0000	000	

# Réseaux de Neurones Profonds, Apprentissage de Représentations

Thierry Artières

ECM, LIF-AMU

January 15, 2018







00	0000	00000
0000	000	

- Main architectures
  - Dense Architectures
  - Autoencoders
  - Convolutional NNs
- Learning
  - SGD
  - Optimization variants
  - Learning DNNs
- 3 Deep architectures
  - Very deep Models
  - What makes DNN work?

Main architectures		
	00000	00
0000	000	00000

# Outline

#### Main architectures

- Dense Architectures
- Autoencoders
- Convolutional NNs

#### Learning



Main architectures		
• 00 0000	00000 0000 000	00 00000
Dense Architectures		

## Dense architecture



Main architectures		
0 • 0	00000 0000	00 00000
Autoencoders	000	

# Autoencoders

#### Principal Component Analysis

- Unsupervised standard (Linear) Data Analysis technique
  - Visualization, dimension reduction
- Aims at finding principal axes of a dataset



- Reconstruct the input at the output via an intermediate (small) layer
- Unsupervised learning
- Non linear projection, distributed representation
- Hidden layer may be larger than input/output layers





Main architectures		
00		
0000	000	
Autoencoders		

## Deep autoencoders

## Deep NN with Diabolo shape

- Extension of autoencoders (figure [Hinton et al., Nature 2006])
- Pioneer work that started the Deep Learning wave



Main architectures		
0 00 ●000	00000 0000 000	00 00000
Convolutional NNs		

#### Motivation

- Exploit a structure in the data
  - Images : spatial structure
  - Texts, audio ; temporal structure
  - videos : spatio-temporal structure

#### Fully connected layers vs locally connected layers



Main architectures		
0 00 0●00	00000 0000 000	00 00000
Convolutional NNs		



Main architectures		
0 00 0●00	00000 0000 000	00 00000
Convolutional NNs		



Main architectures		
0	00000	00 00000
Convolutional NNs		



Main architectures		
0 00 0●00	00000 0000 000	00 00000
Convolutional NNs		



Main architectures		
0 00 0●00	00000 0000 000	00 00000
Convolutional NNs		



Main architectures		
	00000 0000	00 00000
0000	000	
Convolutional NNs		

#### Use of multiple maps



([LeCun and Ranzato Tutorial, DL, 2015])

#### Aggregation layers

- Subsampling layers with aggregation operator
- Max pooling  $\rightarrow$  brings robustness

Main architectures		
0 00 000●	00000 0000 000	00 00000
Convolutional NNs		

# Convolutional models

#### LeNet architecture [LeCun 1997]

• Most often a mix of (convolutional + pooling) layers followed by dense layers



	Learning	
	00000 0000	00 00000
0000	000	

# Outline



#### Learning

2

- SGD
- Optimization variants
- Learning DNNs



	Learning	
	00000	
0000	000	
SGD		

# Learning deep networks

#### Gradient descent optimization as for MLPs

- SGD
- With momentum
- Adagrad, Adam, Adadelta etc

	Learning	
	00000	
0000	000	
SGD		

# Gradient Descent Optimization

#### Gradient Descent Optimization

- Initialize Weights (Randomly)
- Iterate (till convergence)

• Restimate 
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \epsilon \frac{\partial C(\mathbf{w})}{\partial \mathbf{w}}|_{\mathbf{w}_t}$$



	Learning	
	00000	
0000	000	
SGD		

## Gradient Descent: Tuning the Learning rate

# Weight trajectory for two different gradient step settings.

#### Two classes Classification problem





Fig. 11. Weight trajectory and error curve during learning for (a)  $\eta=1.5$  and (b)  $\eta=2.5.$ 

## Images from [LeCun et al.]

	Learning	
	00000	
00	0000	00000
CCD		

## Gradient Descent: Tuning the Learning rate

#### Effect of learning rate setting

- Assuming the gradient direction is good, there is an optima value fir the learning rate
- Using a smaller value slows the convergence and may prevent from converging
- Using a bigger value makes convergence chaotic and may cause divergence



#### Images from [LeCun et al.]

T. Artières (ECM, LIF-AMU)

	Learning	
	00000	
0000	000	
SGD		

# Gradient Descent: Stochastic, Batch and mini batchs

Objective : Minimize  $C(\mathbf{w}) = \sum_{i=1..N} L_w(i)$  with  $L_w(i) = L_w(x^i, y^i, w)$ 

#### Batch vs Stochastic vs Minibatchs

- Batch gradient descent
  - Use  $\nabla C(\mathbf{w})$
  - Every iteration all samples are used to compute the gradient direction and amplitude
- Stochastic gradient
  - Use  $\nabla L_w(i)$
  - Every iteration one sample (randomly chosen) is used to compute the gradient direction and amplitude
  - Introduce randomization in the process.
  - Minimize C(w) by minimizing parts of it successively
  - Allows faster convergence, avoiding local minima etc
- Minibatch
  - Use  $\nabla \sum_{\text{few } j} L_w(j)$
  - Every iteration a batch of samples (randomly chosen) is used to compute the gradient direction and amplitude
  - Introduce randomization in the process.
  - Optimize the GPU computation ability

	Learning	
	00000	00
	0000	
0000	000	
Optimization variants		

## Using Momentum

#### SGD with Momentum

- Standard Stochastic Gradient descent :  $w = w \epsilon \frac{\partial C(w)}{\partial w}$
- SGD with Momentum:

$$v = \gamma v + \epsilon \frac{\partial C(w)}{\partial w}$$
$$w = w - v$$

	Learning	
	00000	00
00	0000	00000
Optimization variants		

# Guiding the learning

#### Regularization

• Constraints on weights (L1 or L2)



- Constraints on activities (of neurons in a hidden layer)
  - L1 or L2
  - Mean activity constraint (Sparse autoencoders, [Ng et al.])
  - Sparsity constraint (in a layer and/or in a batch)
  - Winner take all like strategies
- Disturb learning for avoiding learning by heart the training set
  - Noisy inputs (e.g. Denoising Autoencoder, link to L2 regularization)
  - Noisy labels

	Learning	
	00000	00
	0000	
0000	000	
Optimization variants		

# Guiding the learning

#### Denoising autoencoders and Deep Belief Networks

	N N N		2307447 2307447 2307447 2507447 2507447 2507447 2507447 2507447 2507447 2507447 2507447 2507447 2507447 2507447 2507447 2507447 250777 250777 250777777777777777777777	141	त्यस जाव जाव ज्यस जाव जाव
			和方法可是	222	24 24 24
	1.1	「「「「「「」」」	192002	つわた	C C C
	1 4 1 1		又當個同學生		12051 12051 245
			同業の理由が	スワン	20 25 20
김태양종	1		I H W Z W T	101 101 101	
	*	A State State State	5 51 6. 51 2-00	11/1	122 122 122

Figure 1: Activation maximization applied on MNIST. On the left side: visualization of 36 units from the first (1st column), second (2nd column) and third (3rd column) hidden layers of a DBN (top) and SDAE (bottom), using the technique of maximizing the activation of the hidden unit. On the right side: 4 examples of the solutions to the optimization problem for units in the 3rd layer of the SDAE, from 9 random initializations.

#### • Examples of learned filters with Denoising Autoencoders (top)

	Learning	
	00000	00
00	0000	00000
Ontimization variants		

# Guiding the learning

#### Dropout



- First method that allowed learning rellay deep networks without pretraining and smart initialization
- Related to ensemble of models
- Weights are normalized at inference time

	Learning	
	00000	00
00	0000	00000
	000	
Learning DNNs		

## Learning deep networks: Strategies

#### Few strategies (considering large volumes of unlabeled data)

- Very large labeled training dataset : Fully supervised setting
- Too few labeled training samples for supervised training : Unsupervised feature learning (each layer one after the other) + fine tuning with a classifier on top
- Very few labeled training samples : Unsupervised feature learning (each layer one after the other) + flat classifier learning

	Learning	
0000	000	
Learning DNNs		

## Learning deep networks

#### Unsupervised feature learning layer by layer



T. Artières (ECM, LIF-AMU)

	Learning	
	00000	00
00	0000	00000
	000	
Learning DNNs		

## Learning more general architectures





#### Still optimized with Gradient Descent !!

$$W = W - \epsilon \frac{\partial C(W)}{\partial W}$$

- provided functions implemented by blocks are differentiable
- and derivatives  $\frac{\partial Out(B)}{\partial ln(B)}$  and  $\frac{\partial Out(B)}{\partial W(B)}$  are available for every block



		Deep architectures
0 00	00000	00
0000	000	

# Outline



#### Oeep architectures

- Very deep Models
- What makes DNN work?

		Deep architectures
		0
0000	000	
Very deep Models		

## The Times They Are A Changing



## (slide from [Kaiming He])

T. Artières (ECM, LIF-AMU)	Deep Learning	January 15, 2018 25 / 31

		Deep architectures
		00
00	0000	00000
0000	000	
Very deep Models		

## From shalow to deep



		Deep architectures
		00000
0000	000	
What makes DNN work?		

# Deep vs Shalow ?

Characterizing the complexity of functions a DNN may implement [Pascanu and al., 2014]

- $\bullet\,$  DNNs with RELU activation function  $\Rightarrow$  piecewise linear function
- Complexity of DNN function as the Number of linear regions on the input data
- Case of  $n_0$  inputs and  $n = 2n_0$  hidden cells per HL (k HL) :
  - Maximum number of regions :  $2^{(k-1)n_0} \sum_{i=0}^{n_0} {2n_0 \choose i}$
- Example:  $n_0 = 2$ 
  - Shallow model:  $4n_0$  units  $\rightarrow$  37 regions
  - Deep model with 2 hidden layers with  $2n_0$  units each  $\rightarrow$  44 regions
  - Shallow model:  $6n_0$  units  $\rightarrow$  79 regions
  - Deep model with 3 hidden layers with  $2n_0$  units each ightarrow 176 regions
- Exponentially more regions per parameter in terms of number of HL
- At least order (k-2) polynomially more regions per parameter in terms of width of HL n

		Deep architectures
		00000
0000	000	
What makes DNN work?		

## Deep vs Shalow ?



#### From [Pascanu and al., 2014]

- Left: Regions computed by a layer with 8 RELU hidden neurons on the input space of two dimensions (i.e. the output of previous layer)
- Middle: Heat map of a function computed by a rectifier network with 2 inputs, 2 hidden layers of width 4, and one linear output unit. Black lines delimit regions of linearity of the function
- Right: Heat map of a function computed by a 4 layer model with a total of 24 hidden units. It takes at least 137 hidden units on a shallow model to represent the same function.

		Deep architectures
00	0000	00000
0000	000	
What makes DNN work?		

## The depth alone is not enough

#### Making gradient flow for learning deep models

- Main mechanism : Include the identity mapping as a possible path from the input to the output of a layers
- ResNet building block [He and al., 2015]]



• LSTM (deep in time) [Hochreichter and al., 1998]



		Deep architectures
00	0000	00000
0000	000	

What makes DNN work?

## About generalization, overtraining, local minimas etc

## Traditional Machine Learning

- Overfiting is the enemy
- One may control generalization with appropriate regularization

#### Recent results in DL

- The Overfit idea should be revised for DL [Zhand and al., 2017]
  - Deep NN may learn noise !
  - Regularization may slightly improve performance but is not THE answer for improving generalization
- Objective function do not exhibit lots of saddle points and most local minima are good and close to globale minimas [Choromanska et al., 2015]
  - Not clear what in the DNN may allow to predict its generalization ability

		Deep architectures
00	0000	00000
0000	000	
What makes DNN work?		

## Favorable context

#### Huge training resources for huge models

- Huge volumes of training data
- Huge computing ressources (clusters of GPUs)

#### Advances in understanding optimizing NNs

- Regularization (Dropout...)
- Making gradient flow (ResNets, LSTM, ...)

#### Faster diffusion than ever

- Softwares
  - Tensorflow, Theano, Torch, Keras, Lasagne, ...
- Results
  - Publications (arxiv publication model) + codes
  - Architectures, weights (3 python lines for loading a state of the art computer vision model!)