \gamma, \beta) = \log \sum_{W'} e^{-\beta \mathcal{L}_{NE}(W') - \frac{\gamma}{2}d(W, W')}$$
  $(\mathcal{L}_{NE} \equiv \mathcal{L}_{0-1})$ 

Large-deviation partition function: Z(y,

Assume y integer, and transform the partition function

$$Z(y,\gamma,\beta',\beta) = \sum_{W,\{W_a\}} e^{-\beta' \mathcal{L}_{NE}(W) - \beta \sum_{a=1}^{y} \mathcal{L}_{NE}(W_a) - \frac{\gamma}{2} \sum_{a=1}^{y} d(W,W_a)}$$
center y (real) replicas

CEIIIEI

$$\gamma, \beta', \beta) = \sum_{W} e^{-\beta' \mathcal{L}_{NE}(W) + y \phi(W, \gamma, \beta)}$$

interaction

![](_page_17_Figure_9.jpeg)

### Algorithmic follow up: real replicas

Local free entropy: 
$$\phi(W, \gamma, \beta) = \log \sum_{W'} e^{-\beta \mathcal{L}_{NE}(W') - \frac{\gamma}{2}d(W, W')}$$
  $(\mathcal{L}_{NE} \equiv \mathcal{L}_{0-1})$ 

Large-deviation partition function: Z(y,

Assume y integer, and transform the partition function

$$Z(y,\gamma,\beta',\beta) = \sum_{W,\{W_a\}} e^{-\beta'\mathcal{L}_N}$$

center

$$\gamma, \beta', \beta) = \sum_{W} e^{-\beta' \mathcal{L}_{NE}(W) + y \phi(W, \gamma, \beta)}$$

interaction  $V_E(W) - \beta \sum_{a=1}^y \mathcal{L}_{NE}(W_a) - \frac{\gamma}{2} \sum_{a=1}^y d(W, W_a)$ y (real) replicas

![](_page_18_Figure_9.jpeg)

![](_page_18_Picture_10.jpeg)

- Local Entropy driven Simulated Annealing
- Replicated Message-Passing (aka "focusing Belief Propagation")
- Replicated Stochastic Gradient Descent (SGD)
- Entropy-SGD: Langevin dynamics to estimate local entropy
- Replicated Greedy Algorithms
- Sharpness Aware Minimization
- Stochastic weights + gradient on the probabilities
- Quantum Annealing delocalization mechanism for finding NN ground states
- . . .

Local entropy algorithms

### Local Entropy driven Simulated Annealing

**Objective Function:** (minimize the "energy")

![](_page_20_Picture_2.jpeg)

### 1. SA moves

2. BP method to estimate the local entropy, or use a replicated model

![](_page_20_Figure_5.jpeg)

We performed extensive simulations and studied the scaling properties of EdMC in trast to simulated annealing. Figure 2 is a log-log plot of the number of iterations  $n_{E^{\pm}}$ reach a solution obtained for increasing N at  $\alpha = 0.3$ . A least squares fit $(n_{E=0} \propto N^{2.84})$ firms the evident power law behaviour. Note that even with an extremely low cooling ra

![](_page_20_Figure_7.jpeg)

### Connection with Monasson's 1-rsb formalism (1995)

- 1) We are interested in the out-of-equilibrium regime m > 1
- 2)  $q_1 = 1 2 D$  is the diameter of the hypersphere. We do not need to maximise over  $q_1$
- 3) If we find a value of  $q_1$  that maximises the free entropy then we have an out-of-equilibrium state (m > 1) as in Monasson

$$Z_{1RSB}(\beta', y, D) = \sum_{\{W^a\}} e^{-\beta' \sum_{a=1}^{y} E(W^a)} \prod_{a>b} \delta\left(d\left(W^a, W^b\right) - ND\right)$$

 $Z_{1RSB}\left(\beta', y\right) = \max_{D} Z_{1RSB}\left(\beta', y, D\right)$ 

R. Monasson, Physical review letters, 75(15):2847, 1995. C. Baldassi, F. Pittorino. R. Zecchina, PNAS 117, 161-179 (2020)

![](_page_21_Picture_8.jpeg)

- How are these dense regions composed? • Why they generalise well?
- Do they contain high margin solutions?

![](_page_22_Figure_3.jpeg)

$$\kappa) = \prod_{\mu=1}^{P} \Theta \left( y^{\mu} \sigma_{\text{out}}^{\mu} - \kappa \right)$$

 $E = \sum \Theta \left( -y^{\mu} \sigma_{\text{out}}^{\mu} + \kappa \right)$ 

### Typical high margin solutions are less but tend to be much closer to each other

when  $\kappa = \kappa_{max}$ 

distance  $d=(1-q_1)/2$ between typical solutions for a given margin

![](_page_23_Figure_3.jpeg)

The lines change from solid to dashed when the entropy of solutions becomes negative, i.e.

![](_page_23_Picture_5.jpeg)

# A wide flat minima arises by the coalescence of (atypical) high margin minima

![](_page_24_Picture_1.jpeg)

### Unveiling the structure of wide flat minima in neural networks

Carlo Baldassi,<sup>1</sup> Clarissa Lauditi,<sup>2</sup> Enrico M. Malatesta,<sup>1</sup> Gabriele Perugini,<sup>1</sup> and Riccardo Zecchina<sup>1</sup> <sup>1</sup>Artificial Intelligence Lab, Bocconi University, 20136 Milano, Italy <sup>2</sup>Department of Applied Science and Technology, Politecnico di Torino, 10129 Torino, Italy arXiv:2107.01163

![](_page_24_Picture_4.jpeg)

Binary perceptron: efficient algorithms can find solutions in a rare well-connected cluster

arXiv:2111.03084 Emmanuel Abbe \* Allan Sly<sup>‡</sup> Shuangping Li<sup>†</sup>

![](_page_24_Picture_8.jpeg)

# A wide flat minima arises by the coalescence of (atypical) high margin minima

### liquid regions

![](_page_25_Picture_2.jpeg)

### Unveiling the structure of wide flat minima in neural networks

Carlo Baldassi,<sup>1</sup> Clarissa Lauditi,<sup>2</sup> Enrico M. Malatesta,<sup>1</sup> Gabriele Perugini,<sup>1</sup> and Riccardo Zecchina<sup>1</sup> <sup>1</sup>Artificial Intelligence Lab, Bocconi University, 20136 Milano, Italy <sup>2</sup>Department of Applied Science and Technology, Politecnico di Torino, 10129 Torino, Italy arXiv:2107.01163

![](_page_25_Picture_5.jpeg)

Binary perceptron: efficient algorithms can find solutions in a rare well-connected cluster

arXiv:2111.03084 Emmanuel Abbe \* Shuangping Li<sup>†</sup> Allan Sly<sup>‡</sup>

![](_page_25_Picture_9.jpeg)

## Overparametrization & non convexity

### The non-convex random features model

![](_page_26_Figure_2.jpeg)

![](_page_26_Figure_4.jpeg)

- Gaussian Equivalence Theorem: in the large N, D, P limit (with  $\alpha = P/N$ ,  $\alpha_T = P/D$  fixed)

$$\tilde{x}_i^{\mu} = \sigma \left( \frac{1}{\sqrt{D}} \sum_{k=1}^D F_{ki} x_k^{\mu} \right) \to \mu_0 + \frac{\mu_1}{\sqrt{D}} \sum_{k=1}^D F_{ki} x_k^{\mu} + \mu_\star z_i^{\mu}$$

![](_page_27_Figure_2.jpeg)

- Statistical Mechanics analysis similar to the underparametrized case: geometrical structure of minimizers, large deviation study of dense regions.

,  $z_i^\mu \sim \mathcal{N}(0,1)$ Montanari, Mei (2019); Goldt, Mézard, Krzakala, Zdeborová (2020)

(Baldassi, Malatesta, Lauditi, Pacelli, Perugini, RZ, PRE 2021)

![](_page_27_Figure_8.jpeg)

![](_page_27_Picture_9.jpeg)

### Interpolation Threshold $\alpha_c$ vs Algorithmic Threshold $\alpha_{LE}$ in **discrete** NN

![](_page_28_Figure_1.jpeg)

### Continuous networks have differences

# Let's consider the simplest non convex continuous network: the negative margin spherical perceptron.

 $W_i \in \mathbb{R}$  &

 $1 = \operatorname{Sign}(\Delta t)$ 

![](_page_29_Figure_4.jpeg)

$$\& \sum_{i=1}^{N} W_i^2 = N$$

$$(\mu - k)$$
 or  $\Delta^{\mu} > k$ 

![](_page_29_Figure_7.jpeg)

### T=0 phase diagram

![](_page_30_Figure_1.jpeg)

Simplified model for Jamming: Franz, Parisi, Urbani, Zamponi, ..., and many others

![](_page_30_Figure_3.jpeg)

### LE transition: Appearance of max entropy minima

![](_page_31_Figure_1.jpeg)

•Not directly related to computational hardness as the model is fRSB.

•Still related to geometry and generalization.

![](_page_31_Figure_5.jpeg)

![](_page_31_Figure_6.jpeg)

### Expected train error of the linear interpolation of y replicas

$$w^{\gamma} \equiv rac{\sum_{a=1}^{y} \gamma_a w_a}{\|\sum_{a=1}^{y} \gamma_a w_a\|}$$
 with

![](_page_32_Figure_2.jpeg)

$$\sum_{a} \gamma_{a} = \mathbf{1}$$

![](_page_33_Figure_0.jpeg)

work by B. Annesi, C. Lauditi, C. Lucibello, E. M. Malatesta, G. Perugini, F. Pittorino, L. Saglietti, 2023

![](_page_33_Picture_3.jpeg)

### Atypical solutions are surrounded by an exponentially higher number of solutions wrt typical.

![](_page_34_Figure_1.jpeg)

 $\alpha = 2, k=-0.5$ 

Use BP to compute entropy around solutions

In neural network we expect a generalisation cross over at  $\alpha_{LE}$ 

**Correct regime for** overparametrized NN models

(detailed balance algorithm)

![](_page_34_Picture_7.jpeg)

### Comparison: LE transition in the binary negative k perceptron

capacity

![](_page_35_Figure_2.jpeg)

generalization : learning with k<0, test with k=0

![](_page_35_Figure_4.jpeg)

- In continuous models, overparameterization and negative k both increase (exponentially) the volume of minimisers (aka solutions)
- We need to find the minimisers which have good generalisations
- Barycentre of liquid regions (locally Bayesian)
- Specific large margin solutions are hard to find at zero training error

Large networks generalise optimally at negative stabilities in liquid flat regions

### Large networks generalise optimally at negative stabilities if in high entropy solutions

![](_page_37_Figure_1.jpeg)

$$\alpha = 5.0, \alpha_T = 50$$

### Overparametrized continuous models:

### Teacher: input-noisy perceptron

Baldassi, Lauditi, Malatesta, Perugini, Saglietti, Zecchina, in preparation (2023)

![](_page_37_Picture_7.jpeg)

![](_page_38_Figure_0.jpeg)

![](_page_38_Picture_1.jpeg)

### MLP 3072-1000-1 Cifar10 (bird, deer)

![](_page_39_Figure_1.jpeg)

### On going wok:

- SGD with negative k in the hidden layers
- Smaller architectures with good generalization
- Continual & few-shot learning

Baldassi, Mezard, Zecchina, in preparation (2023) Baldassi, Lauditi, Malatesta, Mezard, Perugini, Saglietti, Zecchina, in preparation (2023)

![](_page_41_Picture_1.jpeg)

# Liquid fixed points in recurrent asymmetric neural network

Baldassi, Mezard, Zecchina, in preparation (2023) Baldassi, Lauditi, Malatesta, Mezard, Perugini, Saglietti, Zecchina, in preparation (2023)

![](_page_42_Picture_2.jpeg)

### Some basic results in asymmetric NN

### **Chaos in Random Neural Networks**

H. Sompolinsky, A. Crisanti & H. J. Sommers, PRL 1988

$$\dot{h}_i = -h_i + \sum_{j=1}^N J_{ij} S_j = -h_i + \sum_{j=1}^N J_{ij} \phi(h_j)$$

**Dynamics of Random Neural Networks with Bistable Units** M. Stern, H. Sompolinsky and L. Abbott, PRE 2014

$$\frac{dx_i}{dt} = -x_i + s \tanh(x_i) + g \sum_{j(\neq i)} J_{ij} \tanh(x_j)$$
$$m_i = \tanh\left[sm_i + g \sum_{j(\neq i)} J_{ij}m_j\right]$$

 $m_i = \tanh(x_i)$ 

![](_page_43_Figure_7.jpeg)

Interactions between the neurons within a cluster are represented in these models by self-coupling s

In the limit  $g \rightarrow \infty$ ,  $s \rightarrow \infty$  with fixed  $J_D = s/g$  we get back the binary model.

# Can an asymmetric random attractor neural network display exponentially many stable attractors (internal representations)?

# We expect liquid fixed points to exist

# Perceptron storing $\alpha N$ patterns $\sigma^{\mu} = \operatorname{Sign}(\sum_{i=1}^{N} \xi_{i}^{\mu} W_{i}) \quad \forall \mu \quad (\mu = 1, ..., \alpha N)$ $W_i = \pm 1$

Asymmetric recurrent network of neurons

$$S_i = \operatorname{Sign}(\sum_{j=1}^N J_{ij}S_j) \quad \forall i \quad (i = 1, ..., N)$$
$$S_i = \pm 1$$

In the recurrent NN we have more constraints. To find similar phase space we need to relax them.

 $J_{ii} = J_D > 0$ One possibility (among others): add a diagonal feedback term

$$P(\xi_i^{\mu}, \xi_j^{\nu}) = P(\xi_i^{\mu}) P(\xi_j^{\nu})$$
$$P(\xi_i^{\mu} = \pm 1) = \frac{1}{2} \quad P(\sigma^{\mu} = \pm 1) = \frac{1}{2}$$

$$P(J_{ij}, J_{i'j'}) = P(J_{ij}) P(J_{i'j'})$$
  
 $P(J_{ij} = \pm 1) = \frac{1}{2}$ 

Random binary couplings

$$P(J_{ij}, J_{ji}) = P(J_{ij})P(J_{ji}) \qquad P(J_{ij} = \pm 1) = \frac{1}{2}$$

Model: fixed point conditions of the update dynamics

$$s_i^{t+1} = \text{Sign}\left(J_D s_i^t + \sum_{i \neq j} J_{ij} s_j^t\right)$$

Factors: count violated fixed point constraints

$$E_J(s) = \sum_{i=1}^N \left[ 1 - s_i \text{Sign} \left( J_D s_i + \sum_{j(\neq i)} J_{ij} s_j \right) \right]$$

### $J_D \ge 0$ local feedback

![](_page_46_Figure_7.jpeg)

![](_page_46_Picture_10.jpeg)

# Analogies with error correcting codes

![](_page_47_Figure_1.jpeg)

Exponential number of correlated attractive configurations (w.r.t. the decoding dynamics)

They belong to some subspace, e.g. linear subspace in the case of linear codes (LDPC)

related work by A. Karbasi, A. H. Salavati, A. Shokrollahi, and L. R. Varshney, ...

![](_page_47_Figure_5.jpeg)

![](_page_48_Figure_1.jpeg)

![](_page_48_Picture_2.jpeg)

![](_page_48_Picture_3.jpeg)

![](_page_48_Picture_4.jpeg)

# In NN noise will also appear in the input before the encoding

![](_page_49_Figure_1.jpeg)

![](_page_49_Picture_2.jpeg)

![](_page_49_Picture_3.jpeg)

![](_page_49_Picture_4.jpeg)

Can asymmetric NN perform like sub-optimal ECC?

Can any input mapped to an internal representation (codeword) ?

Can they correct an extensive number of errors?

Can asymmetric NN perform like sub-optimal ECC?

Can any input mapped to an internal representation (codeword) ?

Can they correct an extensive number of errors?

![](_page_51_Figure_3.jpeg)

![](_page_51_Figure_4.jpeg)

Yes

# Non-linear liquid code

![](_page_52_Picture_1.jpeg)

Coding rate:

$$\mathbf{s}_A$$

![](_page_52_Picture_6.jpeg)

 $\mathbf{S}_B$ 

![](_page_52_Picture_8.jpeg)

# Shannon bound and positioning of simple random ANNs

![](_page_53_Figure_1.jpeg)

numerical result (could be done analytically)

![](_page_54_Figure_0.jpeg)

$$\rightarrow \quad \eta_i = \operatorname{Sign}(J_D \eta_i^t + \sum_{j \neq i} J_{ij} \eta_j) \quad , \quad \eta = \eta_i$$

 $\{\xi^{\mu}\} \rightarrow \{\eta^{\mu}\}$ 

![](_page_54_Picture_6.jpeg)

### Number of configurations which satisfy the fixed point condition, $\beta \rightarrow \infty$

$$Z_J = \sum_{s} e^{-\beta E_J(s)} = e^{-\beta N} \sum_{s} e^{\beta \sum_{i} \operatorname{Sign}(J_D + s_i \sum_{j(\neq i)} J_{ij} s_j)}$$

### First moment gives the **exact** typical number of fixed points:

$$\Phi_1 = -\frac{1}{\beta N} \log \mathbb{E}_J Z_J = 1 - \frac{1}{\beta} \left[ \log 2 + \log \left( e^\beta H(-J_D) + e^{-\beta} H(J_D) \right) \right]$$

where H(x) =

$$\int_{x}^{\infty} \frac{dz}{\sqrt{2\pi}} e^{-z^2/2} = \frac{1}{2} \operatorname{Erfc}(\frac{x}{\sqrt{2}})$$

$$S_1 = \frac{1}{N} \log \bar{Z}_J|_{\beta \to \infty} = \log 2 + \log |$$

![](_page_56_Figure_1.jpeg)

# $[1 - H(J_D)]$

 ${f s}({f t})\cdot{f s}'({f t})$  $\mathbf{s}(\mathbf{0}) = \mathbf{s}'(\mathbf{0}) + 1$  spin perturbation

![](_page_56_Figure_4.jpeg)

4

![](_page_56_Figure_6.jpeg)

![](_page_56_Figure_7.jpeg)

![](_page_56_Figure_8.jpeg)

# RS, 1-RSB and Local Entropy results

- The RS and 1-RSB solution give back the Annealed (first moment results) for the dominating fixed points ( $q_0 = 0$ , things are relatively simple)
- We can use the 1-RSB formalism to analyse existence of liquid fixed points

$$\Phi = -\beta - \frac{1}{2} \left( r_1 \left( 1 - q_1 \right) \right) - \frac{m}{2} q_1 r_1 + \frac{1}{m} \log \left\{ \int Dz \, Dt \left[ e^{z \sqrt{r_1}} \left( e^{\beta} H'_{--} + e^{-\beta} H'_{++} \right) + e^{-z \sqrt{r_1}} \left( e^{\beta} H'_{-+} + e^{-\beta} H'_{+-} \right) \right]^m \right\}$$

$$Q^{ab} = \frac{1}{N} \sum_{k} s^{a}_{k} s^{b}_{k} \text{ for } a < b$$
$$R^{ab} = \frac{1}{N} \sum_{k} \hat{h}^{a}_{k} \hat{h}^{b}_{k} s^{a}_{k} s^{b}_{k} \text{ for } a \leq b$$

$$H'_{ss'} \equiv H\left(\frac{sJ + s't\sqrt{q_1}}{\sqrt{1-q_1}}\right) \qquad (J_D \equiv J)$$

![](_page_58_Figure_1.jpeg)

entropy

distance

![](_page_58_Figure_4.jpeg)

change of concavity:

geometrically impossible, the entropy has to increase because of the interchangeability of the centroid and solutions at distance d

![](_page_58_Figure_7.jpeg)

$0 \le J \le J_c \simeq 0.1$	isolated solutions co-exists $d = \frac{1-2}{2}$
$J > J_c$	liquid fixed points: isola number of extended de
$J > J_a(N)$	more fixed points bec
$J = \infty$	all configurations are

### xist with an exponential number of small cluster

$$\frac{-q}{2}$$
  $q \in [0.998, 1]$ 

ated solutions co-exist with an exponential ense regions which connect the small clusters

### come accessible (strong finite size effects)

trivially fixed points

![](_page_60_Figure_0.jpeg)

Isolated and well separated small clusters, Overlap-Gap-Property (Gamarnik et al.)

# Non chaotic local dynamics: fBP dynamics

factor graph:

![](_page_61_Figure_2.jpeg)

commutative and associative  $\rightarrow$ 

inverse operation  $\rightarrow$ 

Nor

![](_page_61_Picture_6.jpeg)

![](_page_61_Figure_7.jpeg)

$$i, j$$
 neural state  $s_i$   
 $a, b;$  factor (dynamical constraint)  
 $a, b, i, j \in \{1, ..., N\}$   
 $u_{a \to j}, m_{j \to a}, r_i, q_i \in [-1, 1]$ 

given two magnetizations h, k, define

$$h \otimes k \stackrel{2}{=} \frac{h+k}{1+hk} = \tanh(\operatorname{atanh}(h) + \operatorname{atanh})$$

 $\bigotimes_{i} u_{i} = \tanh\left(\sum_{i} \operatorname{atanh}\left(u_{i}\right)\right)$ 

$$h \oslash k = h \otimes (-k) = \frac{h-k}{1-hk} = \tanh \left( \operatorname{atanh} \left( h \right) - \operatorname{atanh} \right)$$

on-cavity: 
$$\mu_i = \sum_{a \neq i} w_i^a m_{i \to a}, \ \sigma_i = \sqrt{\sum_{a \neq i} (1 - m_{i \to a}^2)}$$

Cavity:  $\mu_{a \to i} = \mu_i - w_i^a m_{i \to a}, \ \sigma_{a \to i} = \sqrt{\sigma_i^2 - (1 - m_{i \to a}^2)}$ 

![](_page_61_Picture_17.jpeg)

### $\mathrm{h}\left(k ight)$

# Simplified fBP equations

- 1) Approximate with non-cavity fields for each variable,  $h_i$ .
- 2) We update them iteratively; at each iteration we pick a perceptron *j* and have a special update for  $h_i$  (which has the role of the output) vs the others.

$$\begin{split} i \neq j \qquad h_i^{t+1} &= h_i^t + \frac{\left(\frac{1+m_j}{2}\right)\left(H^{++} - H^{+-}\right) + \left(\frac{1-m_j}{2}\right)\left(H^{--} - H^{-+}\right)}{\left(\frac{1+m_j}{2}\right)\left(H^{++} + H^{+-}\right) + \left(\frac{1-m_j}{2}\right)\left(H^{--} + H^{-+}\right)} \\ H^{s_1s_2} &= H\left(-\frac{s_1(\mu - w_j^i m_i) + J\sqrt{N} + s_2w_j^i}{\sigma}\right) \\ i &= j \qquad h_j^{t+1} = h_j^t + \operatorname{atanh}\left(\frac{r^{+} - r^{-}}{r^{+} + r^{-}}\right) \qquad r^{\pm} = H\left(-\frac{\pm \mu + J}{\sigma}\right) \end{split}$$

$$\mu = \sum_{i \neq j} w_j^i m_i, \ \sigma = \sqrt{\sum_{i \neq j} (1 - m_i^2)}$$

$$m_i = \operatorname{hardtanh}(h_i)$$

# Other dynamical schemes that are effective and simple:

- •Gradient Descent on the log-likelihood <sup>(1)</sup>
- Replicated simple dynamics <sup>(2)</sup>

•

•Entropic greedy dynamics ~ entropy-SGD, with m=1 and T=0,  $^{(3)}$ 

- Saglietti, Tartaglione, Zecchina, Phys. Rev. Lett. **120**, 268103 (2018)
- (2020)
- Chayes, Sagun, Zecchina 2017 ICLR, 2019 JSTAT

(1) Role of Synaptic Stochasticity in Training Low-Precision Neural Networks, Baldassi, Gerace, Kappen, Lucibello,

(2) Shaping the learning landscape in neural networks around wide flat minima, Baldassi, Pittorino, ZecchinaPNAS, 117

(3) entropy-SGD: biasing gradient descent into wide valleys, Chaudhari, Choromanska, Soatto, LeCun, Baldassi, Borgs,

![](_page_64_Figure_0.jpeg)

![](_page_64_Figure_1.jpeg)

dist from  $\eta_1$ 

distance matrix between fixed points; geodesic path (rBP)

![](_page_64_Figure_4.jpeg)

step

### Preliminary numerical results on the dynamics:

- 1) Finds liquid-stable fixed points
- 2) These are stable (also for the naive dynamics)
- 3) One can add external field to find fixed points which are close in Hamming distance the external input.
- 4) Extensive errors are corrected (given a fixed point, add an external field in its direction with a fraction of fields flipped and minimum *J*)
- 5) Fixed points further stabilized by Habbian learning  $\eta \eta$  and  $\xi \eta$
- 6) Internal representations are non trivially correlated, e.g. they possess a higher capacity if learned as attractors,  $\alpha_c \simeq 6$  instead of  $\alpha_c = 2$

# Regularization framework for non convex NNs

# Conclusion

# Use ANN convergent dynamics as learning modules