

From Spin Glasses to Machine Learning

Marc Mézard

Bocconi University, Milan

July 3, 2023

Congrès de la Société Française de Physique
Paris

Strongly interacting disordered systems with many components: a long-term perspective

- Maxwell, Boltzmann, etc., 150 years ago, create statistical physics
 - Give up deterministic description
 - Probabilistic approach

Strongly interacting disordered systems with many components: a long-term perspective

- Maxwell, Boltzmann, etc., 150 years ago, create statistical physics
Give up deterministic description
Probabilistic approach
- 50 years ago, creation of a new branch of statistical physics, strongly disordered systems, posing several formidable challenges
Spin glasses. Major developments in the last four decades, starting with Parisi's replica solution of the SK model in 1979.

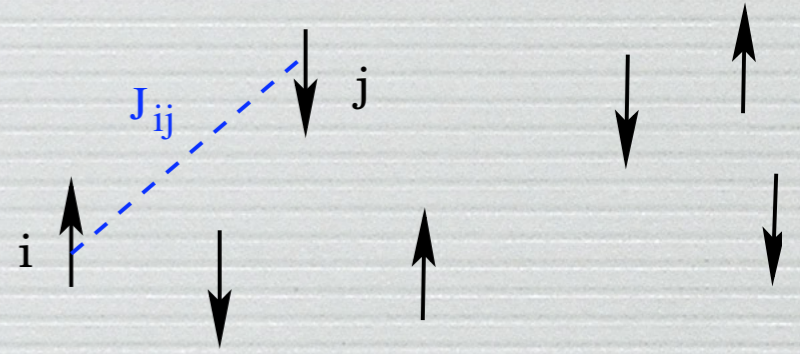
Strongly interacting disordered systems with many components: a long-term perspective

- Maxwell, Boltzmann, etc., 150 years ago, create statistical physics
Give up deterministic description
Probabilistic approach
- 50 years ago, creation of a new branch of statistical physics, strongly disordered systems, posing several formidable challenges
Spin glasses. Major developments in the last four decades, starting with Parisi's replica solution of the SK model in 1979.

Many challenges -> new branch of statistical physics

Challenge 1: ensembles of samples

One sample of a spin glass = set of couplings J between $N \gg 1$ spins.
Boltzmann probability measure on the spins $P_J(S)$



$$s_i = \pm 1$$

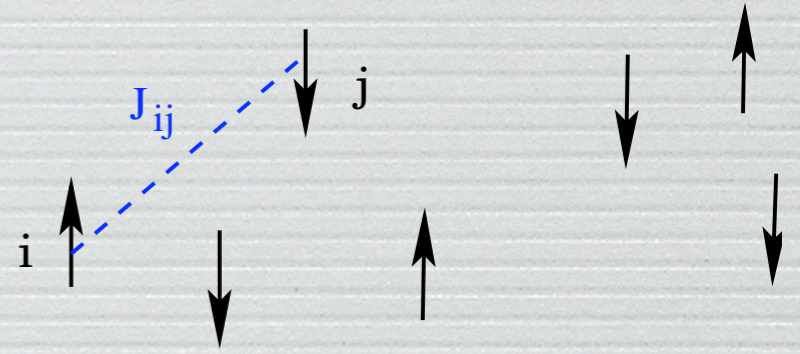
$$J = \{J_{ij}\}$$

$$E_J(S) = - \sum_{(i,j)} J_{ij} s_i s_j$$

$$P_J(S) = \frac{1}{Z_J} e^{-\beta E_J(S)}$$

Challenge 1: ensembles of samples

One sample of a spin glass = set of couplings J between $N \gg 1$ spins.
Boltzmann probability measure on the spins $P_J(S)$



Ensemble of samples = probability distribution on the set of couplings $\mathcal{P}(J)$

$$s_i = \pm 1$$

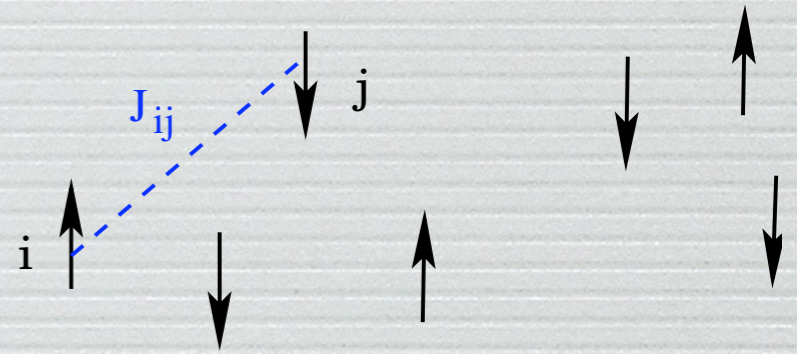
$$J = \{J_{ij}\}$$

$$E_J(S) = - \sum_{(i,j)} J_{ij} s_i s_j$$

$$P_J(S) = \frac{1}{Z_J} e^{-\beta E_J(S)}$$

Challenge 1: ensembles of samples

One sample of a spin glass = set of couplings J between $N \gg 1$ spins.
Boltzmann probability measure on the spins $P_J(S)$



Ensemble of samples = probability distribution on the set of couplings $\mathcal{P}(J)$

Generate a sample with probability $\mathcal{P}(J)$
What are the properties of spin configurations sampled from $P_J(S)$?

$$s_i = \pm 1$$

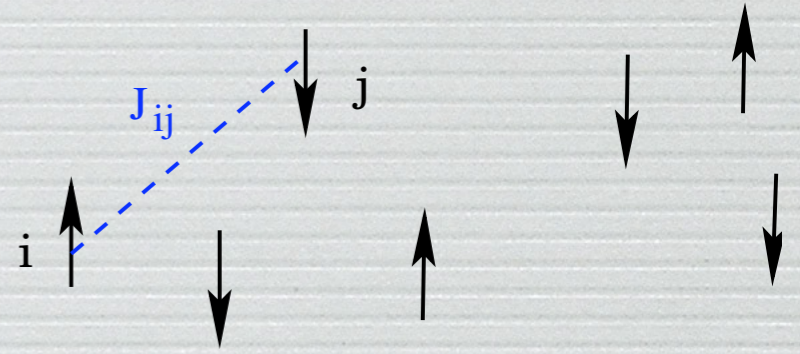
$$J = \{J_{ij}\}$$

$$E_J(S) = - \sum_{(i,j)} J_{ij} s_i s_j$$

$$P_J(S) = \frac{1}{Z_J} e^{-\beta E_J(S)}$$

Challenge 1: ensembles of samples

One sample of a spin glass = set of couplings J between $N \gg 1$ spins.
Boltzmann probability measure on the spins $P_J(S)$



Ensemble of samples = probability distribution on the set of couplings $\mathcal{P}(J)$
Generate a sample with probability $\mathcal{P}(J)$
What are the properties of spin configurations sampled from $P_J(S)$?

$$s_i = \pm 1$$

$$J = \{J_{ij}\}$$

$$E_J(S) = - \sum_{(i,j)} J_{ij} s_i s_j$$

$$P_J(S) = \frac{1}{Z_J} e^{-\beta E_J(S)}$$

Quenched disorder: each sample is different.

Thermal disorder: in a given sample, spins fluctuate.

Challenge 1: ensembles of samples

Disorder: each sample is different.
Study sample **ensembles**. Find « self-averaging » quantities, which are identical in almost all samples.
Understand differences (between samples)

Self-averaging:

$$N \rightarrow \infty \quad \frac{1}{N} \sum_i \langle s_i \rangle \quad \frac{1}{N} \langle E_J(S) \rangle$$

Sample dependent: details of the landscape, ground state

eg Sherrington
Kirkpatrick
model

$$J_{ij} \sim \mathcal{N}(0, 1/N)$$

$$E_J(S) = - \sum_{(i,j)} J_{ij} s_i s_j$$

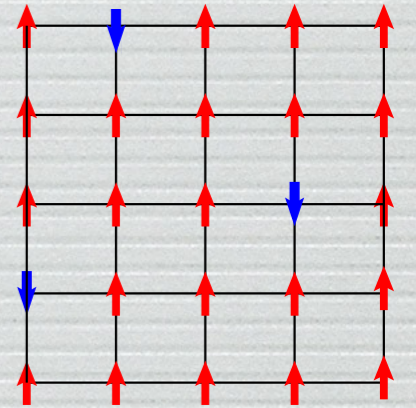
$$P_J(S) = \frac{1}{Z_J} e^{-\beta E_J(S)}$$

Challenge 2: inhomogeneity

Every spin is in a different environment.

Different magnetizations.

No « representative agent ».



Mean-field equations = N coupled equations for the local magnetizations (Thouless, Anderson, Palmer 1976)

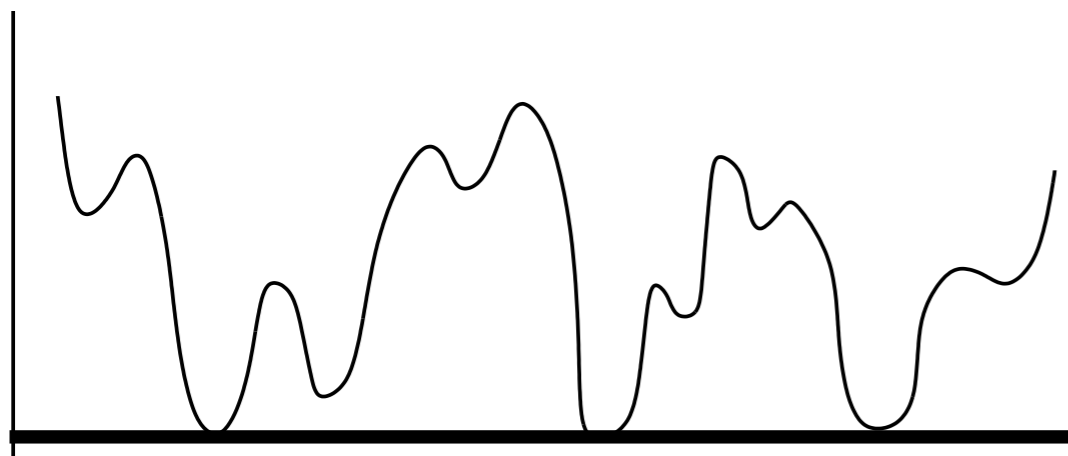
Major simplification from a probability over 2^N configurations

Statistical description of the magnetizations, the local fields:
cavity method (M, Parisi, Virasoro 1986)

Challenge 3: rough landscape

Complicated landscape, many states in which the spin system can freeze. In SK: hierarchical (ultrametric) structure (MPSTV 85)

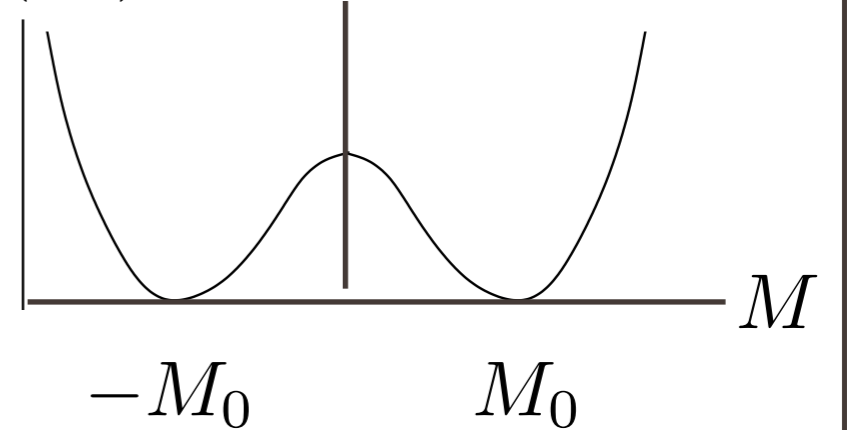
Energy per spin



(sketch in a N -dimensional space)

Spin glass

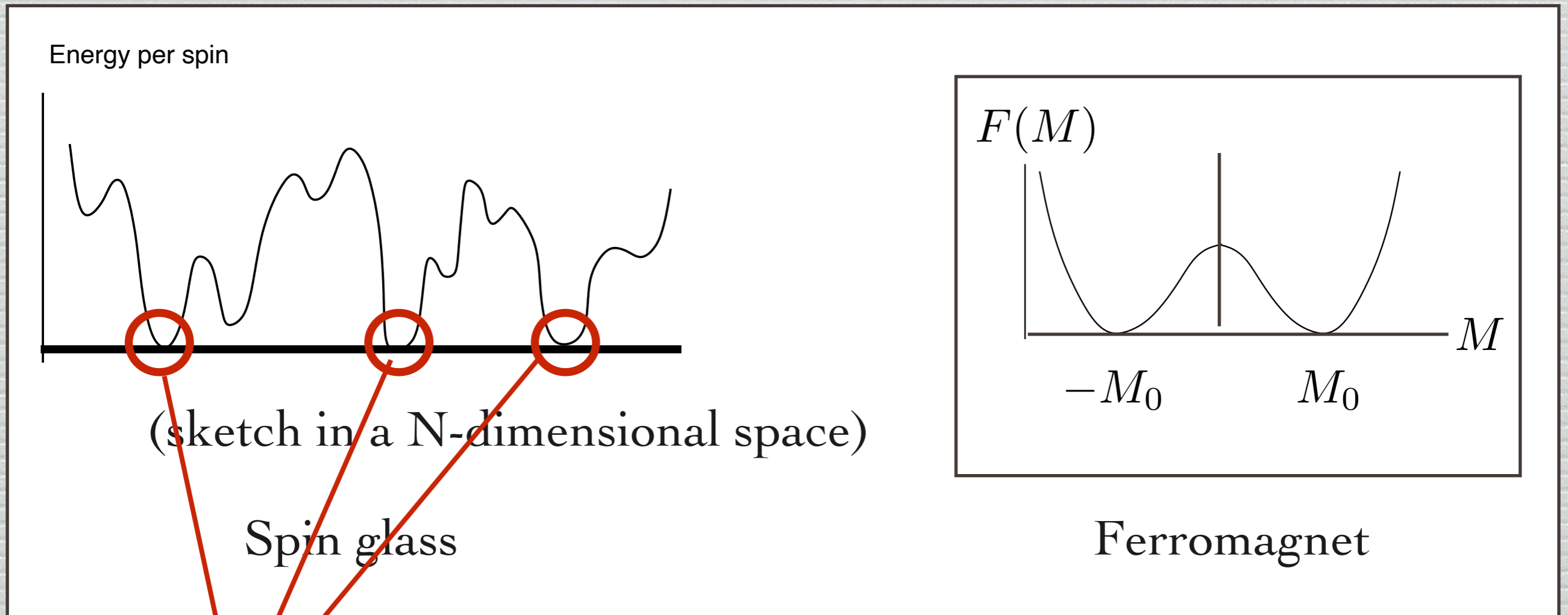
$F(M)$



Ferromagnet

Challenge 3: rough landscape

Complicated landscape, many states in which the spin system can freeze. In SK: hierarchical (ultrametric) structure (MPSTV 85)



Details of the landscape depend on the sample !

A new branch of statistical physics

- ➔ Study ensembles of problems
- ➔ Each spin 'sees' a different local field
 - Spins freeze in random directions
- ➔ Rough landscape: difficult to find min. of E
- ➔ Strong out of equilibrium dynamical effects

NB : beyond the simple mean field theory of the « representative agent »:
Statistics of agents. **Replicas, cavity...**

A new branch of statistical physics

- ➔ Study ensembles of problems
- ➔ Each spin 'sees' a different local field
 - Spins freeze in random directions
- ➔ Rough landscape: difficult to find min. of E
- ➔ Strong out of equilibrium dynamical effects

NB : beyond the simple mean field theory of the « representative agent »:
Statistics of agents. **Replicas, cavity...**

Useless, but « cornucopia »...

SK= Generic model of binary variables interacting by pairs

Machine Learning and Large Dimensional Inference

Machine learning going deep: a decade of technological revolution

1- Image understanding.

In the last ten years, detection, segmentation and recognition of objects and regions in images. Image generation.

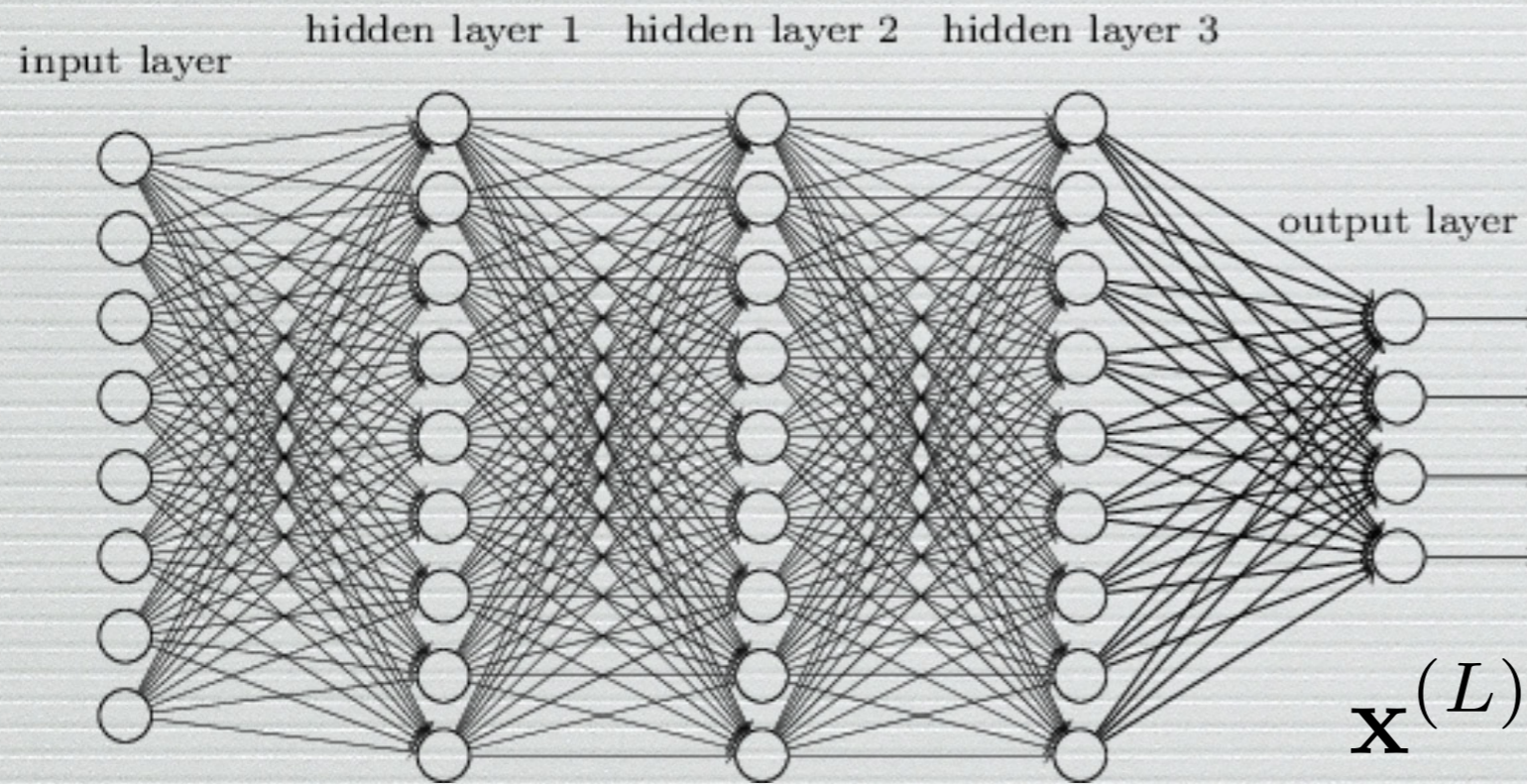
2- Language analysis: topic classification, question answering, language translation. Language generation

3- Science. Protein Folding. Predicting the activity of potential drug molecules. Algorithmic speedup, feature detection in data, quantum computing...

4- Playing games (chess, go, poker, video-games,...)
etc.

waiting for a general theoretical framework

The tool: Deep neural network



$$\mathbf{X}^{(1)} \quad \mathbf{X}^{(2)}$$

$$\mathbf{X}^{(n+1)} = f \left(\mathbb{W}^{(n)} \mathbf{X}^{(n)} \right)$$

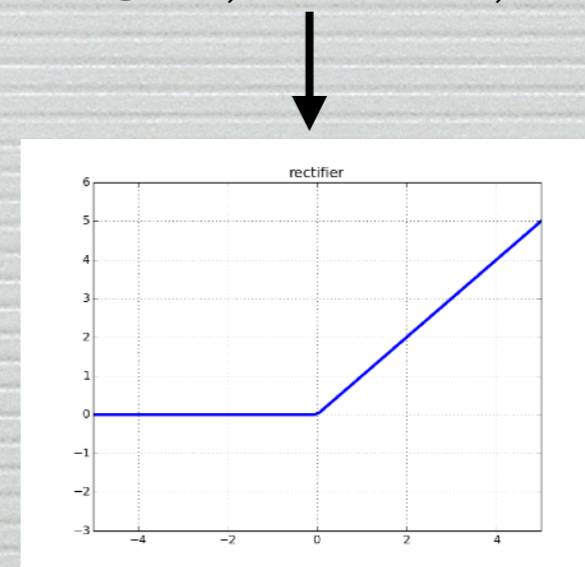
Artificial neuron

$$x_i^{(n+1)} = f \left(\sum_j \mathbb{W}_{ij}^{(n)} x_j^{(n)} \right)$$

NB : component-wise nonlinearity

Parameters to be learnt: weights \mathbb{W}

$f = \text{Sign}, \text{Relu}, \text{tanh} \dots$



Machine learning



$$\xi \in \mathbb{R}^N \rightarrow y = f(W, \xi) \begin{cases} \in \mathbb{R} & \text{or} \\ \in \{0, 1, \dots, q\} \end{cases}$$

Database = P examples of input-output $(\xi_\mu, y_\mu) \quad \mu = 1, \dots, P$

Machine learning



$$\xi \in \mathbb{R}^N \rightarrow y = f(W, \xi) \begin{cases} \in \mathbb{R} & \text{or} \\ \in \{0, 1, \dots, q\} \end{cases}$$

Database = P examples of input-output

$$(\xi_\mu, y_\mu) \quad \mu = 1, \dots, P$$

Desired label (« supervised learning »)

Machine learning



$$\xi \in \mathbb{R}^N \rightarrow y = f(W, \xi) \begin{cases} \in \mathbb{R} & \text{or} \\ \in \{0, 1, \dots, q\} \end{cases}$$

Database = P examples of input-output $(\xi_\mu, y_\mu) \quad \mu = 1, \dots, P$

Desired label (« supervised learning »)

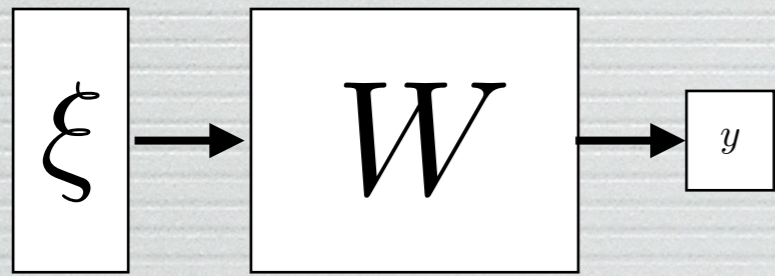
Learning = Optimization

Find W^* that minimizes the training error:
(or other « loss function »)

$$\sum_{\mu=1}^P [f(W, \xi_\mu) - y_\mu]^2$$

Example stochastic gradient descent Very large dimensional landscape.

Machine learning



$$\xi \in \mathbb{R}^N \rightarrow y = f(W, \xi) \begin{cases} \in \mathbb{R} & \text{or} \\ \in \{0, 1, \dots, q\} \end{cases}$$

Database = P examples of input-output $(\xi_\mu, y_\mu) \quad \mu = 1, \dots, P$

Desired label (« supervised learning »)

Learning = Optimization

Find W^* that minimizes the training error:
(or other « loss function »)

$$\sum_{\mu=1}^P [f(W, \xi_\mu) - y_\mu]^2$$

Example stochastic gradient descent Very large dimensional landscape.

The big Challenge: Generalization

Use the optimal W^* , test the machine on new data

Machine learning: learning phase



$$\xi \in \mathbb{R}^N \rightarrow y = f(W, \xi) \begin{cases} \in \mathbb{R} & \text{or} \\ \in \{0, 1, \dots, q\} \end{cases}$$

Database = P examples of input-output (ξ_μ, y_μ) $\mu = 1, \dots, P$

Desired label (« supervised learning »)

Machine learning: learning phase



$$\xi \in \mathbb{R}^N \rightarrow y = f(W, \xi) \begin{cases} \in \mathbb{R} & \text{or} \\ \in \{0, 1, \dots, q\} \end{cases}$$

Database = P examples of input-output (ξ_μ, y_μ) $\mu = 1, \dots, P$

Desired label (« supervised learning »)

Bayesian learning:

$$P(W | \{\xi_\mu, y_\mu\}) = \frac{1}{Z} P^0(W) \exp \left(-\beta \sum_{\mu} [f(W, \xi_\mu) - y_\mu]^2 \right)$$

Unknown Data Prior Loss

Effective inverse temperature allows to tune the importance of data wrt prior (annealing)

Machine learning: learning phase

Disordered system. Database = sample = disorder. For each database, study the properties of the probability measure on the weights

- Specific database, MNIST, CIFAR, etc
- Statistical ensemble of database. Generative models

Bayesian learning:

$$P(W | \{\xi_\mu, y_\mu\}) = \frac{1}{Z} P^0(W) \exp \left(-\beta \sum_{\mu} [f(W, \xi_\mu) - y_\mu]^2 \right)$$

Unknown Data Prior Loss

Effective inverse temperature allows to tune the importance of data wrt prior (annealing)

The (old) ingredients

- * Feedforward neural networks
- * Trained with gradient descent learning, implemented with gradient back propagation

What is new in practice since the 80's ?

- * Availability of very large data bases
- * Much larger computing power
- * Much deeper networks
- * Numerous « tricks »:
 - Accumulated experience on structures (depth, width).
 - First layers = local convolutions
 - Activation functions (ReLU)
 - Stochastic gradient descent
 - Early stopping
 - Transfer learning
 - ...

The (old) ingredients

- * Feedforward neural networks
- * Trained with gradient descent learning, implemented with gradient back propagation

What is new in practice since the 80's ?

- * Availability of very large data bases
- * Much larger computing power
- * Much deeper networks
- * Numerous « tricks »:

- Accumulated experience on structures (depth, width).
- First layers = local convolutions
- Activation functions (ReLU)
- Stochastic gradient descent
- Early stopping
- Transfer learning
- ...

* Generative models

Surprises and questions

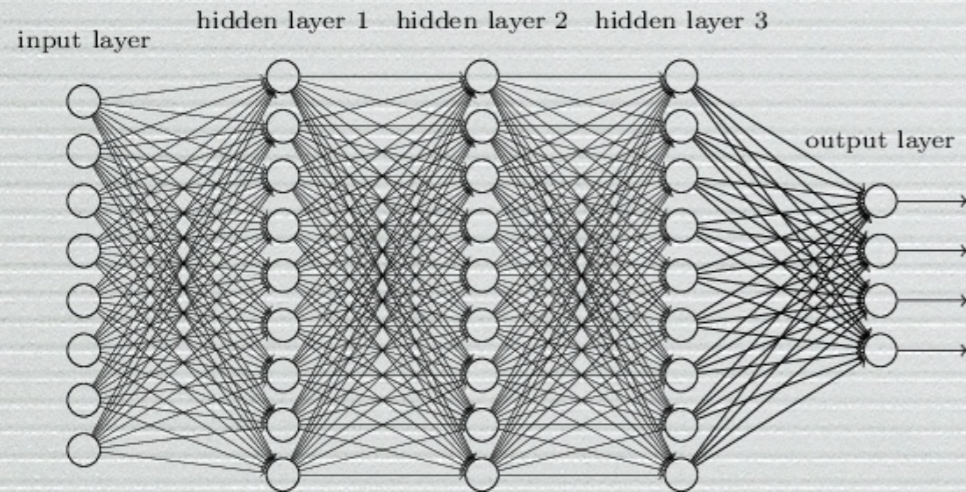
Surprises and questions

Training

Generalization

Mechanism

Surprises and questions



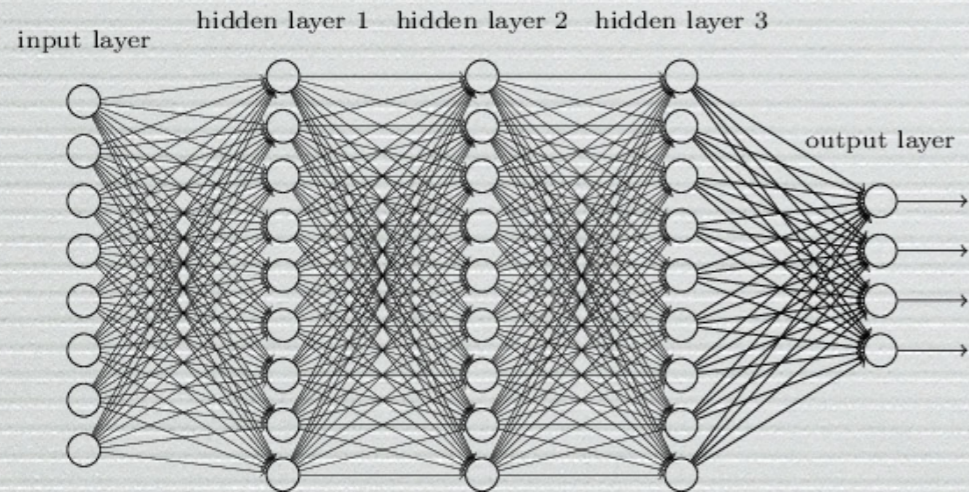
Training

Generalization

Mechanism

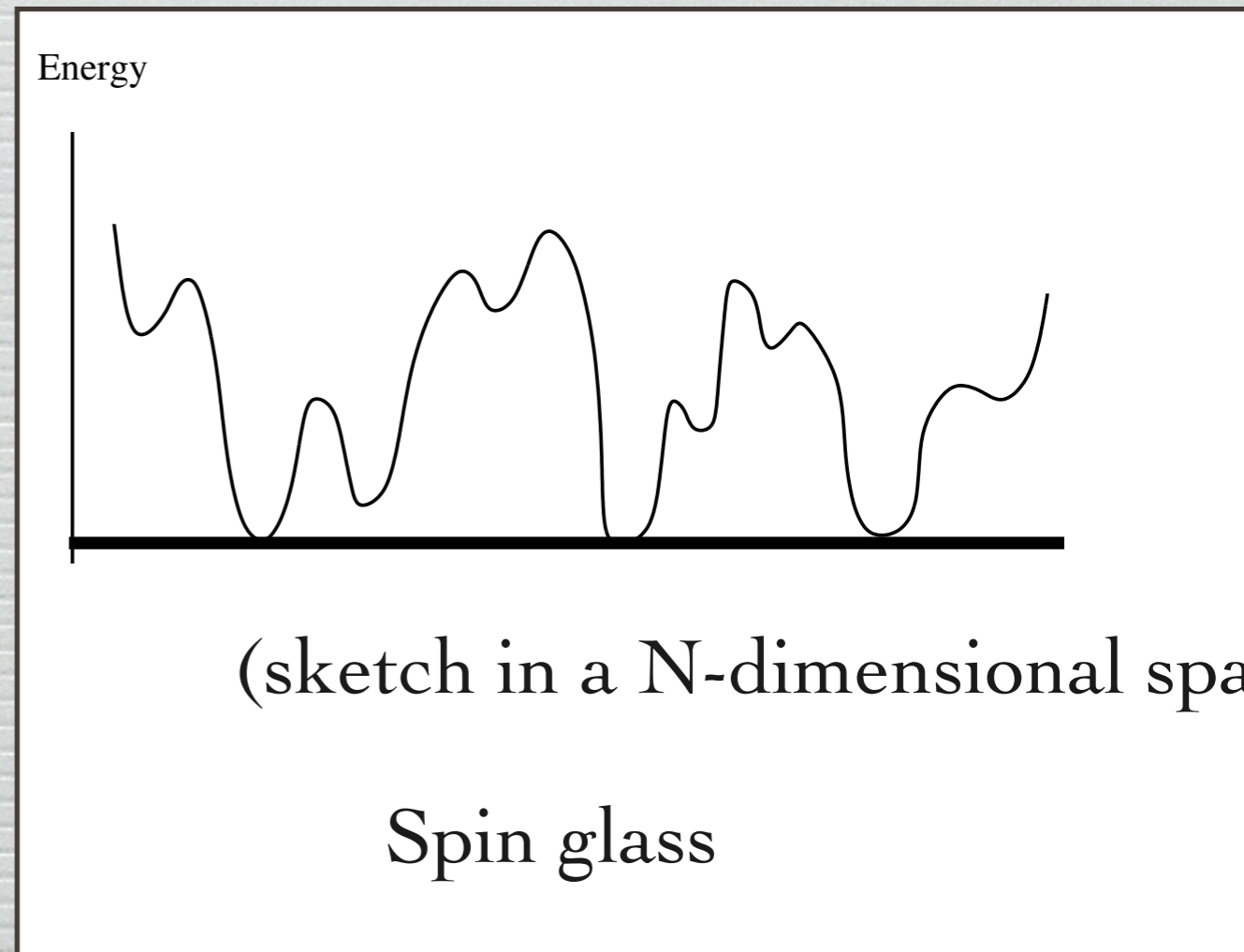
Training = optimization of a disordered system in a large dimensional space

Surprises and questions

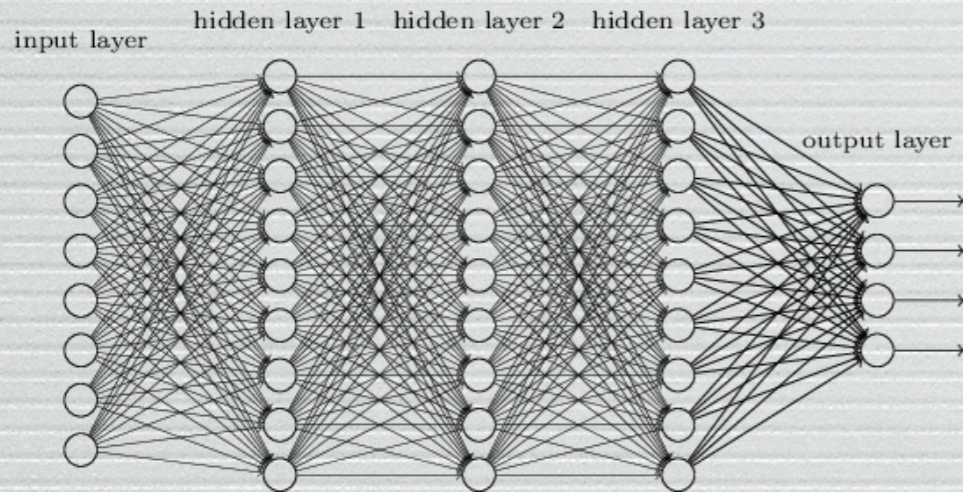


Training
Generalization
Mechanism

Training = optimization of a disordered system in a large dimensional space



Surprises and questions

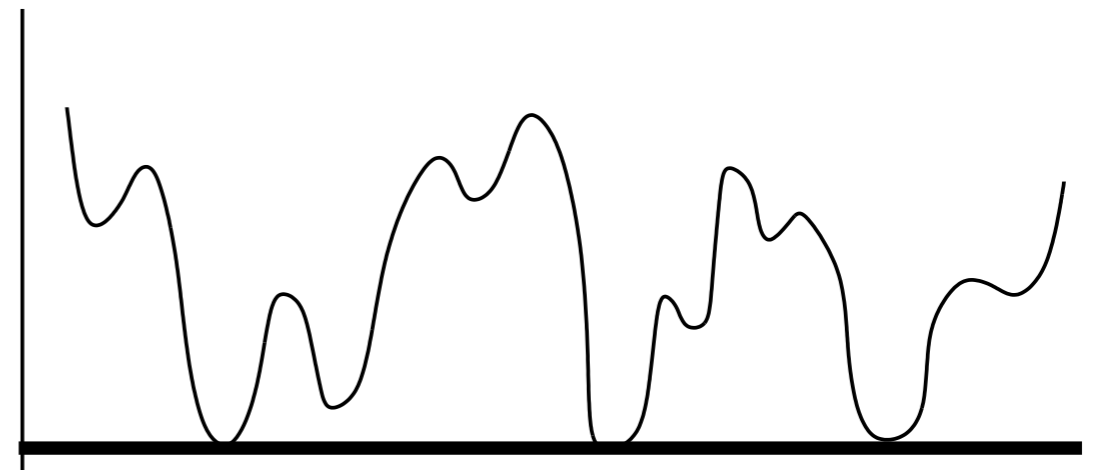


Training
Generalization
Mechanism

Training = optimization of a disordered system in a large dimensional space

Experimentally: one can reach zero training error, using simple stochastic gradient descent, in the neighborhood of any random starting point provided the network is deep enough

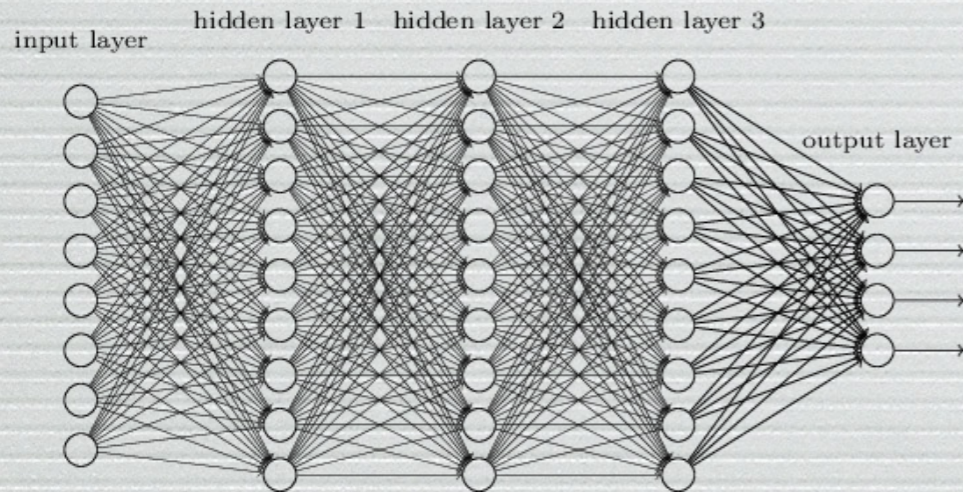
Energy



(sketch in a N-dimensional space)

Spin glass

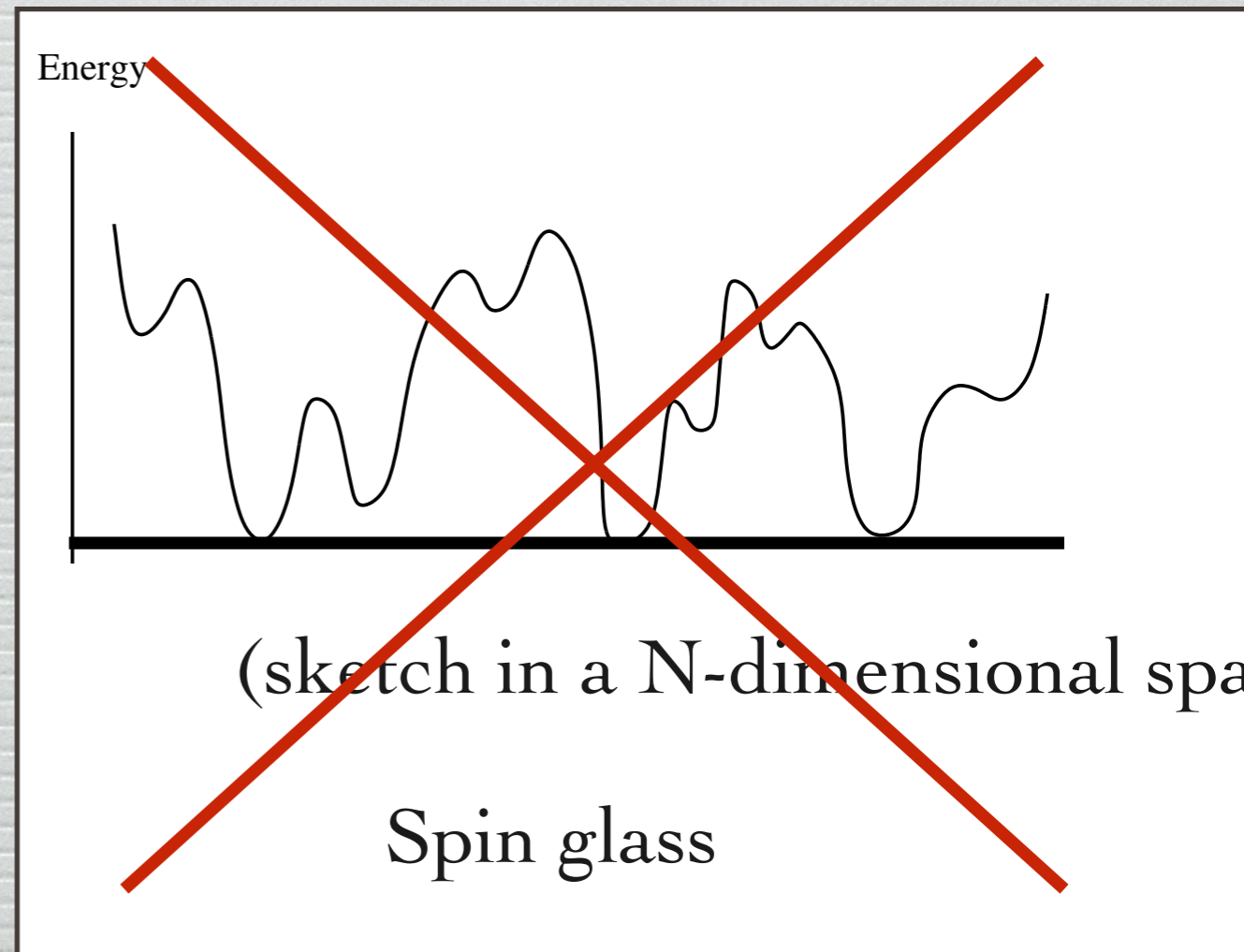
Surprises and questions



Training
Generalization
Mechanism

Training = optimization of a disordered system in a large dimensional space

Experimentally: one can reach zero training error, using simple stochastic gradient descent, in the neighborhood of any random starting point provided the network is deep enough

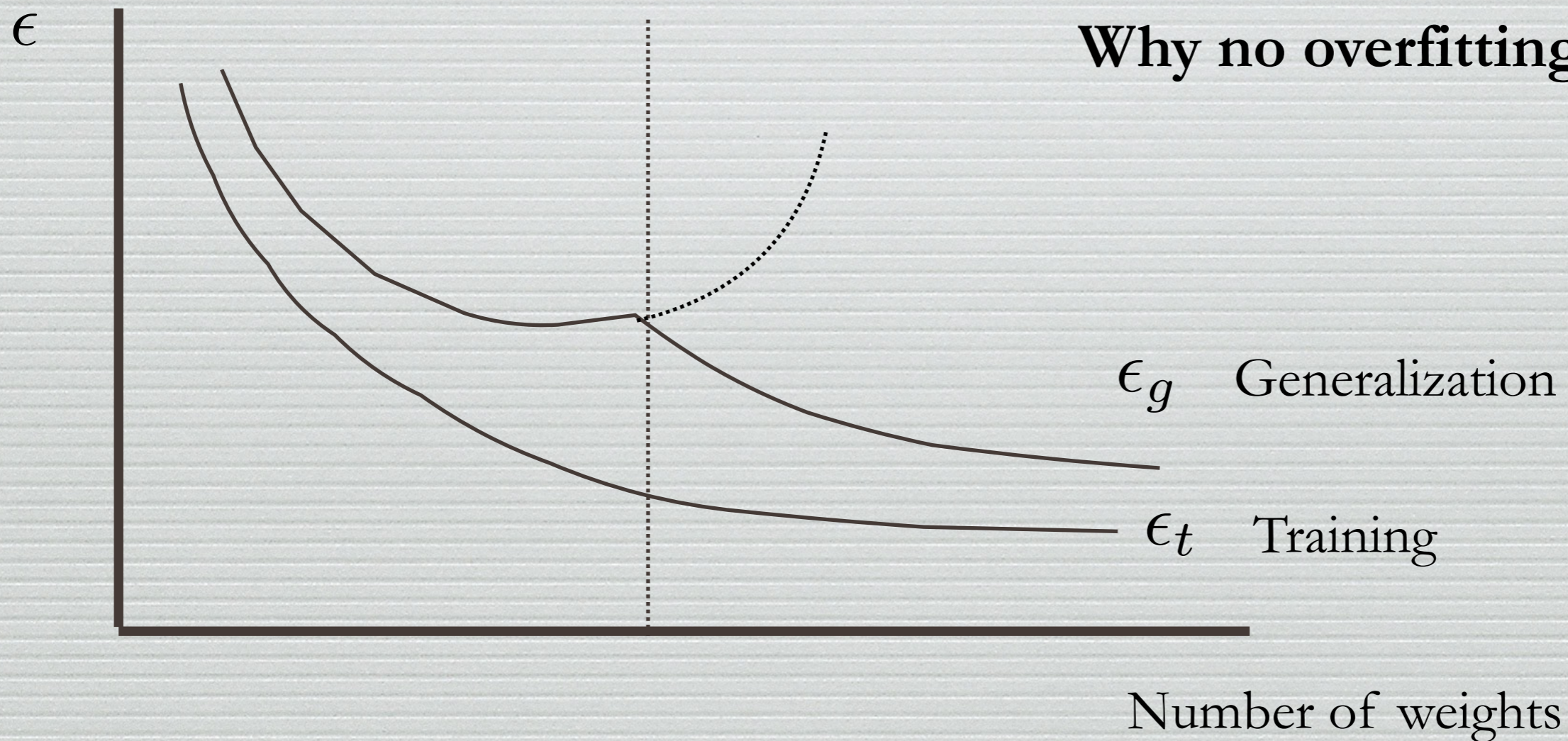


Surprises and questions

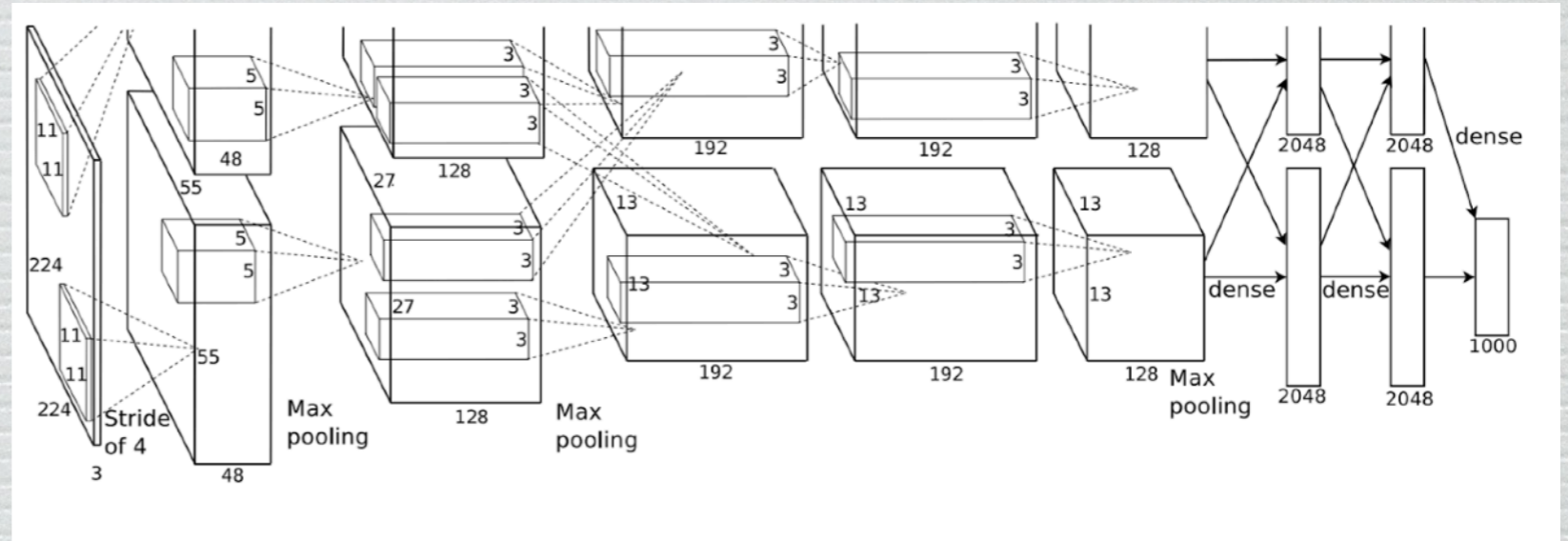
Generalization :

We train with billions of parameters.

Why no overfitting?



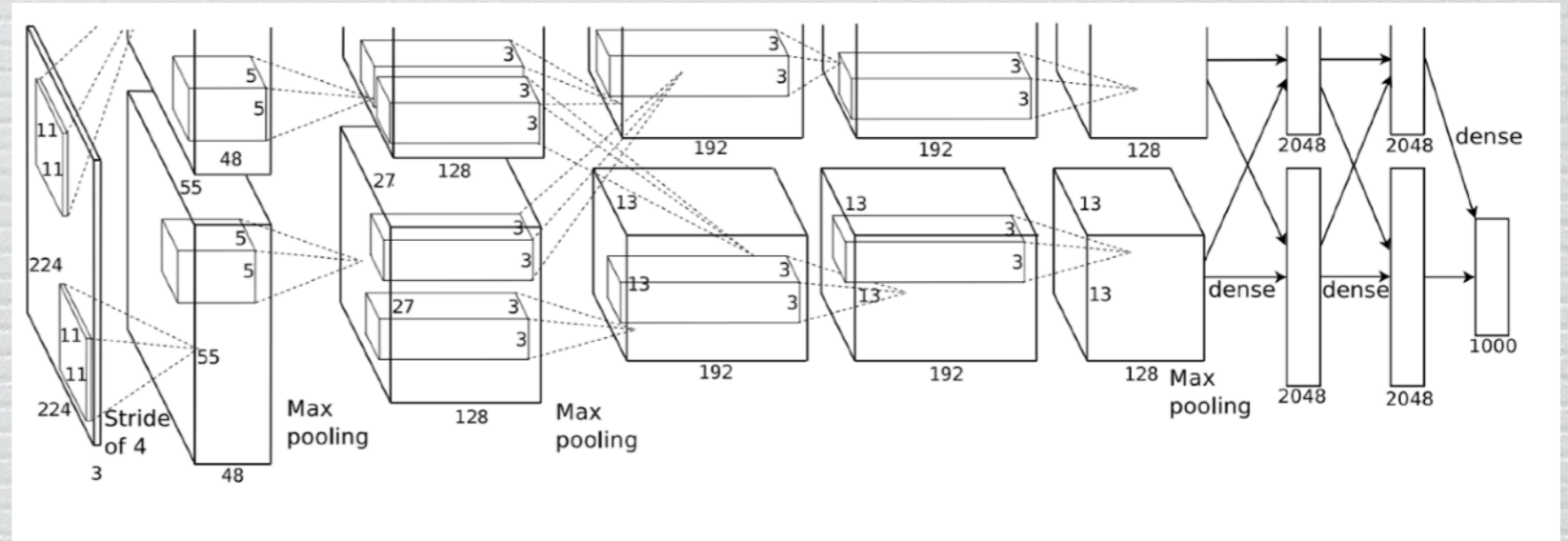
Surprises and questions



We know everything of the trained network
(neuroscientist's dream)

We do not understand much. **Emergent phenomenon**

Surprises and questions



We know everything of the trained network
(neuroscientist's dream)

We do not understand much. **Emergent phenomenon**

No guarantee

No explanation

Ingredients of deep networks

Architecture

Art. Go deep, use convolutions in first layers, use pooling, etc...

Learning algorithms

Art. The (nearly) most naive algorithm, stochastic gradient descent initialized with small weights

« Simple » Data structure

Maybe the tasks that machine learning addresses are easier than expected because data has a lot more structure than our theories (worst case, or typical case with iid data) used so far

The Challenge of Data Structure

Combinatorial

Hierarchical

Semantic

Low-dimensional Manifold

RFI 29/06

Nouvelle nuit de violences en France après la mort de Nahel: le récit des dernières heures

Des violences ont éclaté en marge de la marche blanche organisée ce jeudi en hommage à Nahel. Ce jeudi soir, les violences se sont répandues et des incidents ont éclaté dans des villes aux quatre coins du pays. 40 000 forces

Semantic: eg Large Language Models

Des **violences** ont éclaté en marge de la marche blanche organisée ce jeudi en hommage à Nahel. Ce jeudi soir, les violences **se sont répandues** et des incidents ont éclaté dans des **villes aux quatre coins du pays. 40 000 forces**

Attention Mechanism

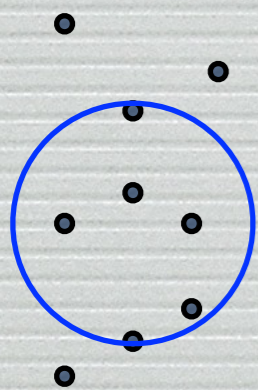
Hidden manifold example: MNIST



Input space: dimension

$$28^2 = 784$$

Manifold of handwritten digits in MNIST:

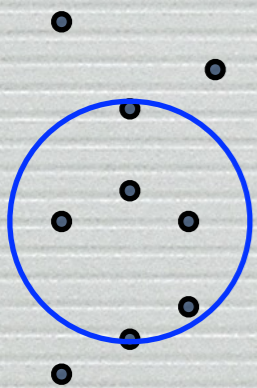
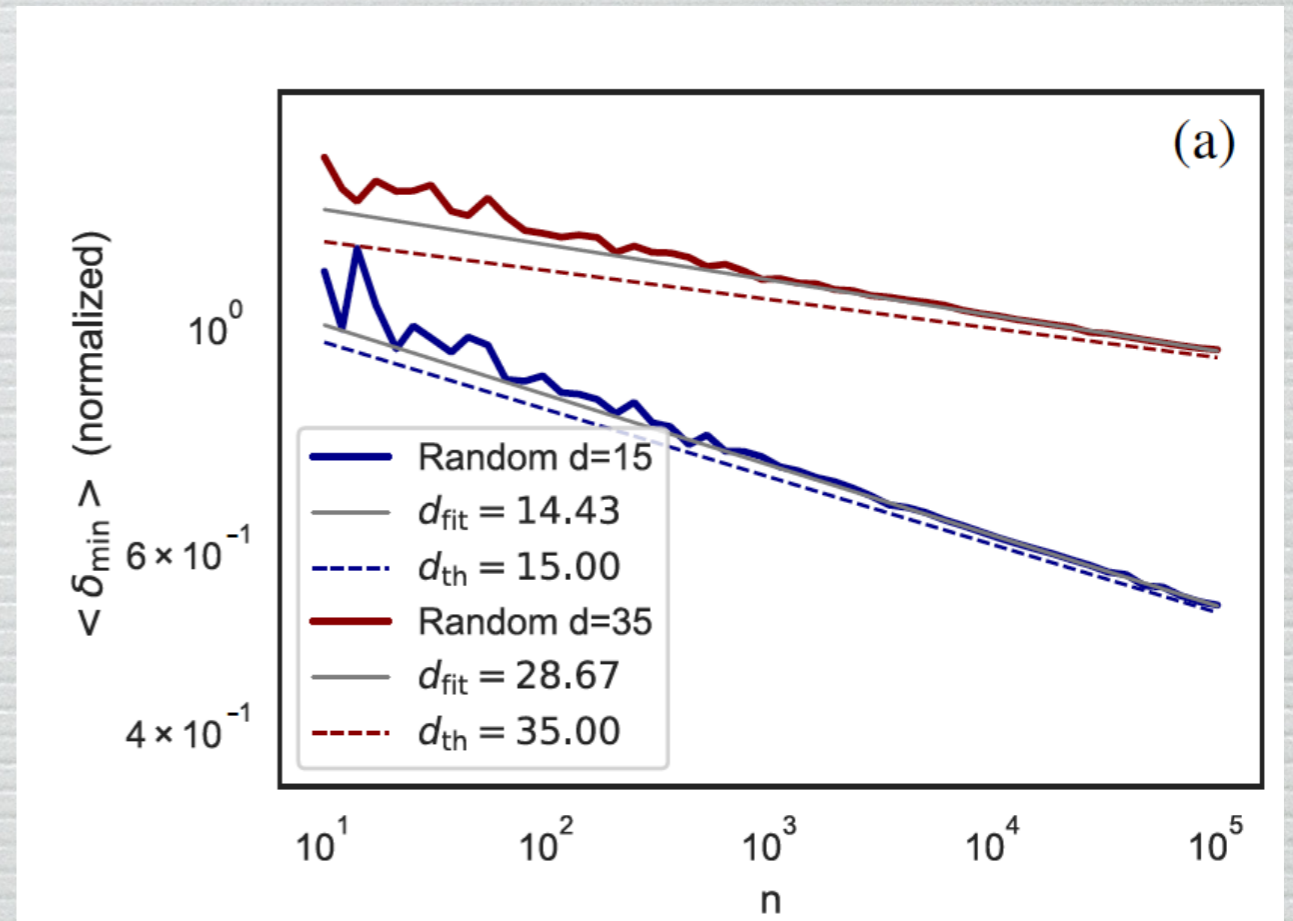


$$p \simeq cR^d$$

Nearest neighbors' distance : $R_{nn} \simeq p^{-1/d}$

Grassberger Procaccia 83, Costa Hero 05, Heinz Audibert 05, Facco et al. 17, Ansuini et al. 19, Spigler et al. 19...

Hidden manifold example: MNIST



$$p \simeq cR^d$$

Nearest neighbors' distance : $R_{nn} \simeq p^{-1/d}$

$$d_{\text{eff}} \simeq 15$$



$$d_{\text{eff}}(5) \simeq 12$$

Table 7. Number of samples and estimated intrinsic dimensionality of the digits in MNIST.

1	2	3	4	5
7877	6990	7141	6824	6903
8/7/7	13/12/13	14/13/13	13/12/12	12/12/12
6	7	8	9	0
6876	7293	6825	6958	6903
11/11/11	10/10/10	14/13/13	12/11/11	12/11/11

Hein Audibert 05

MNIST problem: in the 15-dim manifold of handwritten digits, identify the 10 perceptual submanifolds associated with each digit, of dimensions between 7 and 13...



$$d_{\text{eff}}(5) \simeq 12$$

Table 7. Number of samples and estimated intrinsic dimensionality of the digits in MNIST.

1	2	3	4	5
7877	6990	7141	6824	6903
8/7/7	13/12/13	14/13/13	13/12/12	12/12/12
6	7	8	9	0
6876	7293	6825	6958	6903
11/11/11	10/10/10	14/13/13	12/11/11	12/11/11

Hein Audibert 05

MNIST problem: in the 15-dim manifold of handwritten digits, identify the 10 perceptual submanifolds associated with each digit, of dimensions between 7 and 13...

... from an input in 784 dimensions!

A brief conclusion

Statistical physics of disordered systems has been a major evolution of statistical physics in the last 50 years, opening the way to applications in many other branches of science

When using statistical physics ideas in the study of machine learning, one must be able to face a new challenge, the one of structured data (manifolds, attention, etc.)

The End