

'incent Lostanlen, Stéphane Mallat Joan Bruna, Joaquim Andèn

École Normale Supérieure www.di.ens.fr/data High Dimensional Classification

• Audio signal  $x = (x(1), ..., x(d)) \in \mathbb{R}^d$ :



• Why is it so difficult ?



## **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  $\epsilon^{-d}$  points to cover  $[0,1]^d$  at a Euclidean distance  $\epsilon$  $\Rightarrow ||x - x_i||$  is always large



Wednesday, November 19, 14

#### **Euclidean Embedding**

Representation  $\Phi x \in \mathcal{H}$ Linear Classifier



Equivalent Euclidean metric:  $C_1 \|\Phi x - \Phi x'\| \le \Delta(x, x') \le C_2 \|\Phi x - \Phi x'\|$ 

How to define  $\Phi$  ?

Data:  $x \in \mathbb{R}^d$ 

||x - x'||: non-informative

## **Deep Convolution Neworks**

• The revival of an old (1950) idea: Y. LeCun



Optimize the  $L_k$  with architecture conditions: over 10<sup>9</sup> paramet Exceptional results for *images, speech, bio-data* classification. Products by FaceBook, IBM, Google, Microsoft, Yahoo...

Why does it work so well?

# **Overview of Questions**

- How to build audio signal representations for classification ?
- Why are deep neural networks so efficient?
- Why do wavelets appear in the cochlea and in most classifiers ?
- Why non-linearities ?

#### A Geometric Approach to Timbre

# **Geometric Representation**

- What geometry ?  $\longrightarrow_t$  quite poor...
- **Invariance** to translations  $x_c(t) = x(t-c)$

$$\forall c \in \mathbf{R} , \Phi(x_c) = \Phi(x) .$$

• Stability to deformations  $x_{\tau}(t) = x(t - \tau(t))$ small deformations of  $x \implies$  small modifications of  $\Phi(x)$ 

$$\forall \tau$$
 ,  $\|\Phi(x_{\tau}) - \Phi(x)\| \le C \sup_{t} |\tau'(t)| \|x\|$ .  
● Preserve information deformation size

# **Fourier Translation Invariance**

- Fourier transform  $\hat{x}(\omega) = \int x(t) e^{-i\omega t} dt$  invariance: if  $x_c(t) = x(t-c)$  then  $|\hat{x}_c(\omega)| = |\hat{x}(\omega)|$
- Instabilities to small deformations  $x_{\tau}(t) = x(t \tau(t))$ :  $||\hat{x}_{\tau}(\omega)| - |\hat{x}(\omega)||$  is big at high frequencies



Wednesday, November 19, 14

#### Wavelet Transform

• Dilated wavelets:  $\psi_{\lambda}(t) = 2^{-j/Q} \psi(2^{-j/Q}t)$  with  $\lambda = 2^{-j/Q}$ 





• Choice of Q: sparsity

Q-constant band-pass filters  $\hat{\psi}_{\lambda}$  $Q \approx 16$  for audio

• Wavelet transform: 
$$Wx(t) = \left\{ x \star \phi(t) , x \star \psi_{\lambda}(t) \right\}_{\lambda}$$

• If 
$$|\phi|^2 + \sum_{\lambda} |\hat{\psi}_{\lambda}|^2 = 1$$
 then  $||Wx||^2 = ||x||^2$ .

# Wavelet Translation Invariance



Modulus improves invariance:  $|x \star \psi_{\lambda_1}(x) \dagger \psi_{\lambda_1}(x) \dagger \psi_{\lambda_1}(x) \dagger \psi_{\lambda_1}(x) \dagger \psi_{\lambda_1}(x) \dagger \psi_{\lambda_1}(x) = 0$ 



Second wavelet transform modulus

$$|W_2| |x \star \psi_{\lambda_1}| = \left( \begin{array}{c} |x \star \psi_{\lambda_1}| \star \phi_{2J}(t) \\ |x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}(t)| \end{array} \right)_{\lambda_2}$$



**EXAMPLE 1**  
**Source of Contractions invariance and deformation stability:**  

$$i = x(u) = x(u - \tau(u)) \text{ then}$$

$$\lim_{J \to \infty} \|S_J D_\tau x - S_J x\| \le C \|\nabla \tau\|_{\infty} \|x\|$$

Wednesday, November 19, 14

## Audio Model

Excitation e(t) (pitched or random) Resonator filter h(t)

Amplitude modulation a(t)

$$x(t) = a(t) e \star h(t)$$

If the excitation is stationary then

$$\frac{||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}| \star \phi_{2^J}}{|x \star \psi_{\lambda_1}| \star \phi_{2^J}} \approx \frac{|a \star \psi_{\lambda_2}| \star \phi_{2^J}}{a \star \phi_{2^J}}$$

which characterises the amplitude spectrum.

#### **Amplitude Modulation**

#### Harmonic sound: $x(t) = a(t) e \star h(t)$ with varying a(t)



#### **Stochastic Excitation**

Random source:  $x(t) = a(t) e \star h(t)$  with varying a(t)

#### $|x \star \psi_{\lambda_1}|(t)$



#### Same Power Spectrum

Random source:  $x(t) = a(t) e \star h(t)$  with varying a(t) Spectrum



# Genre Classification (GTZAN)

- Musical genre classification (jazz, rock, classical, ...) 10 classes and 30 seconds tracks.
- Each frame is classified using a Gaussian kernel SVM.

| Feature Set                   | Error (%) |
|-------------------------------|-----------|
| Delta-MFCCs                   | 18.0      |
| Time Scat., order 1           | 19.1      |
| Time Scat., order 2           | 10.7      |
| Time Scat., order 3           | 10.6      |
| Time-frequency Scat., order 2 | 9.4       |

T = 740 ms

• Difficult to analyse the sources of inefficiencies

Time Invariant Scattering Moments

The scattering transform of a stationary process X(t)

$$S_{J}X = \begin{pmatrix} X \star \phi_{2J} \\ |X \star \psi_{\lambda_{1}}| \star \phi_{2J} \\ ||X \star \psi_{\lambda_{1}}| \star \psi_{\lambda_{2}}| \star \phi_{2J} \\ |||X \star \psi_{\lambda_{2}}| \star \psi_{\lambda_{2}}| \star \psi_{\lambda_{3}}| \star \phi_{2J} \\ \dots \end{pmatrix}_{\lambda_{1},\lambda_{2},\lambda_{3},\dots}$$

converges to moments if X is ergodic when  $2^J$  increases

$$\mathbb{E}(SX) = \begin{pmatrix} \mathbb{E}(X) \\ \mathbb{E}(|X \star \psi_{\lambda_1}|) \\ \mathbb{E}(||X \star \psi_{\lambda_1}| \star \psi_{\lambda_2}|) \\ \mathbb{E}(||X \star \psi_{\lambda_2}| \star \psi_{\lambda_2}| \star \psi_{\lambda_3}|) \\ \dots \end{pmatrix}_{\lambda_1, \lambda_2, \lambda_3, \dots}$$

• Does  $\mathbb{E}(SX)$  approximates well enough the distribution of X ?

# Inverse Scattering Transform

Joan Bruna

- Reconstruct x(t) from: N samples
  - 1st order coefficients:  $|x \star \psi_{\lambda_1}| \star \phi_{2^J}(t), \forall \lambda_1$

for 
$$2^J = \infty$$
:  $\int |x \star \psi_{\lambda_1}(u)| du$   
 $Q_1 \log_2 N$  coeffs

- 2nd order coefficients:  $||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}| \star \phi_{2^J}(t), \forall \lambda_1, \lambda_2$ 

for 
$$2^J = \infty$$
:  $\int ||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}(u)| du$   
 $Q_1 Q_2 \log_2^2 N/2$  coeffs

# **Stochastic of Audio Textures**

![](_page_19_Figure_1.jpeg)

![](_page_20_Figure_0.jpeg)

V.Lostanlen

Need to express frequency channel interactions: time-frequency image

![](_page_20_Figure_2.jpeg)

Wednesday, November 19, 14

Wavelets for Images

• Complex wavelet:  $\psi(t) = g(t) \exp i\xi t$ ,  $t = (t_1, t_2)$ rotated and dilated:  $\psi_{\lambda}(t) = 2^{-j} \psi(2^{-j}r_{\theta}t)$  with  $\lambda = (2^j, \theta)$ 

![](_page_21_Figure_2.jpeg)

• Wavelet transform:  $Wx = \begin{pmatrix} x \star \phi_{2^J}(t) \\ x \star \psi_{\lambda}(t) \end{pmatrix}_{\lambda \leq 2^J}$ 

## Scattering Transform in 2D

![](_page_22_Figure_1.jpeg)

ENS

# Ens Ergodic Texture Reconstructions

#### Joan Bruna

#### Original Textures

2D Turbulence

![](_page_23_Picture_4.jpeg)

![](_page_23_Picture_5.jpeg)

![](_page_23_Picture_6.jpeg)

![](_page_23_Picture_7.jpeg)

Gaussian process model with same second order moments

![](_page_23_Picture_10.jpeg)

![](_page_23_Picture_11.jpeg)

![](_page_23_Picture_12.jpeg)

![](_page_23_Picture_13.jpeg)

![](_page_23_Picture_14.jpeg)

For  $2^J = N$ :  $O(\log N^2)$  scattering moments:  $\|x \star \psi_{\lambda_1}\|_1 \approx \mathbb{E}(|x \star \psi_{\lambda_1}|)$ ,  $\||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}\|_1 \approx \mathbb{E}(||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}|)$ 

![](_page_23_Picture_16.jpeg)

![](_page_23_Picture_17.jpeg)

![](_page_23_Picture_18.jpeg)

![](_page_23_Picture_19.jpeg)

![](_page_23_Picture_20.jpeg)

#### **Cortical Transform**

K. Patil, D. Pressnitzer, S. Shamma, M. Elhilali

![](_page_24_Figure_2.jpeg)

# Timber of Musical Instruments

![](_page_25_Figure_1.jpeg)

• Good instrument classification performances with an SVM

#### Harmonic Spiral

• Problem: the wavelet transform of harmonics is not sparse

![](_page_26_Figure_2.jpeg)

• Alignment of harmonics in two main groups. More regular variations along  $(\theta, j)$  than  $\lambda$ 

#### **Spiral Scattering**

V.Lostanlen

![](_page_27_Picture_2.jpeg)

EN

#### Shepard-Risset Glissando

![](_page_28_Figure_1.jpeg)

3D separable Spiral wavelet transform  $W_2$ 

## Separate Pitch, Amplitude, Envelop-

V.Lost an len

• Second order spiral scattering in 5D:

$$Sx(t,\lambda_1,\lambda_t,\lambda_\theta,\lambda_j) = ||x \star \psi_{\lambda_1}| \star \psi_{\lambda_t} \psi_{\lambda_\theta} \psi_{\lambda_j}| \star \phi_{2^J}(t)$$

• Scale separation in  $x(t) = a(t) e \star h(t)$ 

 $\frac{\partial S(t,\lambda_1,\lambda_t,\lambda_\theta,\lambda_j)}{\partial \lambda_t} : \text{depends on the moduluation amplitude}$ 

$$\frac{\partial S(t,\lambda_1,\lambda_t,\lambda_\theta,\lambda_j)}{\partial \lambda_\theta} : \text{depends on the pitch variations}$$

$$\frac{\partial S(t,\lambda_1,\lambda_t,\lambda_\theta,\lambda_j)}{\partial \lambda_j}$$

: depends on the frequency envelop  $\hat{h}(\omega)$ 

# **Remarks and Questions**

- Difficult to study timber through classification experiments only Analysis-Synthesis is a powerful complement
- What relations between image and audio perception ?
  - Auditory scene analysis (Bregman)
  - Can we translate audio in images and reverse ?
- Can we analyse timbre from a geometrical point of view ?
  - Search for regularity transformed into sparsity with wavelets
- Are such tools enough to access to large time scale structures ?
- Do we need to learn Deep Network Filters ?