Colloquium du CERMICS



#### Mathematical Mysteries of Deep Neural Networks

Stéphane Mallat (Collège de France et École Normale Supérieure)

20 septembre 2019



John Zarka, Sixin Zhang

Collège de France École Normale Supérieure

## High-Dimensional Approximations

What regularity properties lead to low-dimensional approximations of f(x) for a high-dimensional  $x \in \mathbb{R}^d$  in physics and machine learning ?

or

• f(x): class of an image x having  $d = 10^6$  pixels



energy of a physical system in a state  $x \in \mathbb{R}^d$ 



• f(x) = p(x) a probability density.





## **Curse of Dimensionality**

• f(x) can be approximated from examples  $\{x_i, f(x_i)\}_i$  by local interpolation if f is regular and there are close examples:



- Need  $n \ge \epsilon^{-d}$  points to cover  $[0, 1]^d$  at a Euclidean distance  $\epsilon$ Problem:  $||x - x_i||$  is always large
- To estimate f(x) when x is in a high-dimensional  $\Omega$ requires strong regularity of f in  $\Omega$ : what regularity ?

#### • Deep Convolutional Network • Deep convolutional neural network to predict y = f(x): $x \in \mathbb{R}^d$ $\rho(a) = \max(a, 0)$ $\rho(L_1)$ $\rho(L$

 $L_j$ : spatial convolutions and linear combination of channels Exceptional results for classification of *images, sounds, language, regressions in physics, signal and image generation...* but not interpretable.

To create simpler interpretable networks:

What underlying regularity is captured and how ?

3 ingredients: Multiscale, Linearize group actions, Sparse

## Scale Separation and Interactions

#### • Dimension reduction:

Interactions de d bodies represented by x(u): particles, pixels...



Interactions across scales

Multiscale regroupement of interactions of d bodies into interactions of  $O(\log d)$  groups.

Scale separation  $\Rightarrow$  wavelet transforms.

How to capture scale interactions ?

Critical harmonic analysis problems since 1970's

## Overview: Simpler Networks

- Scale separation with wavelets and interactions through phase
- Linear scale interaction models and invariants in:
  - Statistical physics for turbulence
  - Quantum chemistry and image classification
- Non-linear scale interactions models with sparse dictionaries in:
  - Classification of complex structures as in ImageNet
  - Generation of non-ergodic processes

## **Scale separation with Wavelets**

• Wavelet filter  $\psi(u)$ : = +i =

rotated and dilated:  $\psi_{\lambda}(u) = 2^{-2j} \psi(2^{-j}r_{\theta}u)$ 



• Wavelet transform: invertible

$$Wx = \left(x \star \psi_{\lambda}\right)_{\lambda}$$



• Zero-mean and no correlations across scales: problem!

$$\sum_{u} x \star \psi_{\lambda}(u) x \star \psi_{\lambda'}^{*}(u) = \sum_{\omega} |\widehat{x}(\omega)|^{2} \psi_{\lambda}(\omega) \psi_{\lambda}(\omega)^{*} \approx 0 \quad \text{if} \quad \lambda \neq \lambda'$$



## **Stat.** Physics of Stationary Proc.

What stochastic models for turbulence ?

$$d = 6 \, 10^4$$



Prior: stationary  $\Leftrightarrow p(x)$  is invariant to translations.

Maximum entropy distribution  $\tilde{p}$  conditioned by M moments

$$\mathbb{E}(\phi_m(x)) = \mu_m \quad \Rightarrow \quad \tilde{p}(x) = \mathcal{Z}^{-1} \exp\left(-\sum_{m=1}^M \beta_m \phi_m(x)\right)$$

With M = d second order moments:

 $\phi_m(x) = \sum_u x(u)x(u-m) \implies \tilde{p}(x)$  is a Gaussian distribution

## **Gaussian Models of Stationary Proc.**

What stochastic models for turbulence ?

$$x = 6 \, 10^4$$

 $\tilde{p}(x)$  is a Gaussian distribution  $\tilde{x}$ 





No correlation is captured across scales and frequencies. Random phases.

How to capture non-Gaussianity and long range interactions? Failure of high order moments. Deep net generations look better. Rectifiers act on Phase

Real wavelets of phase α: ψ<sub>α,λ</sub> = Real(e<sup>-iα</sup> ψ<sub>λ</sub>) Rectifier: ρ(a) = max(a, 0) Ux(u, α, λ) = ρ(x \* Real(e<sup>iα</sup> ψ<sub>λ</sub>)) = ρ(Real(e<sup>iα</sup> x \* ψ<sub>λ</sub>)) x \* ψ<sub>λ</sub> = |x \* ψ<sub>λ</sub>| e<sup>iφ(x\*ψ<sub>λ</sub>)</sup> Homogeneous: ρ(αa) = α ρ(a) if α > 0

$$Ux(u,\alpha,\lambda) = |x \star \psi_{\lambda}| \rho(\cos(\alpha + \varphi(x \star \psi_{\lambda})))$$

#### A Relu computes phase harmonics:

ENS

**Theorem**: Fourier transform along the phase  $\alpha$ :  $\widehat{U}x(u,k,\lambda) = \widehat{\gamma}(k) |x \star \psi_{\lambda}(u)| e^{ik \varphi(x \star \psi_{\lambda}(u))}$ with  $\gamma(\alpha) = \rho(\cos \alpha)$  for any homogeneous non-linearity  $\rho$ .

## **Frequency Transpositions**

Phase harmonics:  $|x \star \psi_{\lambda}(u)| e^{i k \varphi(x \star \psi_{\lambda}(u))}$ 

Performs a non-linear frequency dilation / transposition



Phase Harmonics

Correlated if  $k\lambda \approx \lambda'$ 



#### Real wavelets: $\psi_{\alpha,\lambda} = \operatorname{Real}(e^{-i\alpha}\psi_{\lambda})$ and $\rho(a) = \max(a,0)$





Sixin Zhang



Maximum entropy distribution conditioned by  $M = O(\log^2 d) \text{ wavelet harmonic correlations } \mathbb{E}(\phi_m(x))$   $\phi_m(x) = \sum_u |x \star \psi_\lambda(u)| \, e^{ik\varphi(x \star \psi_\lambda(u))} \, |x \star \psi_{\lambda'}(u)| \, e^{-ik'\varphi(x \star \psi_{\lambda'}(u))}$   $\tilde{p}(x) = \mathcal{Z}^{-1} \, \exp\left(-\sum_{m=1}^M \beta_m \, \phi_m(x)\right)$ 

 ${\mathcal X}$ 

Ergodic Stationary Processes

S. Zhang, J. Bruna, E. Allys, F. Levrier, F. Boulanger  $d = 6 \, 10^4$ 



 $M = 3 \, 10^3$ number of moments

Phase coherence is restored How much physics are these models capturing? What about non-ergodic processes?

 $\tilde{x}$ 

Scattering Wavelet Coefficients

Classification: invariance by translation by spatial averaging

$$\left(\begin{array}{c} x \star \phi_{2^{J}}(2^{J}n) \\ \rho(x \star \psi_{\alpha,\lambda}) \star \phi_{J}(2^{J}n) \end{array}\right)_{\alpha,\lambda}$$

Recover the information loss with a second layer:

$$S_J x = U x \star \phi_J = \left(\begin{array}{c} x \star \phi_{2J}(2^J n) \\ \rho(x \star \psi_{\alpha,\lambda}) \star \phi_J(2^J n) \\ \rho(\rho(x \star \psi_{\alpha,\lambda}) \star \psi_{j',\alpha'}) \star \phi_J(2^J n) \end{array}\right)_{\alpha,\lambda,\alpha',\lambda'}$$

#### - Linearize small deformations

**Theorem** if  $D_{\tau}x(u) = x(u - \tau(u))$  then  $\lim_{J \to \infty} \|S_J D_{\tau}x - S_J x\| \le C \|\nabla \tau\|_{\infty} \|x\|$  Quantum Chemistry: N-Body Problem

• Can we learn the interaction energy f(x) of a system with  $x = \{ \text{positions, charges} \}$ ?

Symmetries:



f(x) is invariant to translations and rotations,

multiscale interactions: chemical bounds, Van der Waal forces...

The energy depends upon the electronic density (Kohn-Sham)

Ground state electronic density computed with Schroedinger







- We do not know the electronic density at equilibrium.
- The molecular state  $\{z_k, r_k\}_{k \leq d}$  is represented by Diracs located at  $r_k$  weighted by charges  $z_k$ :

$$x(u) = \sum_{k=1}^{d} z_k \,\delta(u - r_k)$$

Electronic density

Dirac density x(u)





## Harmonic Wavelet Interferences -

$$x = \sum_{k} z_k \delta(u - r_k) \Rightarrow \rho\left(x \star \psi_{2^j,\ell}(u)\right) = \rho\left(\sum_{k} z_k \psi_{2^j,\ell}(u - r_k)\right)$$

$$\ell = 0 \qquad \ell = 1 \qquad \ell = 2 \qquad \ell = 3$$





ENS



$$j = 5$$



QM9: Data basis of 130.000 organic molecules with C, H, O, N, Fwith DFT atomisation energies

Regression error  $\sim 0.5$  kcal/mol  $\sim$  Deep Nets.

But small molecules with at most 29 atoms and 9 heavy ones



**Sparse Dictionary Representation** 

• Need to learn "sparse informative patterns"

Pattern representations with sparse dictionary expansions:



• How to minimise this convex cost ?



• Homotopy algorithms decrease the multiplier  $\alpha_k$ :

$$z = \arg\min_{\bar{z}} \|\bar{x} - D\bar{z}\|_{2}^{2} + \alpha_{k} \|z\|_{1}$$

with an iterated soft-threshold decreasing thresholds:

$$z_{k+1} = T_{\alpha_k} (D^t \bar{x} + (I - D^t D) z_k) \xrightarrow[k \to \infty]{\alpha_k \sim \gamma^{-k}} z_k$$

where  $T_{\alpha}(a) = \operatorname{sign}(a) \max(|a| - \alpha, 0)$  is a soft-thresholding.

 $\bullet$  Implemented with a convolutional network of depth K:



## **Dictionary Learning for Classification**

• Deep network with sparse coding and classification:



• Optimize the dictionary D and the classifier to minimize the classification loss over a supervised data basis  $\{x_i, y_i\}_i$ :

$$Loss(D) = \sum_{i} loss(y_i, \tilde{f}(x_i))$$

• Stochastic gradient descent



# Non Ergodic Processes autoencoder: trained on n examples $\{x_i\}_{i \le n}$ Encoder Gaussian white x $\mu$ $\mu$

Network trained on bedroom images:  $w = \alpha w_1 + (1 - \alpha)w_2$ 

 $w_1$ 

#### Network trained on faces of celebrities:



 $W_2$ 

## Generation as an Inverse Problem

#### Tomas Angles



#### **Inversion:**



U has a linear inverse  $U^{-1}$ :  $\rho(a) + \rho(-a) = a$ L is non-invertible linear projector

Regularization: inversion in a dictionary D where Ux is sparse Compute z such that Ux = Dz where z is sparse Non-linear multiscale model

## Generation as an Inverse Problem

#### **Inversion:**

Tomas Angles



U has a linear inverse  $U^{-1}$ , L is non-invertible linear projector

Inversion in a dictionary D where Ux is sparse:

$$Ux = Dz \Rightarrow w = LUx = LDz$$



• How to optimise the dictionary D ?

Learn the dictionary D by minimizing  $\sum_{i} ||x_{i} - \tilde{x}_{i}||^{2}$ with a stochastic gradient descent on a training set  $\{x_{i}\}_{i}$ 

#### **Training Reconstruction**

#### Celebrities Data Basis

ENS

Tomas Angles









 $x_i$ 

 $\tilde{x}_i$ 



## Image: Second second







Tomas Angles

Syntheses with different input noises



## Random Generations from Noise

Tomás Angles



#### Celebrities

ENS



#### Bedrooms





- Deep neural network are complex computational machines whose flexibility can be compared with Turing machines.
- A Relu on multiscale filters can produce scale interactions: creates phase harmonics, it may also be used to compute sparse representations, or piecewlinear approximations.
- One can define structured networks which are interpretable: similar to a structured program, with state of the art results.
- Still need functional analysis models and approximation theorems with decay rates.