# Introduction au Deep Learning

Thierry Artières

Equipe QARMA d'apprentissage automatique Pole Science des Données - LIS lab Ecole Centrale Marseille

July 4, 2019







000000

# Outline

- Learning DNNs
  - Learning is not so easy!
  - SGD for DNNs
  - Architecture design
  - Adversarial

Learning DNNs

Output

Adversarial O

0000000

# From shalow to deep



Learning DNNs
000000000
00000000
SGD for DNNs



### While NN optimization is reasonably easy

# (Easier) analysis for ReLU DNNs [LeCun Statlearn tutorial]

- ... but expectation of generalized results to other activation functions
- With ReLu and MaxPooling operators one may formalize what happens on a path from an input to the output
- The output may be computed as :

$$\hat{y} = \sum_{P} \delta_{P}(W, X) (\prod_{(ij) \in P} w_{ij}) x_{j_{start}}$$

- $\delta_P(W,X)$  : 1 if active path, 0 otherwise
- Implemented function is piece-wise linear



Learning DNNs
000000000
00000000
SGD for DNNs



#### While NN optimization is reasonably easy

# (Easier) analysis for ReLU DNNs [LeCun Statlearn tutorial]

• Objective function : piece wise polynomial (degree = number of hidden layers) with partially random coefficients

$$C(W) = \sum_{P} C_{P}(X, Y, W)(\prod_{(ij) \in P} w_{ij})$$

• Hint from results on distribution of critical points for polynomials with random coefficients



Learning DNNs	
0 0 <b>00000000</b> 00000000	0 0000000 000000
SGD for DNNs	

### Deep ReLu networks

#### Easier analysis... [LeCun Statlearn tutorial]

- Experiments by [Choromanska and al., 2015]: Train 2-layers nets on Mnist from multiple initializations and measure loss on the test set
- Many close local minimas for large nets
- Objective function do not exhibit lots of saddle points and most local minima are good and close to globale minimas



[Choromanska, Henaff, Mathieu, Ben Arous, LeCun 2015]

Learning DNNs		
0 <b>000000000</b> 000000000		
SGD for DNNs		

#### Adversarial 0 00000000 000000

#### Yet the depth alone is not enough

#### While SGD works well for NNs, the optimization of DNNs requires careful design and tricks

- Make the gradient flow with activation normalization (Batch Normalization)
- Make the gradient flow with structural constraints (Identity mapping)
- Regularization (Dropout)



# Problems with activity propagation in deep NNs [He et al., 2016]

Few slides from Fei Fei Li

#### Standard initialization schema for MLPs

- 10 layers networks (500 neurones each, with tanh)
- Initialization : gaussian random with small (std=0.01) values (what if all null initialization?)
- All activations at 0
- What about the gradient ?



Learning DNNs
0000000000
00000000
SGD for DNNs

#### **Batch Normalization**

#### Main idea

- Usually inputs to neural networks are normalized to either the range of [0, 1] or [-1, 1] or to mean=0 and variance=1
- BN essentially performs Whitening to the intermediate layers of the networks.
- Usually placed before nonlinearities



Learning DNNs	
0 0000000000 00000000	0 00000000 000000
SGD for DNNs	

# Identity mapping in Residual Networks

#### Principle

- Include identity mapping in the model
- ResNet building block [He and al., 2015]]
- Every layer becomes close to the output ( $\Rightarrow$  not far in the backpropagation process)



Learning DNNs	
o ooooo•oooo ooooooo	0 00000000 000000
SGD for DNNs	

# Identity mapping in Residual Networks

#### Principle

- Include identity mapping in the model
- ResNet building block [He and al., 2015]]
- Every layer becomes close to the output ( $\Rightarrow$  not far in the backpropagation process)



Learning DNNs
00000000000
00000000
SGD for DNNs



# Identity mapping with LSTM units in RNNs

#### What it does

- Main behaviour
  - If  $f_t == 1$  and  $i_t == 0$  use previous cell state
  - If  $f_t == 0$  and  $i_t == 1$  ignore previous cell sate
  - If  $o_t == 1$  output ois set to cell sate
  - If  $o_t == 0$  output is set to 0

#### With...

- Cell state ct
- Forget gate  $f_t$
- Input gate it
- Output o<sub>t</sub>
- Hidden state to propagate to upper layers *h*<sub>t</sub>



Learning DNNs	
o ⊙⊙⊙⊙⊙⊙⊙⊙⊙ ⊙○○○○○○○	0 00000000 000000
SGD for DNNs	

#### Auxiliary loss on intermediate layers

#### Google net

• Auxiliary loss brings some gradient to first layers



Learning DNNs	
o oooooooooo oooooooo	0 00000000 000000
SGD for DNNs	

# Regularizing with Dropout [Hinton 2012]

#### Principle



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

Learning DNNs	
0 000000000● 00000000	o 00000000 000000
SGD for DNNs	

# Dropout [Hinton 2012]

Do not ever fear overfitting !



#### Examples of architectures

#### AlexNet [Krizhevsky and al., 2012] (top) and NetworkInNetwork [Lin and al., 2013] (bottom)





Learning DNNs
0000000
Architecture design



# Looking for a good architecture: Lego game

How to reach such an architecture (GoogleNet 2014) ?





Searching for a good architecture requires making choices !!

# Looking for a good architecture: Lego game

#### Gridsearch

- Standard Machine Learning models
  - Very few hyperparameters (regularization tradeoff, kernel width or degree etc)
  - Easier optimization problem
  - Usually much less data and much simpler models
  - $\bullet \ \Rightarrow \mathsf{Quite \ exhaustive \ gridsearch}$
- Large deep networks
  - Many choices (sequence of layers, width of layers, convolution kernel's size and strides, activation function, optimization routine and its parameters...): Not many theoretical hints
  - Harder optimization problem
  - Each try is expensive
  - $\Rightarrow$  Reuse of others' architectures whenever possible
  - $\bullet \ \Rightarrow \ \mathsf{Gain} \ \mathsf{experience} \ \mathsf{on} \ \mathbf{how} \ \mathbf{to} \ \mathsf{design}$

# Looking for a good architceture



# Looking for a good architceture

#### Illustration [Verbeek 2017]

- Simple search (but for a large network)
  - 19 convolution layers and 5 pooling layers to set
  - Question: where to put the poooling layers?  $\rightarrow$  40 000 architectures !!
  - No question about layers' dimensions, activation function, kernels' size, pooling type etc

#### Remember

- 1 hour GPU on AWS = 1  $\$
- Learning 1 model = Few hours
   ⇒ Expensive design !!!
- Not much alternatives





#### Looking for a good architecture: use others' !!

Deep Models for High resolution images [Radford 2015]

Historical attempts to scale up GANs using CNNs to model images have been unsuccessful. This motivated the authors of LAPGAN (Denton et al.) [2015) to develop an alternative approach to iteratively upscale low resolution generated images which can be modeled more reliably. We also encountered difficulties attempting to scale GANs using CNN architectures commonly used in the supervised literature. However, after extensive model exploration we identified a family of architectures that resulted in stable training across a range of datasets and allowed for training higher resolution and deeper generative models.



Learning DNNs
0000000000
00000000
Architecture design



# Architecture design from prior knowledge (HEP example)

#### What if unstructured data: is DNN useful ?

- Modeling collision at CERN : features = physical features of jets during the collision (speed, energy, angle with collision axis...)
- Main approaches
  - Use non deep machine learning models
  - Represent data as images and use Deep NNs
  - Design DNN architecture from knowledge on the considered process
- Example : Learn to aggregate features of jets using a tree structure inspired from data knowledge [Louppe et al., 2018]





# Architecture design from prior knowledge (HEP example)

#### Ziyu Guo's thesis (with Y. Coadou at CPPM)

- Deep learning in the search for ttH with the ATLAS experiment at the Large Hadron Collider
- Rely on the physical process to design the NN structure
- Better results than DNN, on par with state of the art models in HEP (BDTs)



Adversarial

# Outline



#### Learning DNN

#### 2 Adversarial

- Generative models
- GANs
- Examples

00000000	
Generative models	



# Generative models

#### Goal

- · Learn to generate complex and realistic data
- Statistical viewpoint : learn a model of the density of data / able to sample with this density
  - Postulate a parametric model : Usually not complex enough
  - Postulate a parametric form and perform optimization (e.g. Maximum Likelihood) : Intractable for complex forms  $p(x) = \frac{F(x)}{Z(x)}$  with  $Z(x) = \sum_{x} F(x)$



Maximum Likelihood Estimation (MLE)

GANs

# Adversarial Learning: Generator

#### Determinitic NN as a generative model



Using a deterministic NN as a generative model

- Let note the function implemented by the model as G
- Let note the input  $z \to$  The NN computes G(z)
- Assume z obeys a prior (noise) distribution,  $p_z$ , e.g. Gaussian distribution
- then the output x of the NN follows a distribution

$$\Rightarrow p_G(x) = \int_{z \text{ s.t. } G(z)=x} p_z(z) dz$$

00000	000

#### GANs



# Le principe de l'adversarial learning [Goodfellow and al., 2014]

#### Principle

- Jeu à deux joueurs: un générateur et un discriminateur
  - le discriminateur veut distinguer les exemples générés des vrais exemples
  - Le générateur veut tromper le discriminateur





### Adversarial Learning criterion

Criterion from [Goodfellow and al., 2014]

- Generator G and Discriminator D are two NNs
  - Whose parameters are noted  $\theta_g$  and  $\theta_d$
- Distributions
  - p<sub>data</sub> stands for the empirical distribution of the data from the training set
  - $p_z$  is a prior noise distribution, e.g. a Gaussian distribution
  - On convergence we want  $p_g = p_{data}$
- Learning criterion:

$$min_g max_d v(\theta_g, \theta_d) = \mathbf{E}_{x \sim \rho_{data}} \left[ log D(x) \right] + \mathbf{E}_{z \sim \rho_z} \left[ log (1 - D(G(z))) \right]$$

- Assume G is fixed: D is trained to distinguish between fake and true samples
- Assume D is fixed : G is trained to generate samples as realistic as possible

	Adversarial
00000000	0000000
0000000	000000
CANE	

# Adversarial Learning theory: What happens during Learning



00000000

Adversarial

#### GANs

#### Good Examples



0000000

GANs

Adversarial 0 00000000 000000

# Bad examples







000000	

#### GANs

#### About GANs'

#### Known problems

- DIfficult learning
- Very long learning
- Missing modes
- Evaluation measures

#### Many many variants

- Conditional
- Disantangling
- Image editing



Adversarial O OOOOOOOOO Learning DNNs 0 0000000000 00000000

GANs

Adversarial 0 0000000

# DCGANS [Radford 2015]





T. Artières (ECM - LIS / AMU )







# About using additional discriminators [Ganin et al, ICML 2015]

#### Other use of adversarial constraints

- Adversarial constraints may be used for structuring what we learn in hidden layers
- For instance : favoring some representation do not include a specific information





# Conditional GANs [Mirza and al., 2014]

#### Learning criterion

Citerion

$$min_g max_d v(\theta_g, \theta_d) = \mathbf{E}_{x, y \ P_{data}} \left[ log D(x, y) \right] + \mathbf{E}_{z \ P_z, y' \ P_y} \left[ log (1 - D(G(z, y'), y')) \right]$$



	Adversarial
0 000000000 00000000	0 00000000 <b>00000</b> 0
Examples	

# Image editing with Invertible Conditional GANs [Perarnau and al., 2016]



Learning DNNs	Adversarial
000000000	00000
	000000
Examples	

# Disentangling factors of variation [Chen et al., 2018]

Generating images under various styles





0 000000000 0000000
Examples



# Disentangling factors of variation [Chen et al., 2018]

Transfering styles between images



	Adversarial
0 000000000 00000000	0 00000000 0000000
Examples	

## Transfering styles between motion capture sequences [Qi et al., 2017]



00000000
Examples

Adversarial 0 0000000 000000

#### Motion capture data

