

Reconnaissance des formes

Deep Learning

Alexandre BOULCH

ONERA

THE FRENCH AEROSPACE LAB

Plan de la séance

From 90's to Deep learning

Deep learning

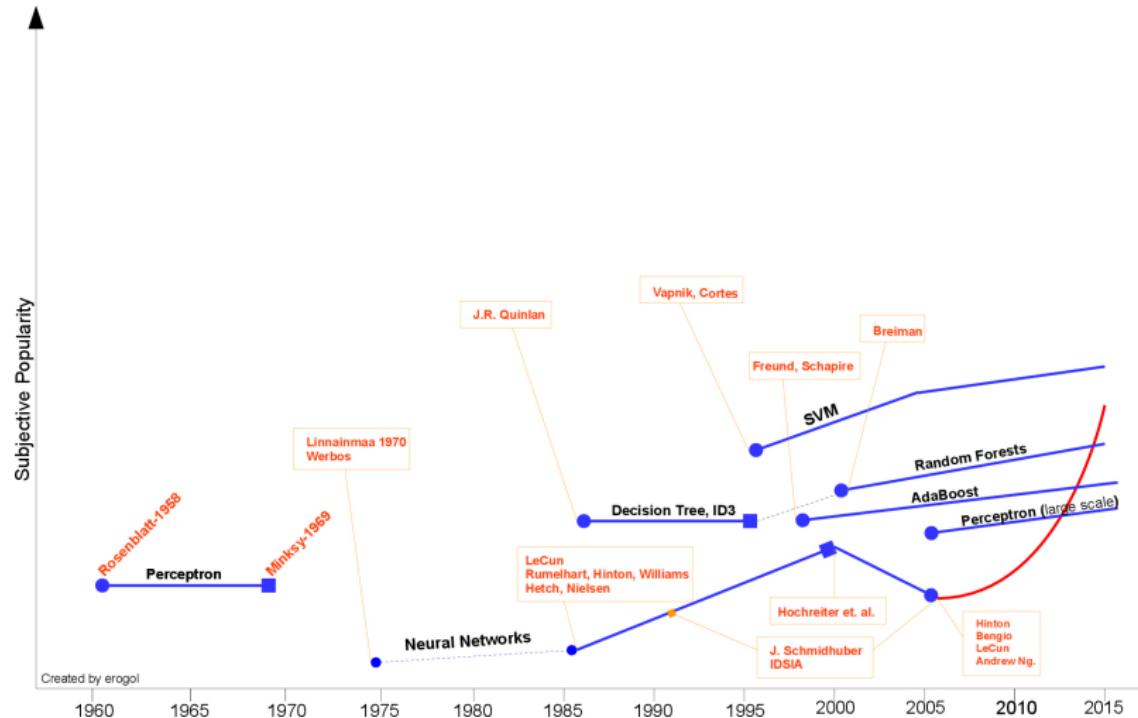
2000's : better, faster, stronger

Do not forget the classics

What can we do with neural networks ?

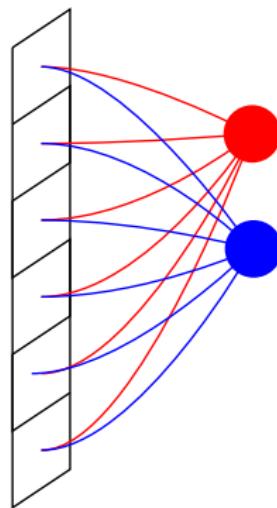
Frameworks

Introduction



Fully connected

- ▶ Perceptron
- ▶ A neuron is connected to every input.
- ▶ High number of parameters
 $(|inputs| * |outputs|)$



Plan de la séance

From 90's to Deep learning

Deep learning

Concept

Convolutional networks

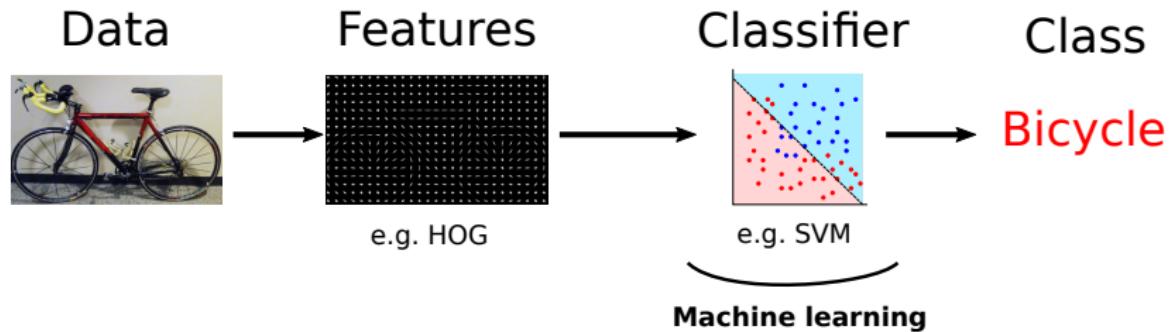
2000's : better, faster, stronger

Do not forget the classics

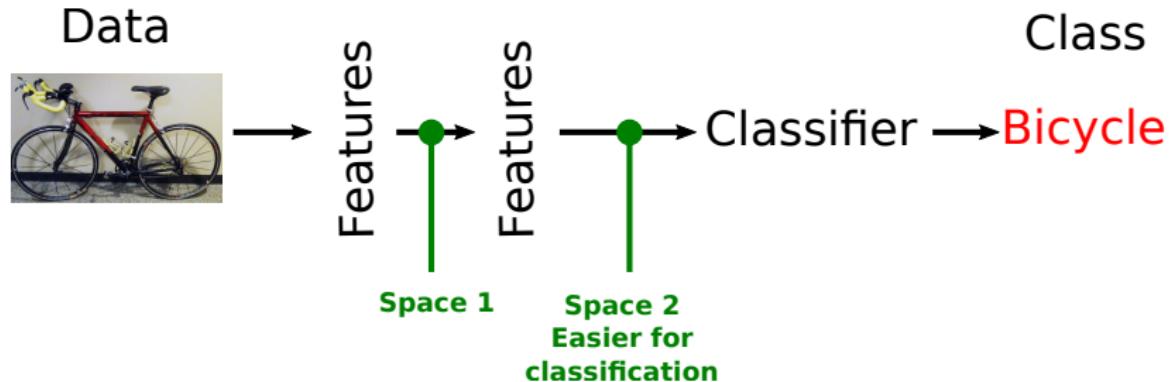
What can we do with neural networks ?

Frameworks

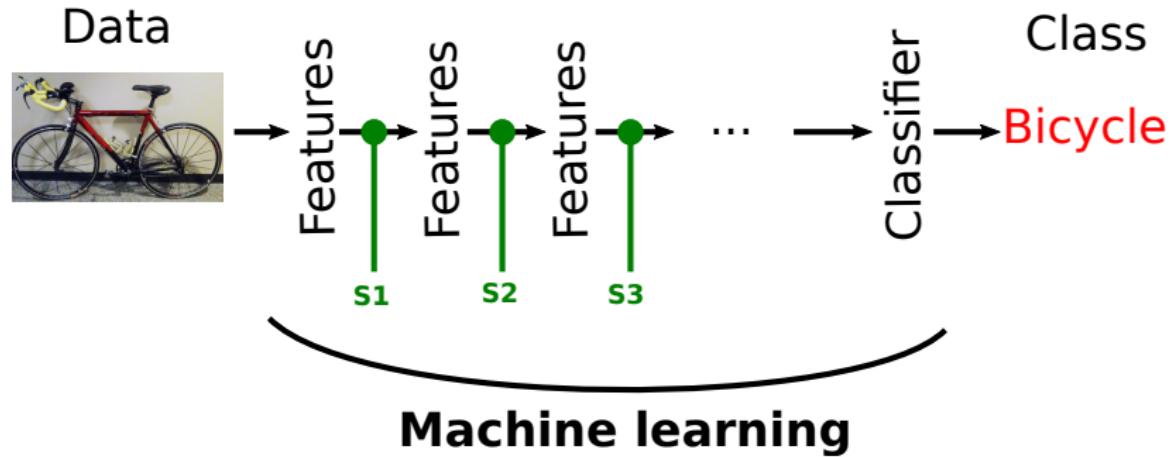
Deep learning concept



Deep learning concept



Deep learning concept



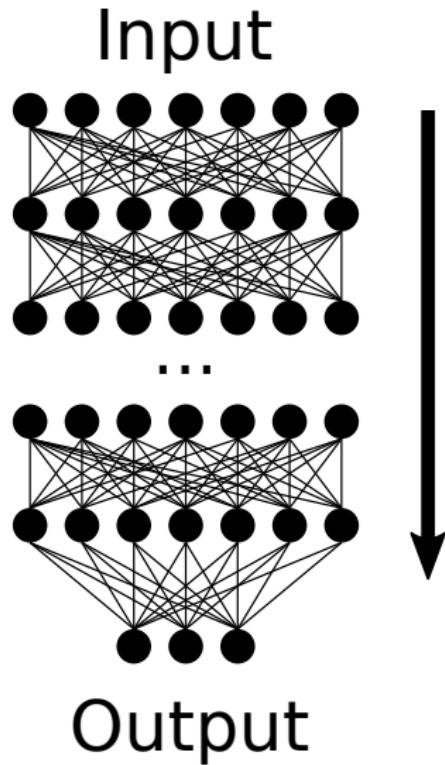
Concept

Deep learning is about learning representations.

Multi-layer perceptron

MLP de plus en plus profonds (avant 2005)

- ▶ Très grosses difficultés d'optimisation
- ▶ Convergence difficile
- ▶ Peu de données
- ▶ Entraînement très long (Abandon progressif au profit des SVMs)
 - ▶ Simples à utiliser
 - ▶ Preuves de convergence
 - ▶ Rapides



First convolutional networks

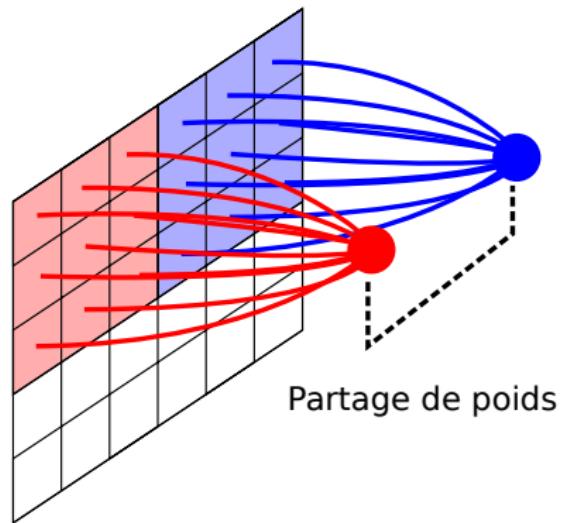
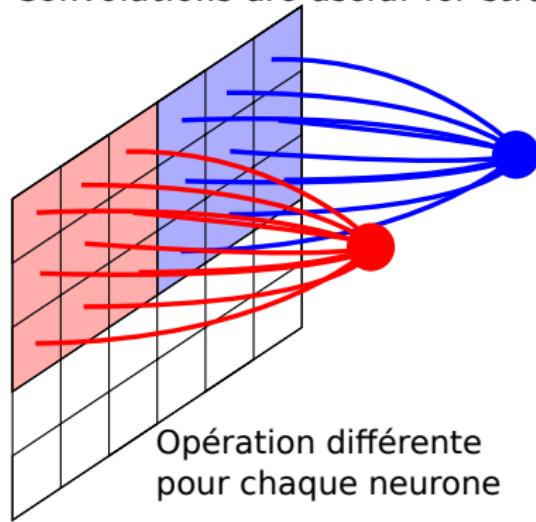
Idea

What is important with images ?

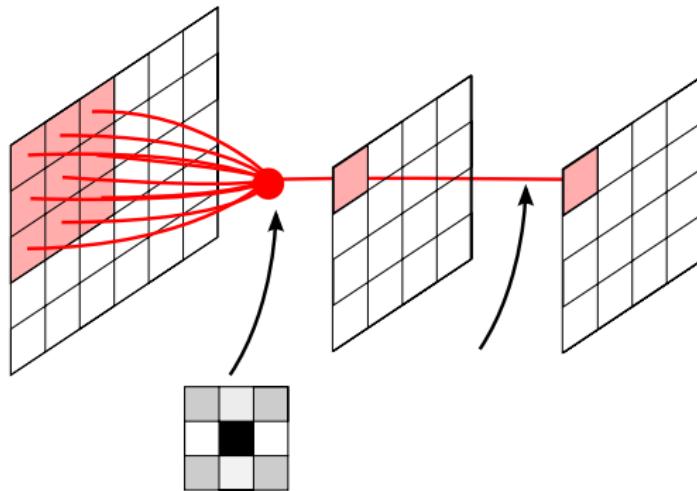
1. Neighborhood (objects have limited sizes)
2. Interest part can be everywhere (same operations everywhere)

Convolutions

Convolutions are useful for structured data.



Convolutions



Forward

$$x_{i,j}^l = \sum_a \sum_b \omega_{a,b} y_{i+a,j+b}^l \quad (1)$$

$$y_{i,j}^{l+1} = \sigma x_{i,j}^l \quad (2)$$

Convolution

Backward weight update

$$\frac{\partial E}{\partial \omega_{a,b}} = \sum_i \sum_j \frac{\partial E}{\partial x_{i,j}^l} \frac{x_{i,j}^l}{\omega_{a,b}} = \sum_i \sum_j \frac{\partial E}{\partial x_{i,j}^l} y_{i,j}^{l-1} \quad (3)$$

And

$$\frac{\partial E}{\partial x_{i,j}^l} = \frac{\partial E}{\partial y_{i,j}^l} \frac{\partial y_{i,j}^l}{\partial x_{i,j}^l} = \frac{\partial E}{\partial y_{i,j}^l} \sigma'(x_{i,j}^l) \quad (4)$$

Finally

$$\frac{\partial E}{\partial \omega_{a,b}} = \sum_i \sum_j \frac{\partial E}{\partial y_{i,j}^l} \sigma'(x_{i,j}^l) y_{i,j}^{l-1} \quad (5)$$

Convolution

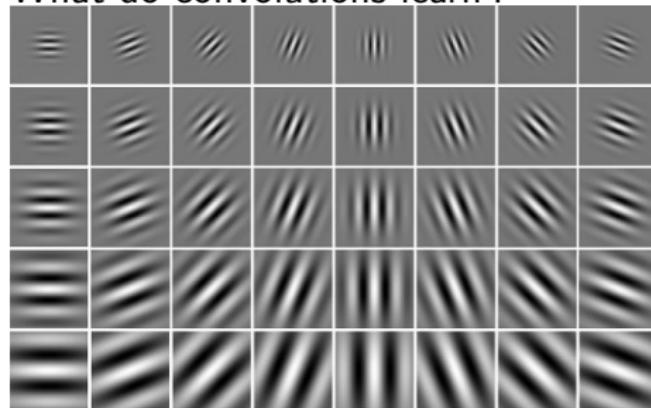
Backward gradient propagation

$$\frac{\partial E}{\partial y_{i,j}^{l-1}} = \sum_a \sum_b \frac{\partial E}{\partial x_{i-a,j-b}^l} \frac{\partial x_{i-a,j-b}^l}{y_{i,j}^{l-1}} \quad (6)$$

$$= \sum_a \sum_b \frac{\partial E}{\partial x_{i-a,j-b}^l} \omega_{a,b} \quad (7)$$

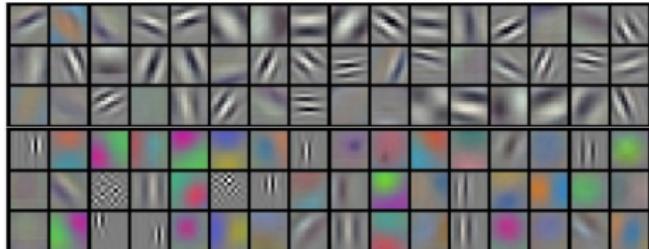
Convolution

What do convolutions learn ?



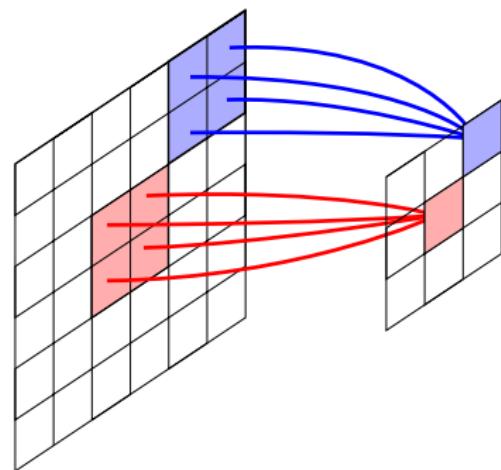
Gabor filters.

First layer of AlexNet.

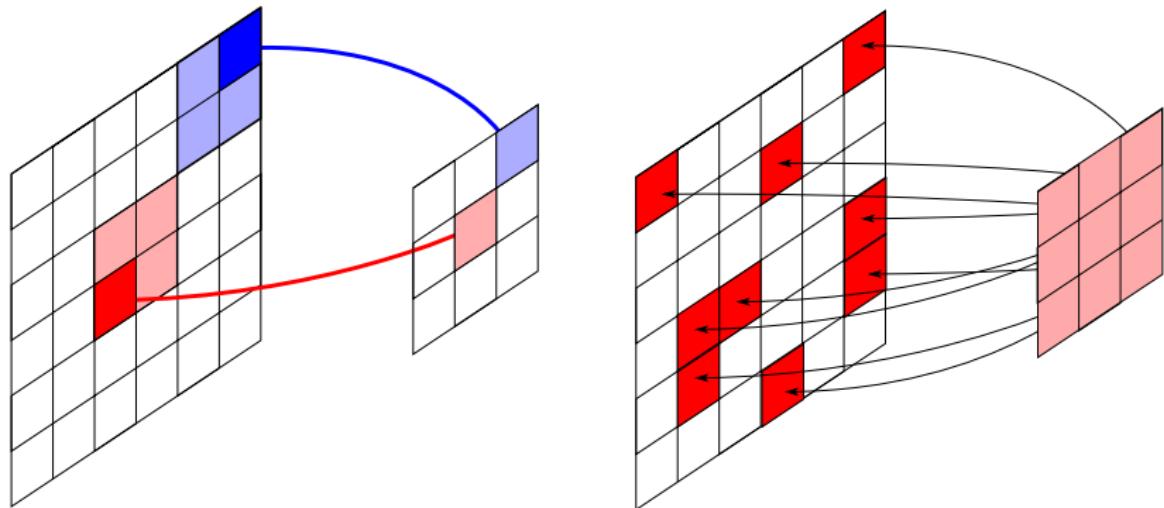


Max Pooling

- ▶ dimension reduction
- ▶ translation invariability



Max Pooling

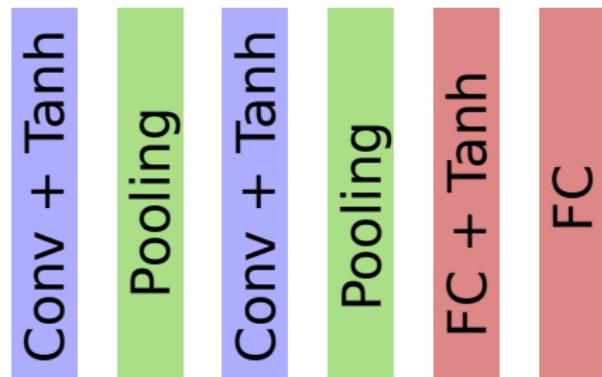


Forward : max signal

Backward : Gradient transmission to max signal origin, zero otherwise

LeNet (1990)

LeNet (1990)
Images 28x28



Issues

- ▶ Learning speed
- ▶ Exploding or vanishing gradients
- ▶ Overfitting
- ▶ Local minima

Limitations

- ▶ Architecture
- ▶ Initialization
- ▶ Computing power
- ▶ Data
- ▶ Optimization

Plan de la séance

From 90's to Deep learning

Deep learning

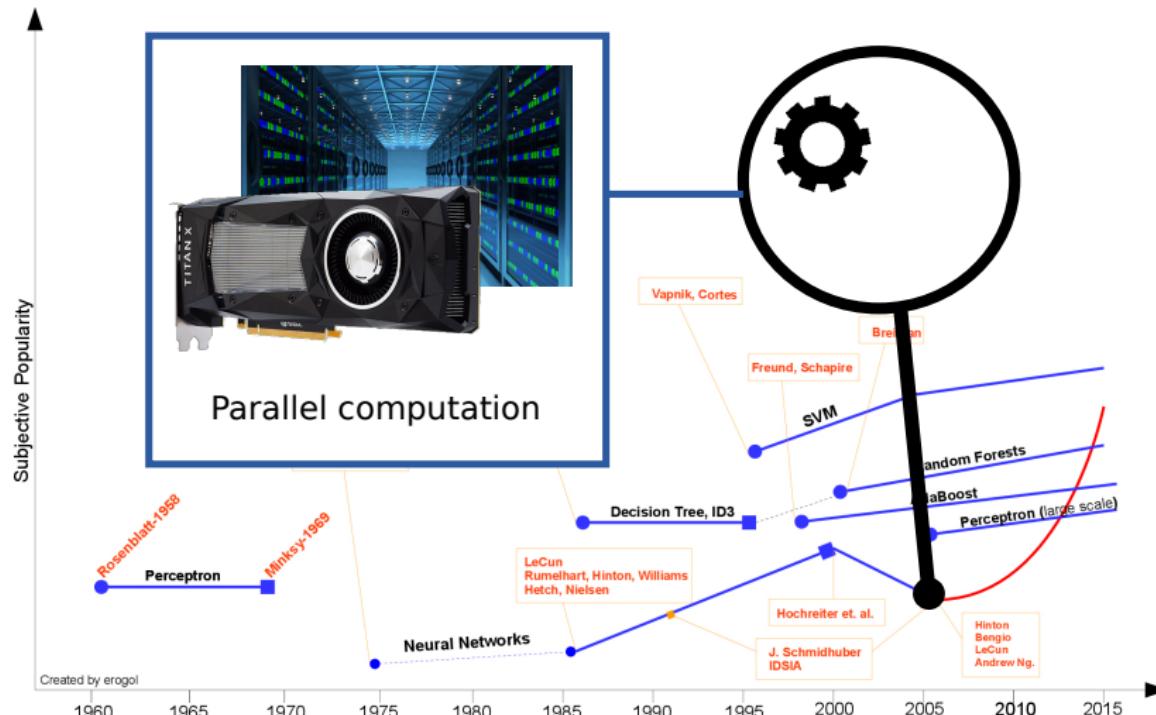
2000's : better, faster, stronger

Do not forget the classics

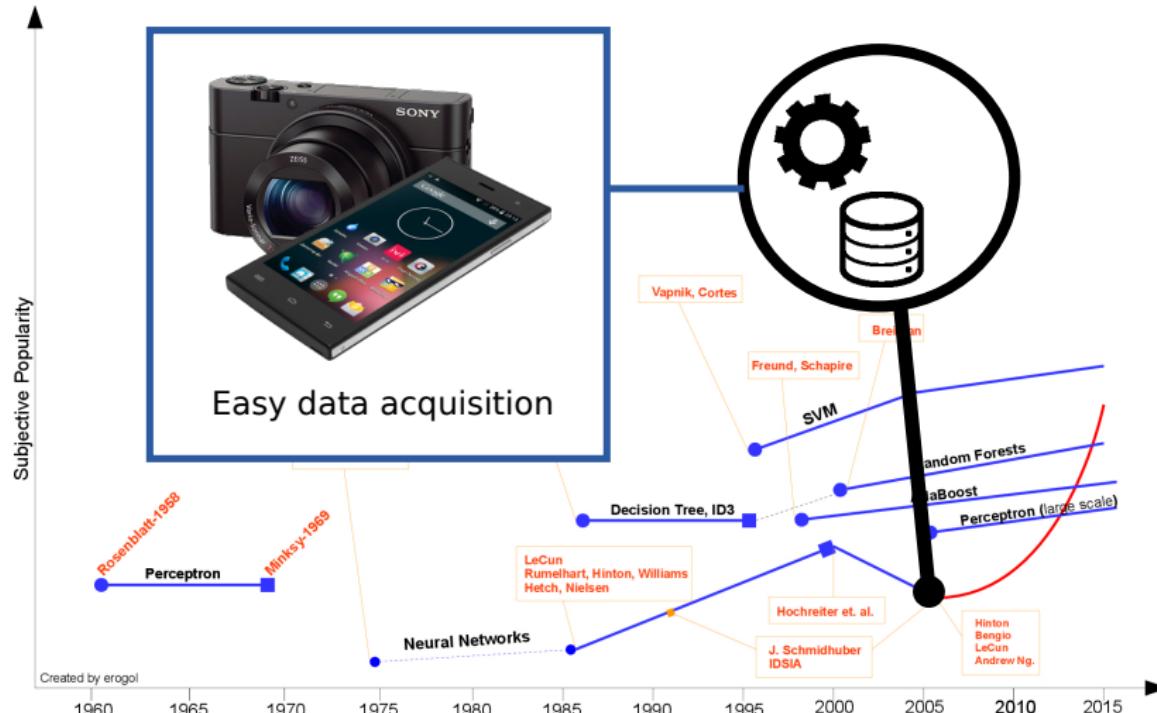
What can we do with neural networks ?

Frameworks

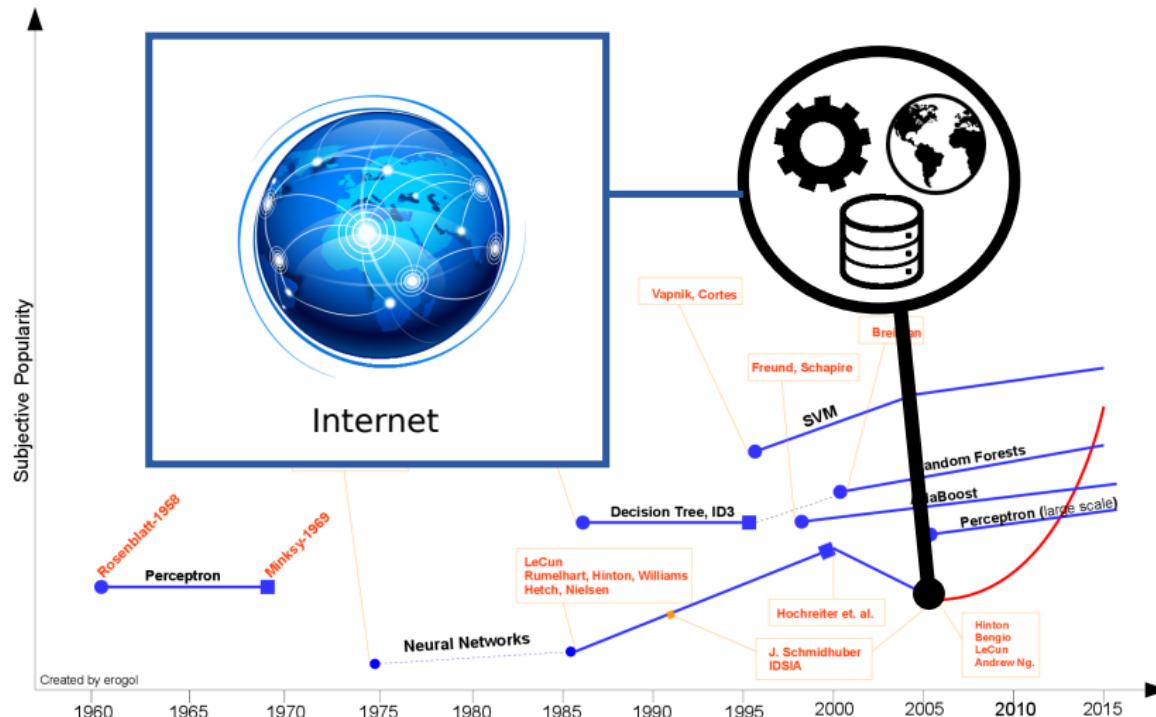
Massively data driven approaches



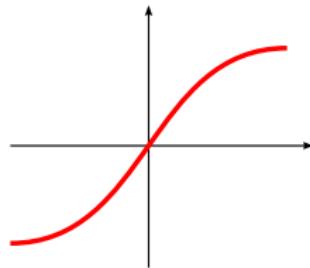
Massively data driven approaches



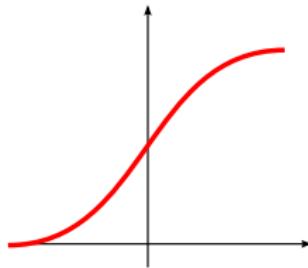
Massively data driven approaches



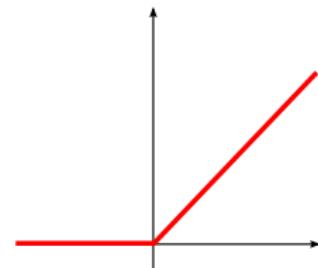
Activations



Tangente
hypébolique



Sigmoïde

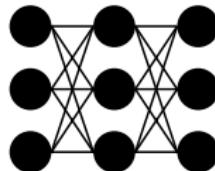


Rectified Linear
Unit (ReLU)

Rectified linear unit

1. Faster gradient computation
2. Similar convergence

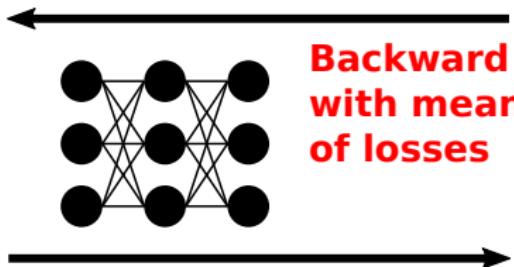
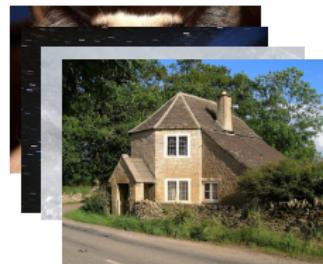
Mini-batches



Loss image 1



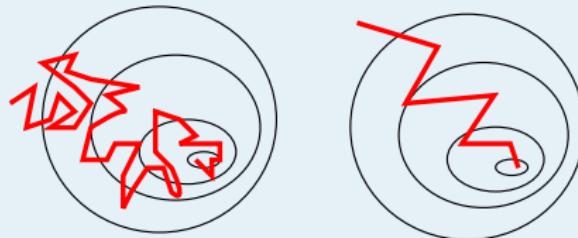
Mini-batches



**Backward
with mean
of losses**

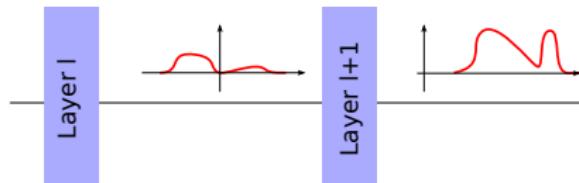
Loss image 1
Loss image 2
Loss image 3
Loss image 4
...

Gradient smoothing



Smoother gradient converges faster.

BatchNorm



Changes in the signal dynamic make the model more difficult to optimize : exponential or vanishing gradients.

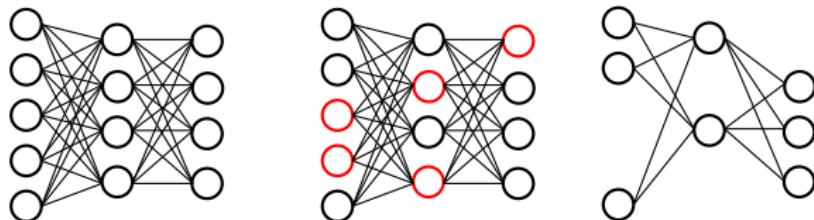
Objective : control the signal distribution :

$$y^{I*} = \frac{y^I - \mu}{\sqrt{\sigma^2 + \epsilon}} \gamma + \beta \quad (8)$$

γ and β are learnt, μ and σ are computed (mean and standard deviation).

Learning is faster (iteration number) but slower (statistics computation).

Dropout



Training

- ▶ Shutdown neurons with probability $1 - p$
- ▶ Equivalent to training multiple sub-networks

Inference

- ▶ Use all neurons, with p ponderation

Dropout give robustness to overfitting.

Weight initialization

Weights have great influence on convergence speed. They are randomly initialized.

- ▶ too small weights : vanishing signal
- ▶ too high : exploding signal

Objective

Conservation of signal properties.

$$\text{Var}(Y) = \text{Var}(X)$$

Weight initialization

Xavier initialization

$X \in \mathbb{R}^n$, weights W and output $Y \in \mathbb{R}$

$$Y = W_1 X_1 + W_2 X_2 + \cdots + W_n X_n$$

X_i and W_i independent :

$$\text{Var}(W_i X_i) = E[X_i]^2 \text{Var}(W_i) + \text{Var}(X_i) \text{Var}(W_i) + \text{Var}(X_i) E[W_i]^2$$

$E[X_i] = 0$ and $E[W_i] = 0$:

$$\text{Var}(W_i X_i) = \text{Var}(X_i) \text{Var}(W_i)$$

finally $\text{Var}(Y) = n \text{Var}(X_i) \text{Var}(W_i)$ and we chose

$$\text{Var}(W_i) = \frac{1}{n}$$

Optimization - Stochastic Gradient Descent

See neural network class

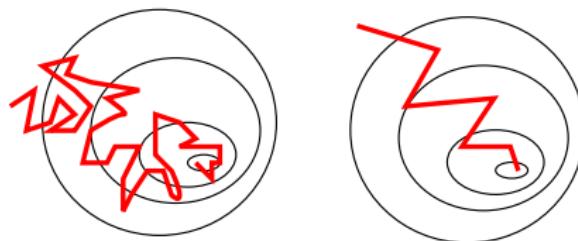
Update rule : $w_{t+1} = w_t + \alpha \Delta w$ (learning rate α)

Learning rate variation

- ▶ Step decrease
- ▶ Exponential decrease
- ▶ ...

Optimization - SGD with Momentum

Same idea as mini batch : smooth gradient in the good direction



Momentum

Use previous gradient to ponderate the direction of the new gradient.

$$v_t = \gamma v_{t-1} + \alpha \Delta w$$

$$w_t = w_{t-1} - v_t$$

γ is the momentum.

Optimization - Adaptative methods

- ▶ **Adagrad**

$$w_{t+1,i} = w_{t,i} - \frac{\alpha}{\sqrt{G_{t,i} + \epsilon}} \Delta w_i$$

with $G_{t,i}$ sum of gradients of w_i

- ▶ **Adam (adaptive momentum estimation)**

$$m_t = \frac{\beta_1}{1 - \beta_1} m_{t-1} + \Delta w_t$$

$$v_t = \frac{\beta_2}{1 - \beta_2} v_{t-1} + \Delta w_t^2$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{v_t + \epsilon}} m_t$$

Plan de la séance

From 90's to Deep learning

Deep learning

2000's : better, faster, stronger

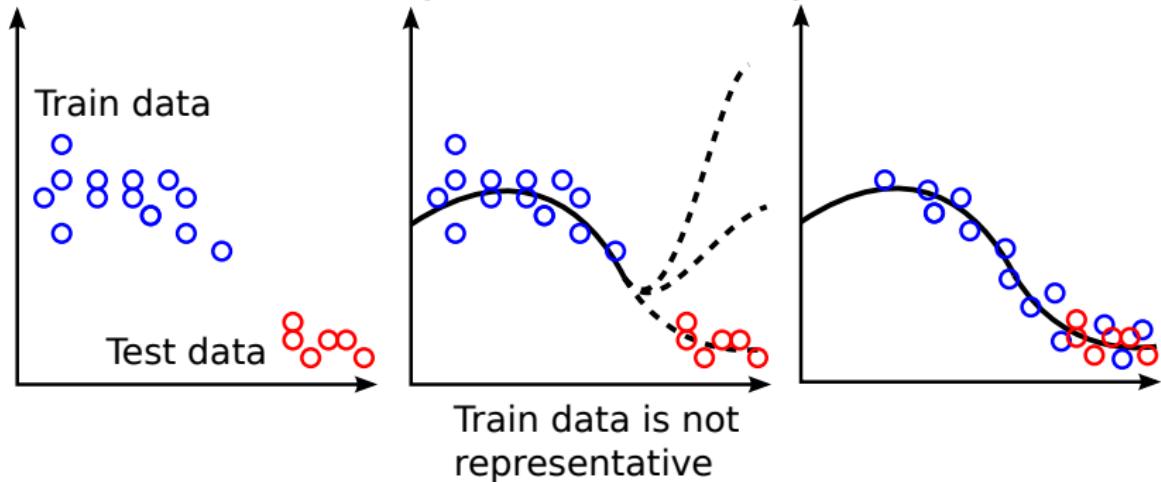
Do not forget the classics

What can we do with neural networks ?

Frameworks

- ▶ Representative
- ▶ Data augmentation
- ▶ Data normalization

Train data must be representative of the problem



Data normalization

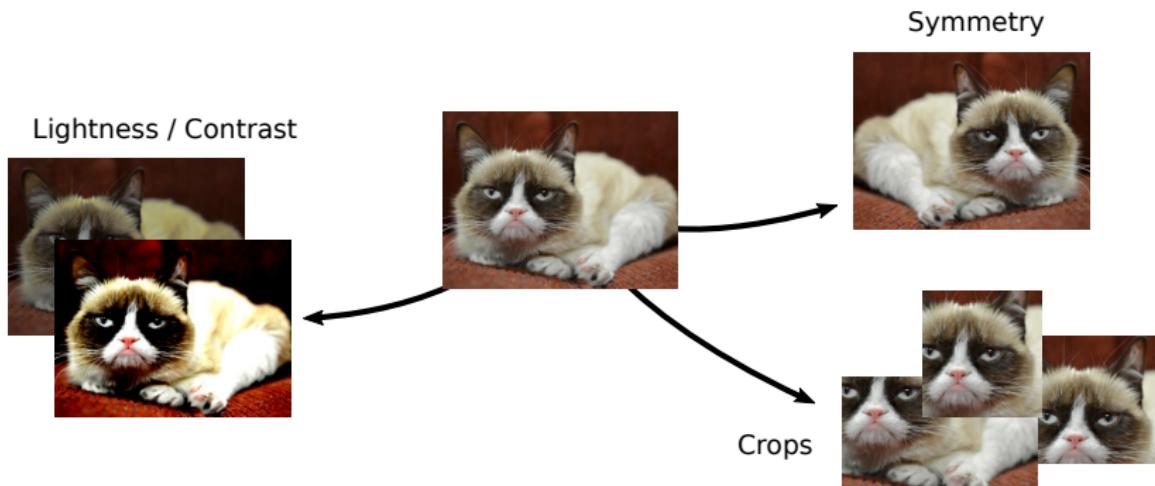
- ▶ Compute mean μ and standard deviation σ on the train set.
- ▶ Normalize input I (train and test) :

$$\hat{I} = \frac{I - \mu}{\sigma}$$

Data augmentation

Random variations of input parameters (images : lightness, contrast ...)

- ▶ train on a more representative set
- ▶ avoid learning on unwanted features



Problems and partial solutions

Problems

- ▶ Small amount of data
- ▶ Low computational power

Solutions ?

- ▶ Use classical approaches (Perceptron, SVM, ...)
- ▶ Use pre-trained deep features and learn classifier (**TP**)

Plan de la séance

From 90's to Deep learning

Deep learning

2000's : better, faster, stronger

Do not forget the classics

What can we do with neural networks ?

Frameworks

Image classification - LeNet vs VGG

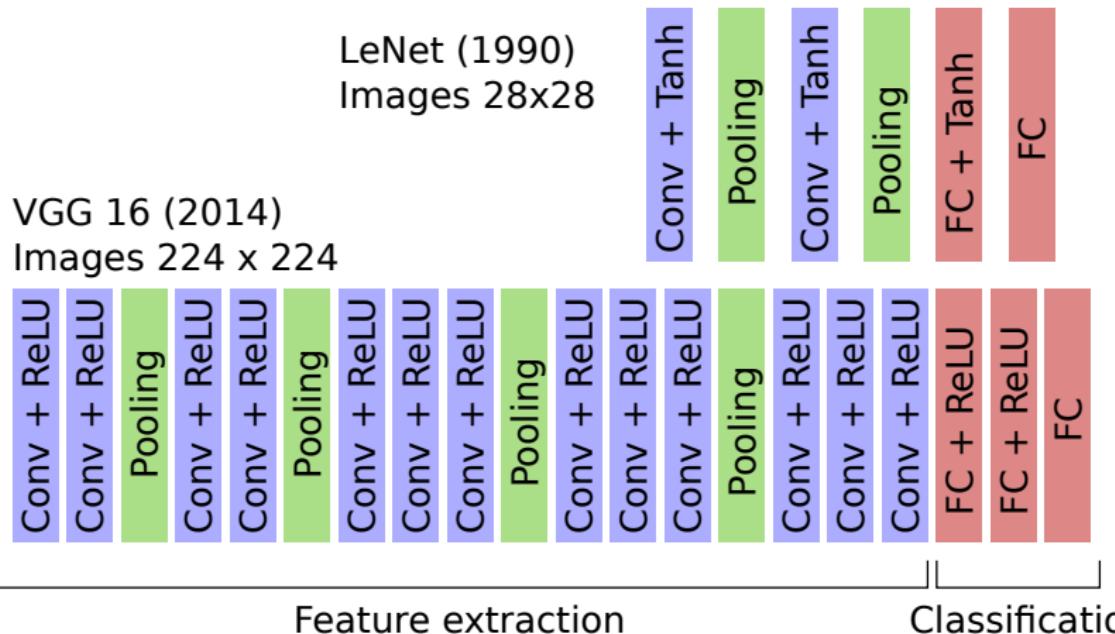
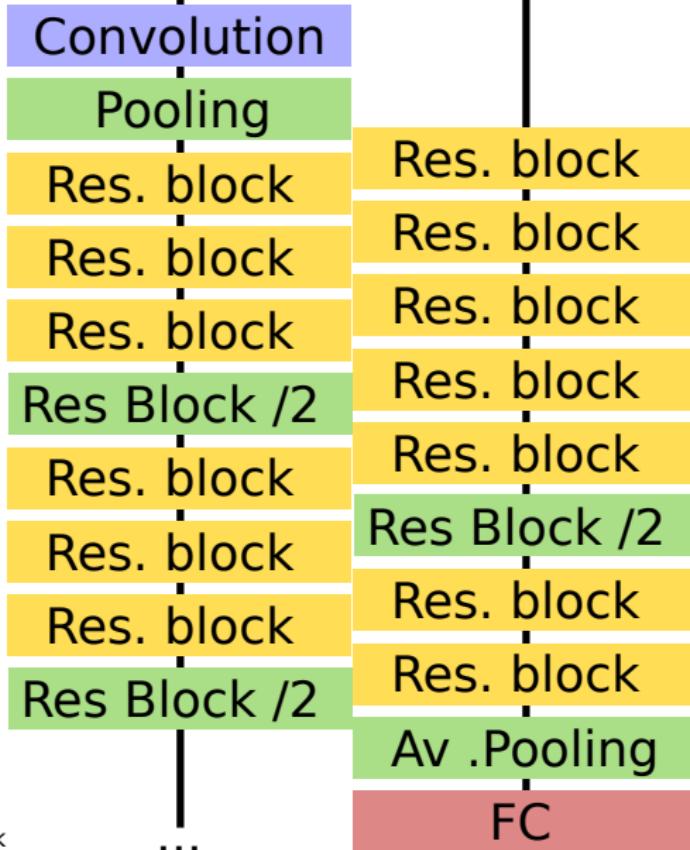
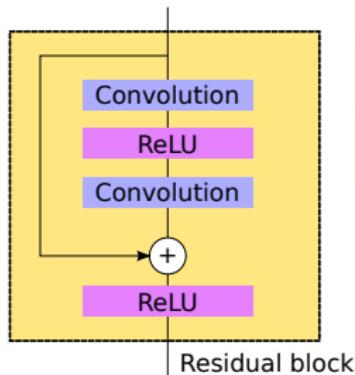


Image classification - ResNet



Semantic segmentation

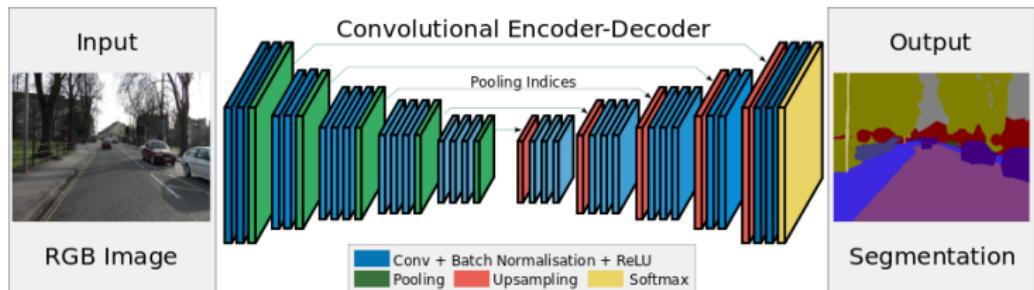
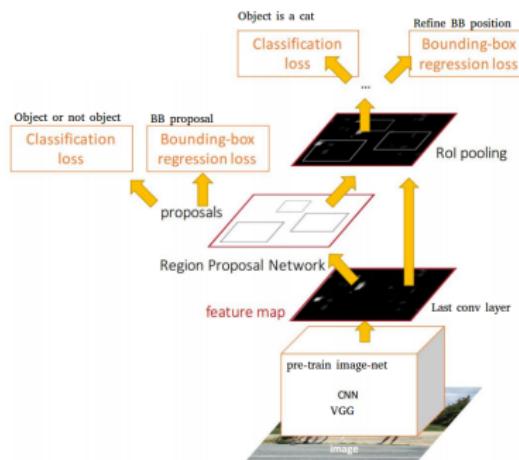
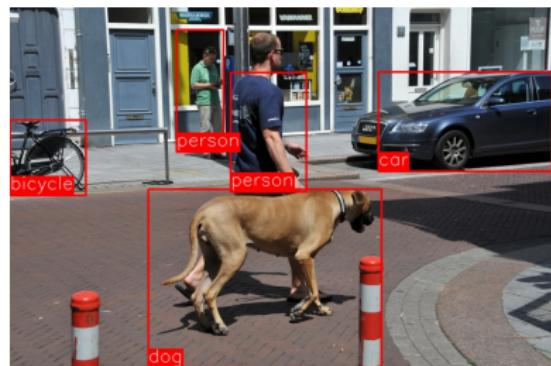


Image : <http://mi.eng.cam.ac.uk/projects/segnet/>

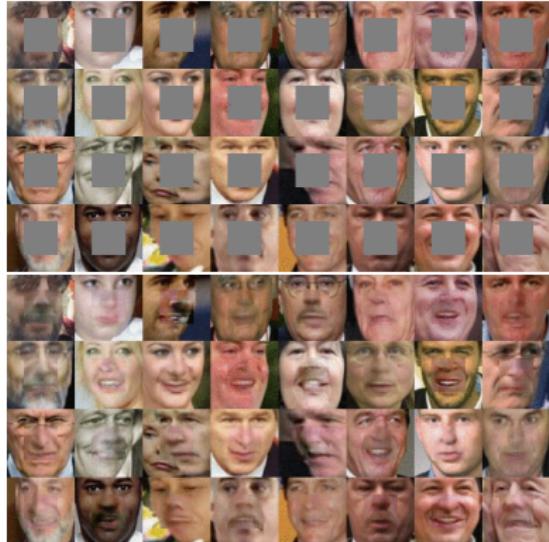
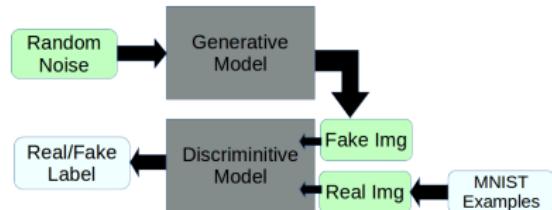
Detection



Data generation

Generative Adversarial Networks

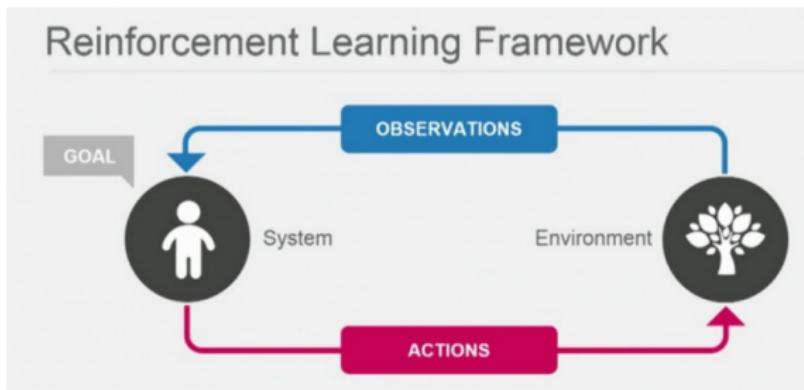
2 networks : one for data generation, one trying to discriminate real and generated data.



Reinforcement learning

Try and error learning

Reward and Penalty.



1

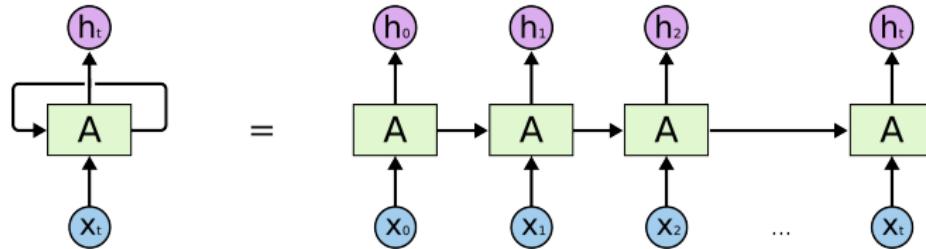
1. Image
[http://visteon.bg/2017/03/02/
machine-learning-algorithms-autonomous-cars/](http://visteon.bg/2017/03/02/machine-learning-algorithms-autonomous-cars/)

<http://visteon.bg/2017/03/02/>

ONERA

THE FRENCH AEROSPACE LAB

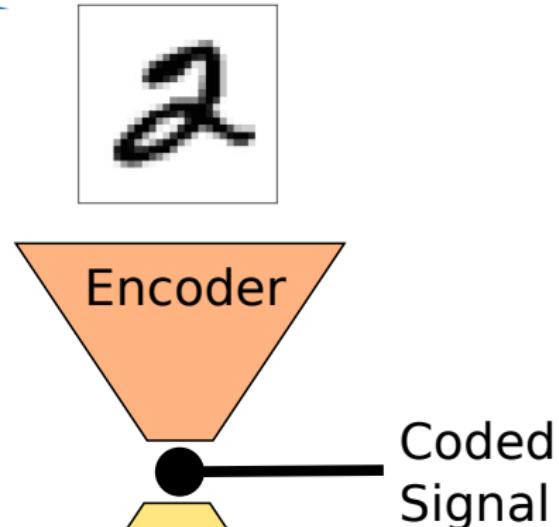
Recurrent neural networks



Dealing with sequences

Apply the same operation on each data and use a memory to propagate information in time.

Autoencoders



Plan de la séance

From 90's to Deep learning

Deep learning

2000's : better, faster, stronger

Do not forget the classics

What can we do with neural networks ?

Frameworks

Frameworks



Tensorflow and Keras

Principaux frameworks Python

Pytorch and Tensorflow : CPU and GPU, good maintenance, tutorials.

Keras run on top of Tensorflow (and Theano, CNTK...). It provides a high level API.

Conclusion

- ▶ Fast overview
- ▶ Few applications
- ▶ Practical work : finetuning and small recurrent neural networks