



# Deep Learning

*Thierry Artières*

*Ecole Centrale Marseille - Option DIGITALE*

September 19, 2019





## 1 Introduction

## 2 MLPs

- Basics
- Deeper in MLPs
- Computation graph
- Regularization

## 3 Programming

## 4 DNNs



# Outline

## 1 Introduction

## 2 MLPs

## 3 Programming

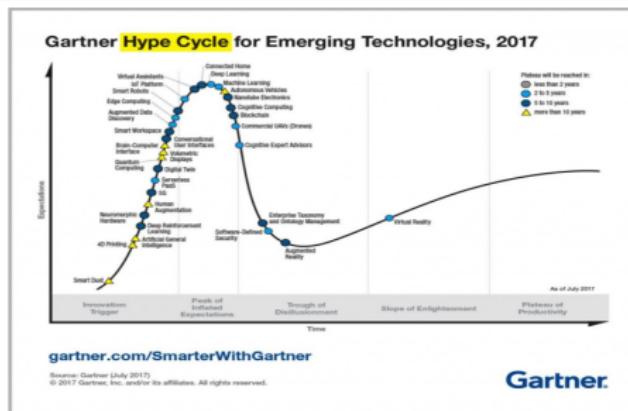
## 4 DNNs



## Where is AI?

### Constat

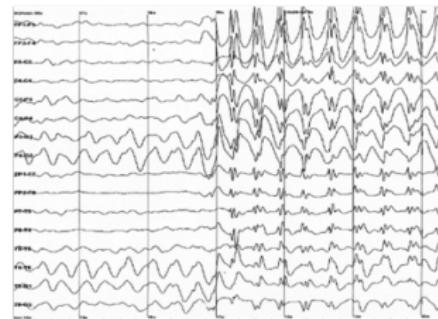
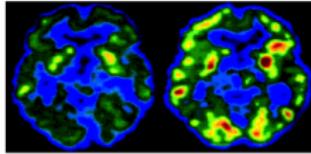
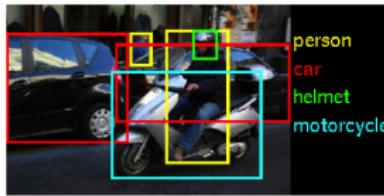
- The original AI was the General AI (IA forte)
- Today 60 years after the Dartmouth meeting
  - We have achieved some NIA results (IA faible)
  - We can start thinking more seriously about GAI
  - These are just the premises.





## At the heart of AI: Machine Learning (and Deep Learning)

Which algorithms to solve these tasks ?





# Machine Learning

## What is it for?

- Writing programs that solve a task while we don't even know how to write the algorithm
- Where a program takes some input and produce a corresponding output
- The program is learned from labeled data = pairs of (input, output)

## What is it?

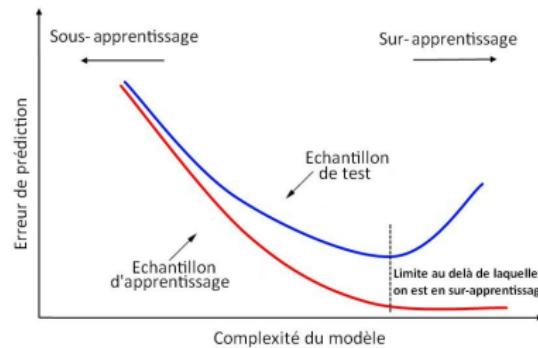
- Algorithms that enable learning a function  $f : x \in X \rightarrow y \in Y$  from a training dataset of samples
- The function must generalize well to data unseen at training time
- $x$  and  $y$  may be discrete, continuous, vectors, matrices, tensors, sequences ...



## Main difficulty

### Generalization

- It is “easy” to learn models that are perfect on training data
- But is is useless





# History of Neural Networks

## Key dates

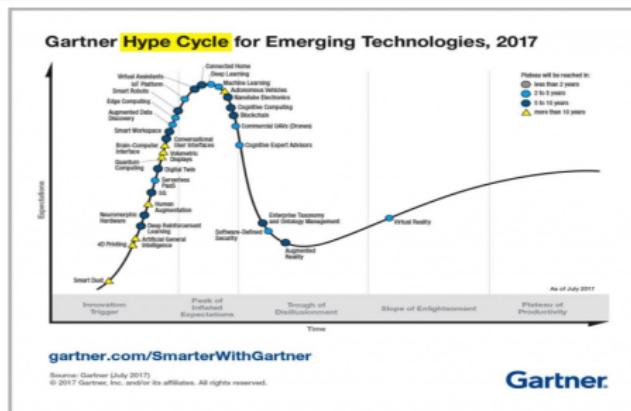
- 1943 : Formal neuron [McCulloch-Pitts]
- 1950 : Organization of neurons and learning rules [Hebb]
- 1960 : Perceptron [Rosenblatt]
- 1960 : Update rule [Widrow Hoff]
- 1969 : Limitations of the Perceptron [Minsky]
- 1980s : Back-propagation [Rumelhart and Hinton]
- 1990s : Convolutional Networks [LeCun and al.]
- 1990s: Long Short Term Memory networks [Hochreiter and Schmidhuber]
- 2006 : Paper on Deep Learning in Nature [Hinton and al.]
- 2012 : Imagenet Challenge Win [Krizhevsky, Sutskever, and Hinton]
- 2013 : First edition of ICLR
- 2013 : Memory networks [Weston and al.]
- 2014 : Adversarial Networks [Goodfellow and al.]
- 2014 : Google Net [Szegedy and al.]
- 2015 : Residual Networks [He et al.]
- ...



## AI today

Where are we?

- The original AI (Dartmouth workshop) was General AI (IA forte)
- Today 60 years after Dartmouth
  - We have achieved NIA results (IA faible)
  - We can start thinking more seriously about GAI
  - These are just the premises.



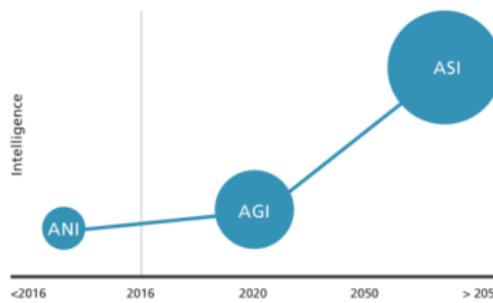


## AI today

### Where are we?

- The original AI (Dartmouth workshop) was General AI (IA forte)
- Today 60 years after Dartmouth
  - We have achieved NIA results (IA faible)
  - We can start thinking more seriously about GAI
  - These are just the premises.

The evolution of artificial intelligence



Source: UBS, as of 15 August 2016



## Deep Learning today

Spectacular breakthroughs - fast industrial transfer

- Images, Videos, Audio, Speech, Texts
- Successful setting
  - Structured data (temporal, spatial...)
  - Huge volumes of data
  - Huge models (millions of parameters)
  - Huge storage and computing resources (GPU, TPU)

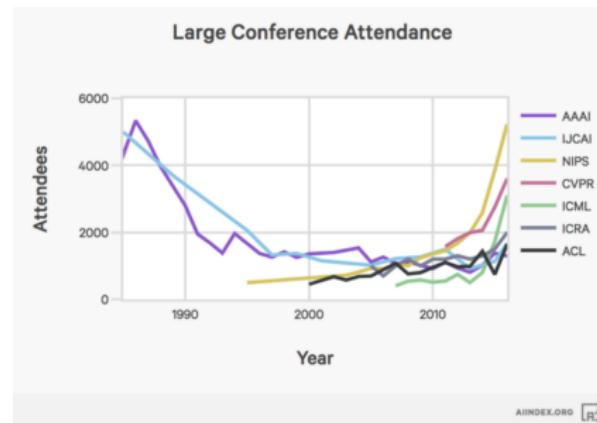
	VGGNet	DeepVideo	GNMT
Used For	Identifying Image Category	Identifying Video Category	Translation
Input	Image 	Video 	English Text 
Output	1000 Categories	47 Categories	French Text
Parameters	140M	~100M	380M
Data Size	1.2M Images with assigned Category	1.1M Videos with assigned Category	6M Sentence Pairs, 340M Words
Dataset	ILSVRC-2012	Sports-1M	WMT'14



## Machine Learning and Deep Learning today

### Spectacular diffusion and activity

- Machine Learning and Deep Learning Conferences sold out early
- More attendees than ever seen in computer science conferences
- Exponential growth
- Semantic change in what AI means



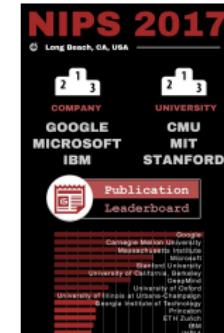


# Machine Learning and Deep Learning today

Topics, trends and who's who?

- Mix between academics and companies
- Extreme popularity of Deep Learning topics
- Birth of the International Conference on Learning Representation (2014)

Rank	Day	Name	Marked	Like
1	1	Tutorials Hall A	2789	287
2	2	Deep Learning, Applications	2364	289
3	3	Deep Learning	1831	163
4	3	Reinforcement Learning, Deep Learning	1592	140
5	1	Optimization	1522	130
6	1	Tutorials Hall C	1344	135
7	1	Algorithms	1307	137
8	2	Theory	1288	83
9	2	Algorithms Optimization	1223	107
10	4	Deep Learning, Algorithms	1202	113
11	4	Deep Reinforcement Learning	1202	43
12	2	Invited talk: Kate Crawford: The Trouble with Bias	1162	71
13	3	Reinforcement Learning, Algorithms, Applications	1156	134
14	3	Invited talk: Pieter Abbeel: Deep Learning for Robotics	1087	61
15	1	Tutorials Grand Ballroom	1082	132





## DL research is going very fast !!

Example of an emerging topic: Generative Adversarial Networks

- First publication : 2014 by Ian J. Goodfellow, and al.
- Hundreds of publications (close to a thousand) papers since

New publication mode

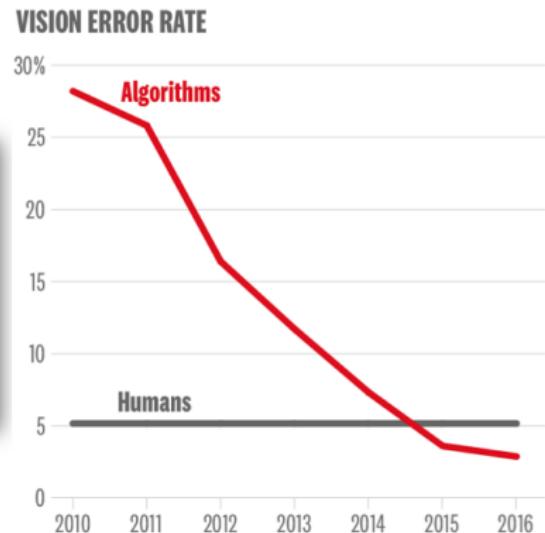
- Wasserstein GANs, Martin Arjovsky and al.
  - Published on arXiv : Jan 2017
  - Published at ICML in Aout 2017
- Improved Training of Wasserstein GANs by Ishaan Gulrajani and al.
  - Published on arXiv : March 2017
  - Published at NIPS in December 2017
- Improving the Improved Training of Wasserstein GANs: A Consistency Term and Its Dual Effect by Xiang Wei and al.
  - Published on Openreview : Oct 2017
  - Accepted as poster at ICLR in 2018 (April 2018)



# Spectacular breakthrough

## Computer vision

### Real time Object recognition



SOURCE ELECTRONIC FRONTIER FOUNDATION © HBR.ORG



## Spectacular breakthrough

### Speech

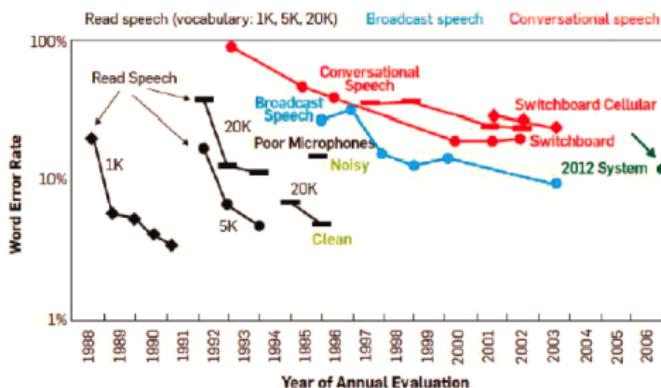


FIGURE 2.7 Historical progress on reducing the word error rate in speech recognition systems for different kinds of speech recognition tasks. Recent competency for the “difficult switchboard” task (human conversation in the wild) is marked with the green dot. SOURCE: X. Huang, J. Baker, and R. Reddy, 2014, A historical perspective of speech recognition, *Communications of the ACM* 57(1):94-103, doi:10.1145/2500887. © 2014, Association of Computing Machinery, Inc. Reprinted with permission.



## Spectacular breakthrough

### Games

- BackGammon, Chess, Go

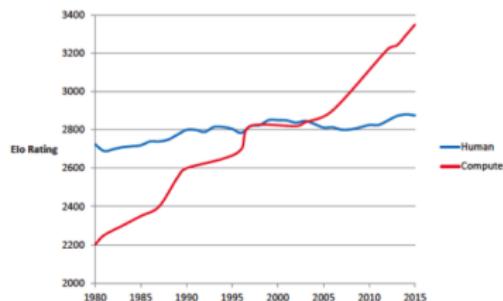


FIGURE 2.3 Elo scores—a measure of competency in competitive games—showing the chess-playing competency of humans and machines, measured over time. SOURCE: Courtesy of Murray Campbell.





# Spectacular breakthrough

## Image generation

### Recent Nvidia results





## Spectacular breakthrough

Should we still trust what we see?





## Why now ?

Huge training resources for huge models

- Huge volumes of training data
- Huge computational resources (clusters of GPUs)

Advances in understanding optimizing NNs

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

Faster diffusion than ever

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



# Outline

1 Introduction

2 MLPs

- Basics
- Deeper in MLPs
- Computation graph
- Regularization

3 Programming

4 DNNs



## Where DL comes from

F. Rosenblatt 1958 : The perceptron



## Basics

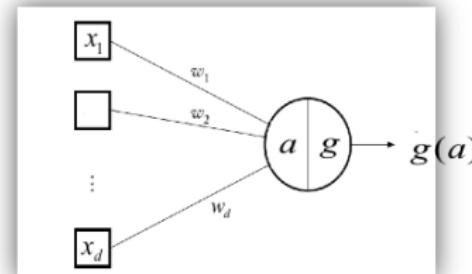
# A single Neuron

## One Neuron

- Elementary computation

$$\text{activation} = w^T \cdot x = \sum_j w_j x_j + w_0$$

$$\text{output} = g(a(x))$$

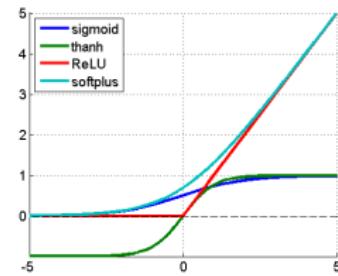


## Non linearity : $g$

- Sigmoide, Hyperbolic tangent, Gaussian
- Rectified Linear Unit (ReLU)

$$f(x) = 0 \text{ if } x \leq 0$$

$$= x \text{ otherwise}$$



## Basics



# Multi Layer Perceptron (MLP)

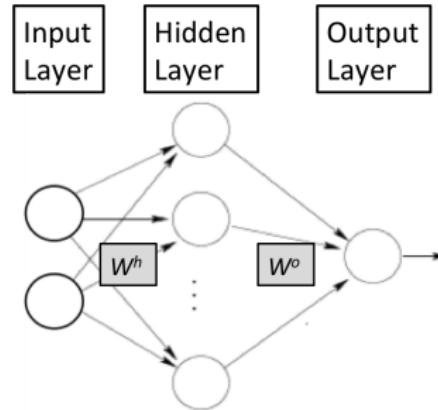
## Structure

- Organization in successive layers
  - Input layer
  - Hidden layers
  - Output layer

## Function implemented by a MLP

$$g(W^o \cdot g(W^h x))$$

- Inference: Forward propagation from input to output layer
  - Fill the input layer with  $x$ :  $h_0 = x$
  - Iterate from the first hidden layer to the last one
    - $h^l = W^l \times h^{l-1}$
    - $h^l = g(h^l)$



## Basics



# Learning a MLP

## Learning as an optimization problem

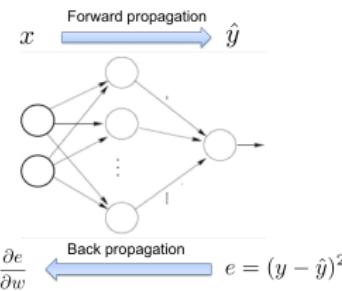
- Objective function of parameters set  $w$  for a given training set  $T$

$$\begin{aligned} C(w) &= F(w) + R(w) \\ &= \sum_{(x,y) \in T} L_w(x, y, w) + \|w\|^2 \end{aligned}$$

- Gradient descent optimization:  $w = w - \epsilon \frac{\partial C(w)}{\partial w}$

## Backpropagation

- Use chain rule for computing derivative of the loss with respect to all weights in the NN





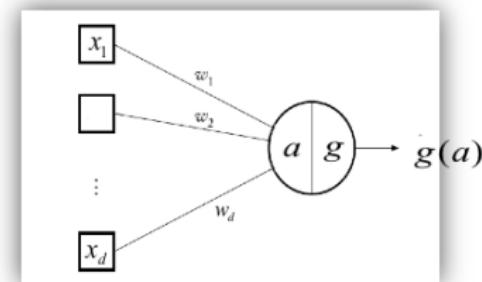
## A single ReLU Neuron

### One Neuron

- Elementary computation

$$\text{activation} = w^T \cdot x = \sum_j w_j x_j + w_0$$

$$\text{output} = \text{ReLU}(a(x))$$





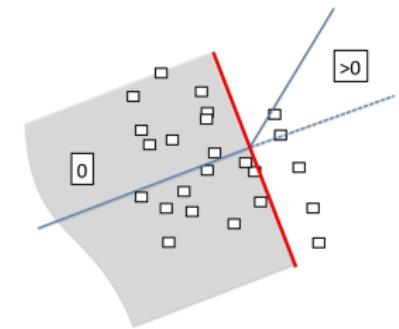
## A single ReLU Neuron

### One Neuron

- Elementary computation

$$\text{activation} = w^T \cdot x = \sum_j w_j x_j + w_0$$

$$\text{output} = \text{ReLU}(a(x))$$



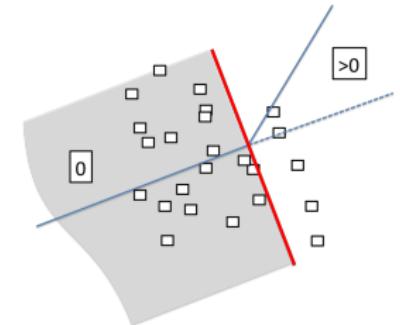
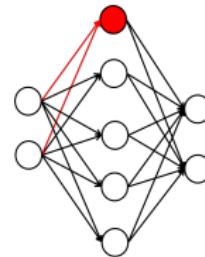
## Deeper in MLPs

## What a MLP may compute



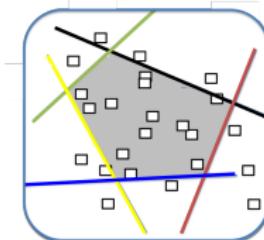
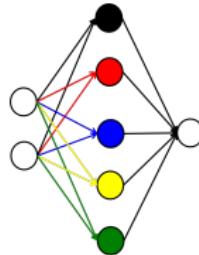
What does a hidden neuron

- Divides the input space in two



Combining multiple hidden neurons

- Allows identifying complex areas of the input space
- New (distributed) representation of the input

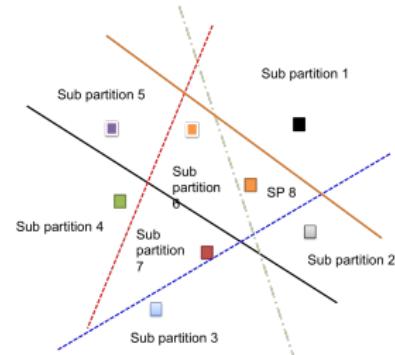
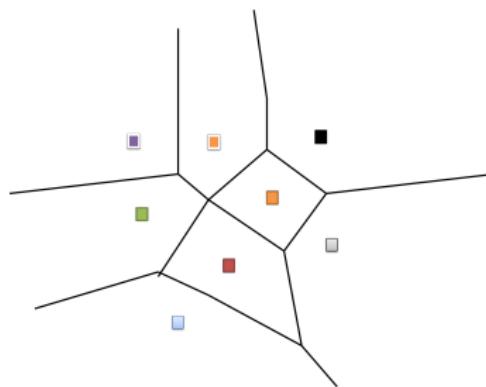


## Deeper in MLPs

## MLP compute distributed representations

Might be much more efficient than non distributed ones

Somewhat the number of regions in which a NN architecture may divide the input space is a measure of its capacity





## MLP = Universal approximators

One layer is enough !

- Theorem [Cybenko 1989]: Let  $\phi(\cdot)$  be a nonconstant, bounded, and monotonically-increasing continuous function. Let  $I_m$  denote the  $m$ -dimensional unit hypercube  $[0, 1]^m$ . The space of continuous functions on  $I_m$  is denoted by  $\mathcal{C}(I_m)$ . Then, given any  $\epsilon > 0$ , there exists an integer  $N$ , such that for any function  $f \in \mathcal{C}(I_m)$ , there exist real constants  $v_i, b_i \in \mathbb{R}$  and real vectors  $w_i \in \mathbb{R}^m$ , where  $i = 1, \dots, N$ , such that we may define:

$$F(x) = \sum_{i=1}^N v_i \phi(w_i^T x + b_i)$$

as an approximate realization of the function  $f$  where  $f$  is independent of  $\phi$ ; that is :  $|F(x) - f(x)| < \epsilon$  for all  $x \in I_m$ . In other words, functions of the form  $F(x)$  are dense in  $\mathcal{C}(I_m)$ .

- Existence theorem only
- Many reasons for not getting good results in practice

Deeper in MLPs

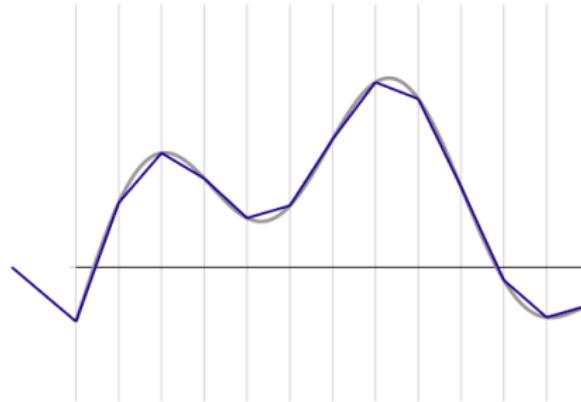


## MLP = Universal approximators: Proof (From Francos Fleuret slides)

simple real function

- Any function in  $C([a, b], \mathbb{R})$  may be approximated with any desired precision with a linear combination of translated / scaled Relu functions.

$$f(x) = R(w_1x + b_1) + R(w_2x + b_2) + \dots$$



## Deeper in MLPs



## MLP = Universal approximators: Proof (From Francois Fleuret slides)

More general functions in  $\Psi \in \mathcal{C}([0, 1]^D, \mathbb{R})$

- Previous result holds for sin function

$$\forall A > 0, \epsilon > 0, \exists N, (\alpha_n, a_n) \in \mathbb{R} \times \mathbb{R}, n = 1, \dots, N,$$

$$\text{s.t. } \max_{x \in [-A, A]} |\sin(x) - \sum_{n=1}^N \alpha_n R(x - a_n)| \leq \epsilon$$

- Density of Fourier series

$$\forall \Psi, \delta > 0, \exists M, (v_m, \gamma_m, c_m) \in \mathbb{R}^D \times \mathbb{R} \times \mathbb{R}, m = 1, \dots, M,$$

$$\text{s.t. } \max_{x \in [0, 1]^D} |\Psi(x) - \sum_{m=1}^M \gamma_m \sin(v_m \cdot x - c_m)| \leq \delta$$

- Result easily follows

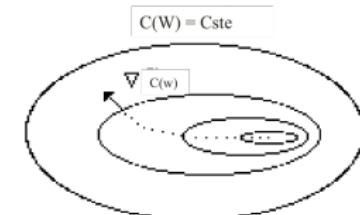
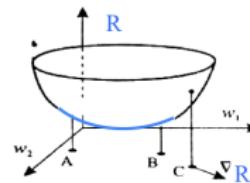
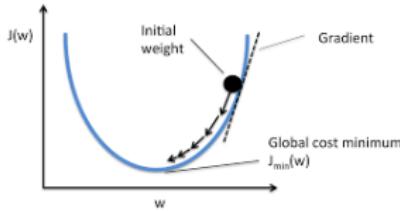


## Gradient Descent Optimization

### Gradient Descent Optimization

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

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



⇒ Few illustrations in these slides are taken from [LeCun et al, 1993], [Fei Fei Li lecture 6], and from S. Ruder's blog

## Gradient Descent: Tuning the Learning rate

Two classes Classification problem

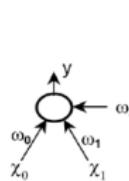
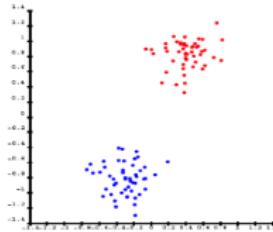


Fig. 9. Simple linear network.



Weight trajectory for two different gradient step settings.

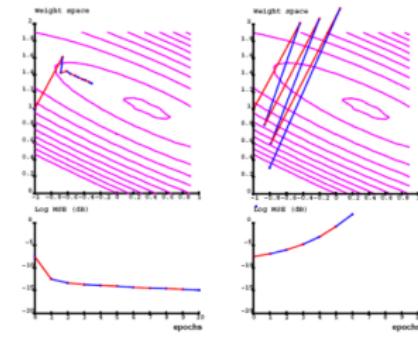


Fig. 11. Weight trajectory and error curve during learning for (a)  $\eta = 1.5$  and (b)  $\eta = 2.5$ .

Images from [LeCun et al.]



Deeper in MLPs

## Gradient Descent: Tuning the Learning rate

### Effect of learning rate setting

- Assuming the gradient direction is good, there is an optima value for 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

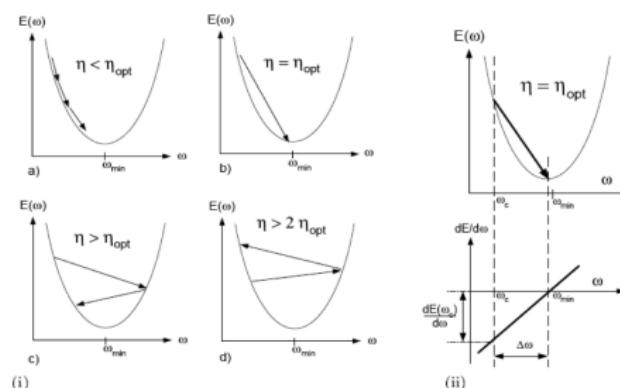


Fig. 6. Gradient descent for different learning rates.

Images from [LeCun et al.]



## Optimal learning rate and convergence speed

### Second order point of view

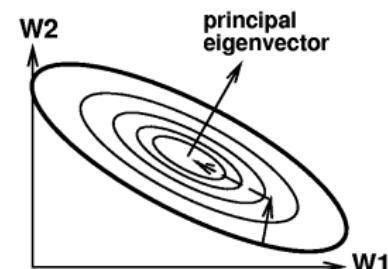
- Taylor expansion, noting  $\nabla^2 C(w)$  the Hessian (a  $N \times N$  matrix with  $N$  a model with parameters )

$$\nabla C(w)|_{w'} = \nabla C(w)|_w + \nabla^2 C(w)(w' - w)$$

- $\rightarrow$  optimum rule (setting  $\nabla C(w)|_{w'}$  to 0):

$$w' = w - (\nabla^2 C(w))^{-1} \nabla C(w)$$

- Optimal move not in the direction of the gradient
- In Order 1 Gradient descent the optimal the optimal value of  $\epsilon$  depends on eigen values of the Hessian  $\nabla^2 C(w)$



[Lecun et al, 93]



## Gradient Descent: Stochastic, Batch and mini batchs

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

### 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  $C(\mathbf{w})$  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.

Deeper in MLPs



## Optimization routines

Many SGD variants popular in DL

- SGD with Momentum
- Nesterov accelerated gradient
- Adagrad
- Adadelta
- RmsProp
- ...



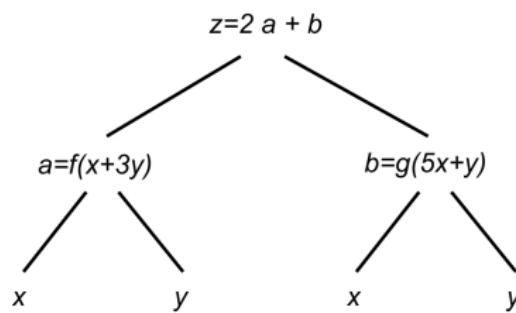
## Computation graph

## Gradient Computation: Chain rule

## Gradient of a function

$$z = 2 \times f(x + 3 \times y) + 6 \times g(5 \times x + y)$$

$$\Rightarrow \frac{\partial z}{\partial x} \Big|_{x,y} = 2 \times f'(x + 3 \times y) + 30 \times g'(5 \times x + y)$$



## Equivalent computation with the Chain rule

Set  $a = f(x + 3 \times y)$  and  $b = g(5 \times x + y)$

$$\Rightarrow z = 2 \times a + 6 \times b$$

$$\Rightarrow \frac{\partial z}{\partial x} = \frac{\partial z}{\partial a} \times \frac{\partial a}{\partial x} + \frac{\partial z}{\partial b} \times \frac{\partial b}{\partial x}$$

With:

$$\frac{\partial z}{\partial a} = 2 \text{ and } \frac{\partial z}{\partial b} = 6$$

$$\frac{\partial a}{\partial x} = f'(a \times x + 3 \times y)$$

$$\frac{\partial b}{\partial x} = 5 \times g'(5 \times x + y)$$

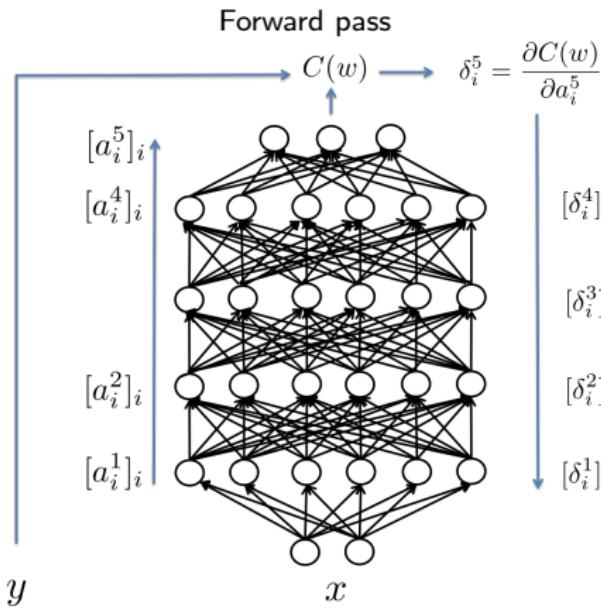
$$\frac{\partial a}{\partial y} = 3 \times f'(a \times x + 3 \times y)$$

$$\frac{\partial b}{\partial y} = 5 \times g'(5 \times x + y)$$



## Computation graph

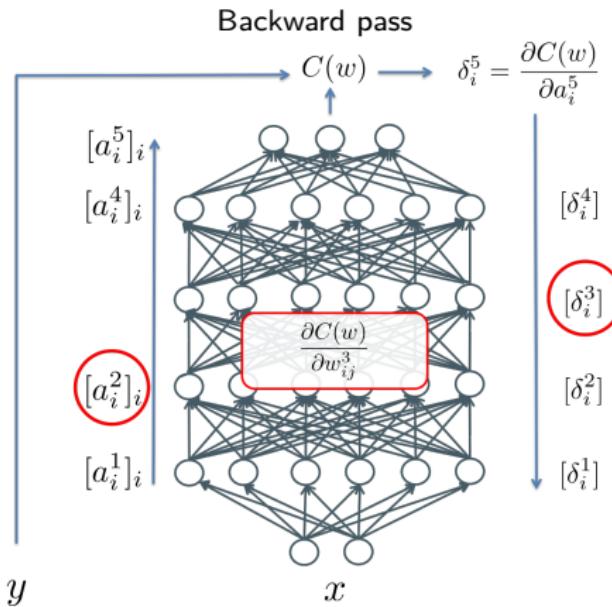
## BackPropagation





## Computation graph

## BackPropagation





## Computation graph

## Gradient computation in MLPs: Stochastic case

## Notations

- Activation function on every layer:  $g$  — Number of layer :  $L$
- Activity of neuron  $i$  in layer  $l$ ,  $a_i^l$  — Output of neuron  $i$  in layer  $l$ ,  $h_i^l = g(a_i^l)$ , and  $o_i^L = g(a_i^L)$
- Weight from a neuron  $j$  of layer  $l - 1$  to neuron  $i$  in layer  $l$  :  $w_{ij}^l$
- Example considered for computing gradient  $(x, y)$
- Squared loss :  $C(w) = \|\mathbf{o}^L - \mathbf{y}\|^2$

## Gradient wrt. last layer weights

- Gradient wrt cell's output  $\frac{\partial C(w)}{\partial o_i^L} = 2(o_i^L - y_i)$
- Gradient wrt cell's activity  $\delta_i^L = \frac{\partial C(w)}{\partial a_i^L} = \frac{\partial C(w)}{\partial o_i^L} \frac{\partial o_i^L}{\partial a_i^L} = 2(o_i^L - y_i)g'(a_i^L)$
- Gradient wrt weights arriving to output cells

$$\frac{\partial C(w)}{\partial w_{ij}^L} = \frac{\partial C(w)}{\partial a_i^L} \frac{\partial a_i^L}{\partial w_{ij}^L} = \delta_i^L \times h_j^{L-1}$$



## Computation graph

## Gradient computation in MLPs: Stochastic case (continues)

## Gradient wrt. last hidden layer (LHL) weights

- Gradient wrt LHL cell's activity  $\delta_j^{L-1} = \frac{\partial C(w)}{\partial a_j^{L-1}} = \sum_i \frac{\partial C(w)}{\partial a_i^L} \frac{\partial a_i^L}{\partial a_j^{L-1}} = \sum_i \delta_i^L w_{ij}^L g'(a_j^{L-1})$
- Gradient wrt weights arriving to a LHL cell

$$\frac{\partial C(w)}{\partial w_{jk}^{L-1}} = \delta_j^{L-1} \times h_k^{L-2}$$



## Computation graph

# Gradient computation in MLPs

Forward propagation of activities, for an input example  $x$

- Fill the input layer with  $x$ :  $h^0 = x$
- Iterate from the first hidden layer to the last one
  - $h^l = W^l \times h^{l-1}$
  - $h^l = a(h^l)$

Backward computation of the error

- Compute the output error  $\delta^L$
- Iterate from the last hidden layer to the first one
  - Compute  $\delta^L$  from  $\delta^{L-1}$

Computing gradient

- For each weight  $w_{jk}^l$  of every layer compute the gradient using  $\delta_j^l$  and  $o_k^{l-1}$



## Guiding the learning through regularization

### Regularization

- Constraints on weights (L1 or L2)
- Constraints on activities (of neurons in a hidden layer) → induces sparsity
  - 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
- Early stopping

## Regularization



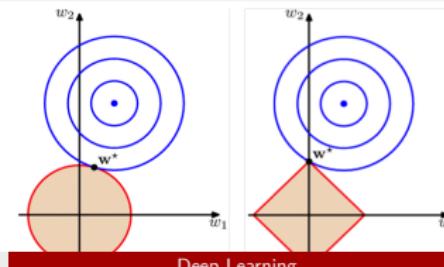
## Constraints on weights

### L2 norm on weights (known as Weight Decay)

- Penalizing the weights through adding a weighted L2 norm  $\lambda\|w\|^2$  to the loss
- It is equivalent to defining a family of models such that  $\|w\|^2 \leq C_\lambda$  with  $C_\lambda$  increasing when  $\lambda$  decreases
- L2 norm penalization  $\leftrightarrow$  diminishing the space of functions implemented with the network architecture

### L2 and L1 norms

- L2 norm move useless weights to 0 (without reaching 0)
- L1 norm set useless weights to 0

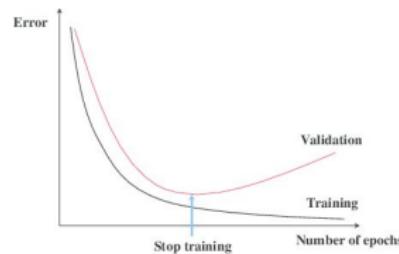




## Regularization

## Early stopping and callbacks

### Principle



- Early stopping monitors performance (loss) on validation set
- Stopped before it reaches a plateau and starts increasing
- Related to the idea that the implemented model's capacity increases with the number of iteration
  - Think of small weights initialization and sigmoid activation
  - ⇒ at the beginning the model is a linear one !

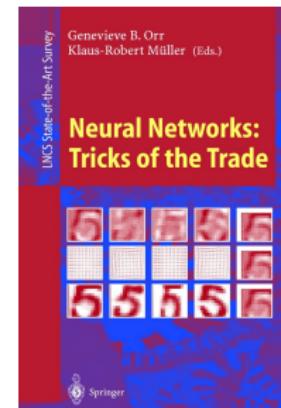
## Regularization



## Lots of tricks to favor convergence

And more...

- Weight Initialization
- Gradient step setting
- ...
- ⇒ Despite appearances not all is automatic





## Outline

1 Introduction

2 MLPs

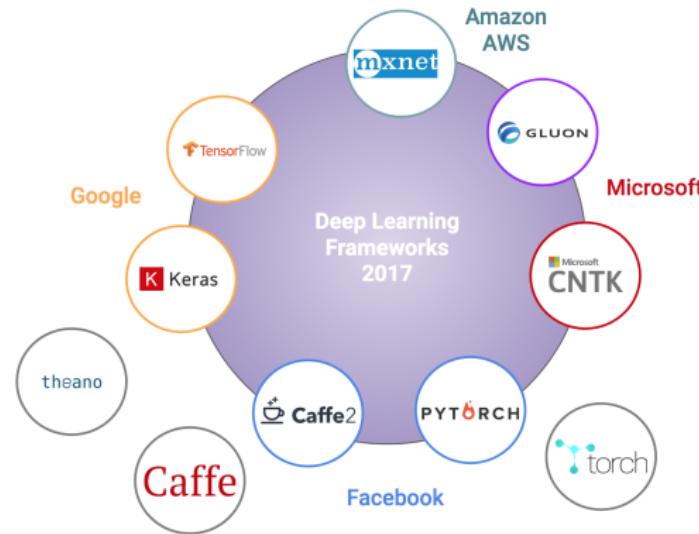
3 Programming

4 DNNs



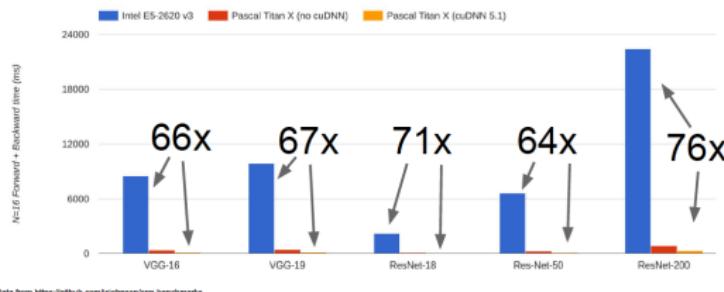
## Platforms

Large and active community (forums, models are available when published...) for each of these





## GPU and CPU



Model  
is here



Data is here

If you aren't careful, training can bottleneck on reading data and transferring to GPU!

### Solutions:

- Read all data into RAM
- Use SSD instead of HDD
- Use multiple CPU threads to prefetch data

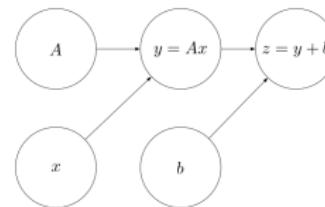


## Computation graph

### Computation graph for a calculus

One may build a computation graph from the calculus definition

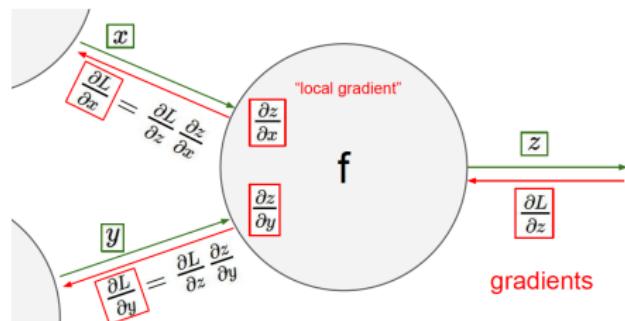
$$z = Ax + b$$



### Automatic differentiation

From a computation graph one may automatically compute the backward differentiation graph !

- Different rules to apply according to the operation yielding  $z$  from  $x$  and  $y$

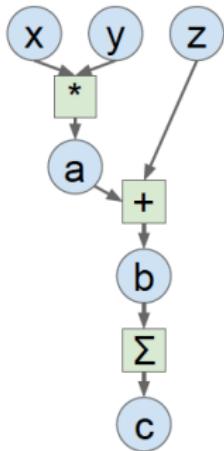


From Fei Fei Li slides



## Computation graph and Pytorch

### Computational Graphs



```

[3] import torch
from torch.autograd import Variable

N, D = 3, 4
x=Variable(torch.randn(N,D).cuda(),requires_grad=True)
y=Variable(torch.randn(N,D).cuda(),requires_grad=True)
z=Variable(torch.randn(N,D).cuda(),requires_grad=True)

[6] a = x*y
b = a + z
c = torch.sum(b)

c.backward()

print (x.grad.data)
print (y.grad.data)
print (z.grad.data)

[7]
tensor([[-2.4946, -1.7749, -2.8303, -1.0450],
       [ 1.8087, -0.8123,  1.4324, -0.7497],
       [ 0.4153, -0.7573, -0.3054,  1.8146]], device='cuda:0')
tensor([[[-0.2363, -1.8247, -3.2515,  4.3729],
       [-0.6283,  1.9725, -3.6697, -1.4272],
       [ 0.8991, -0.2417, -0.2456, -2.4684]], device='cuda:0')
tensor([[2., 2., 2., 2.],
       [2., 2., 2., 2.],
       [2., 2., 2., 2.]], device='cuda:0')
  
```



## Example (pytorch)

### Mnist Classifier (model definition)

```
class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv1 = nn.Conv2d(1, 10, kernel_size=5)
        self.conv2 = nn.Conv2d(10, 20, kernel_size=5)
        self.conv2_drop = nn.Dropout2d()
        self.fc1 = nn.Linear(320, 50)
        self.fc2 = nn.Linear(50, 10)

    def forward(self, x):
        x = F.relu(F.max_pool2d(self.conv1(x), 2))
        x = F.relu(F.max_pool2d(self.conv2_drop(self.conv2(x)), 2))
        x = x.view(-1, 320)
        x = F.relu(self.fc1(x))
        x = F.dropout(x, training=self.training)
        x = self.fc2(x)
        return F.log_softmax(x, dim=1)
```



## Example (pytorch)

### Mnist Classifier (model training)

```
def train(args, model, device, train_loader, optimizer, epoch):
    model.train()
    for batch_idx, (data, target) in enumerate(train_loader):
        data, target = data.to(device), target.to(device)
        optimizer.zero_grad()
        output = model(data)
        loss = F.nll_loss(output, target)
        loss.backward()
        optimizer.step()
        if batch_idx % args.log_interval == 0:
            print('Train Epoch: {} [{}/{} ({:.0f}%)]\tLoss: {:.6f}'.format(
                epoch, batch_idx * len(data), len(train_loader.dataset),
                100. * batch_idx / len(train_loader), loss.item()))
```



## Outline

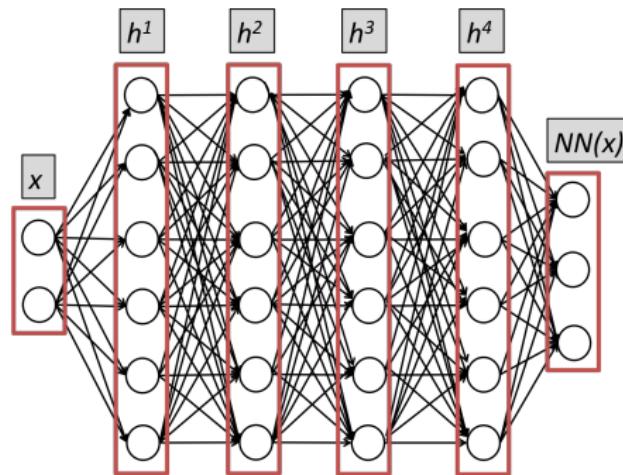
- 1 Introduction
- 2 MLPs
- 3 Programming
- 4 DNNs



## Deep Learning = Representation Learning

Hierarchy of representation spaces by successive hidden layers

$$h^i(x) = g(W^i \times h^{i-1}(x))$$

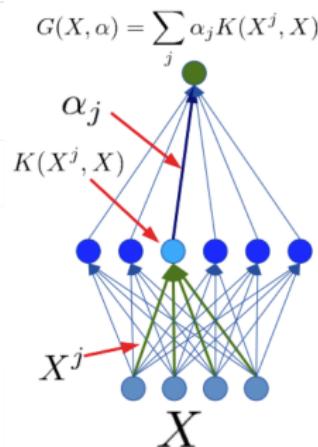




## Shallow vs deep models

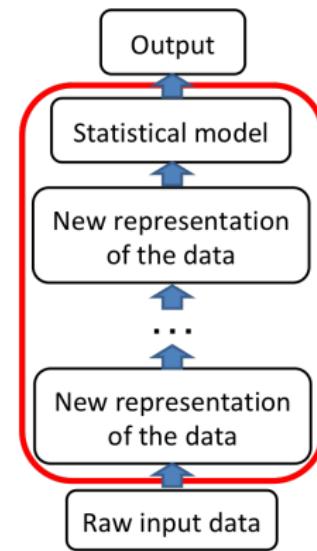
From [LeCun tutorial Statlearn]

- Neural Networks with one hidden layer are shallow models
- SVMs are shallow models



DL

- Joint learning of a hierarchy of representations and of a prediction model

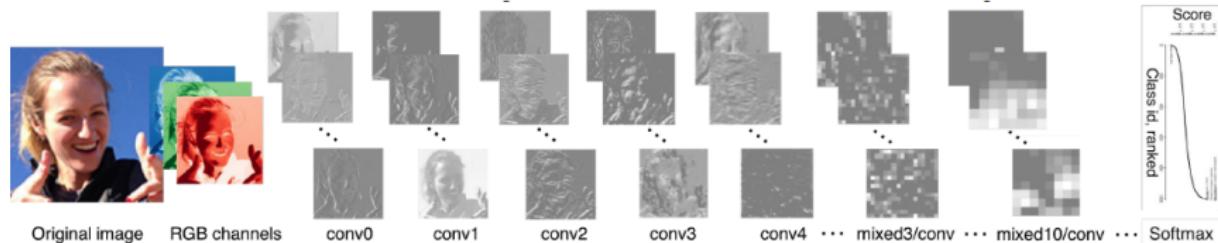




## Feature hierarchy : from low to high level

What feature hierarchy means ?

- Low-level features are shared among categories
- High-level features are more global and more invariant



[From Taigman et al., 2014]