

Bases de l'IA

Traitement d'image

Réseaux neuronaux convolutifs

Elena CABRIO

elena.cabrio@univ-cotedazur.fr

Plan pour cette séance

- Classification des images
 - Classificateurs linéaires
 - Réseaux neuronaux convolutifs

Classification des images : une tâche essentielle de la vision par ordinateur

Supposons un ensemble donné d'étiquettes discrètes, par exemple {chat, chien, vache, pomme, tomate, camion, ... }

$f(\text{apple})$ = “apple”

$f(\text{tomato})$ = “tomato”

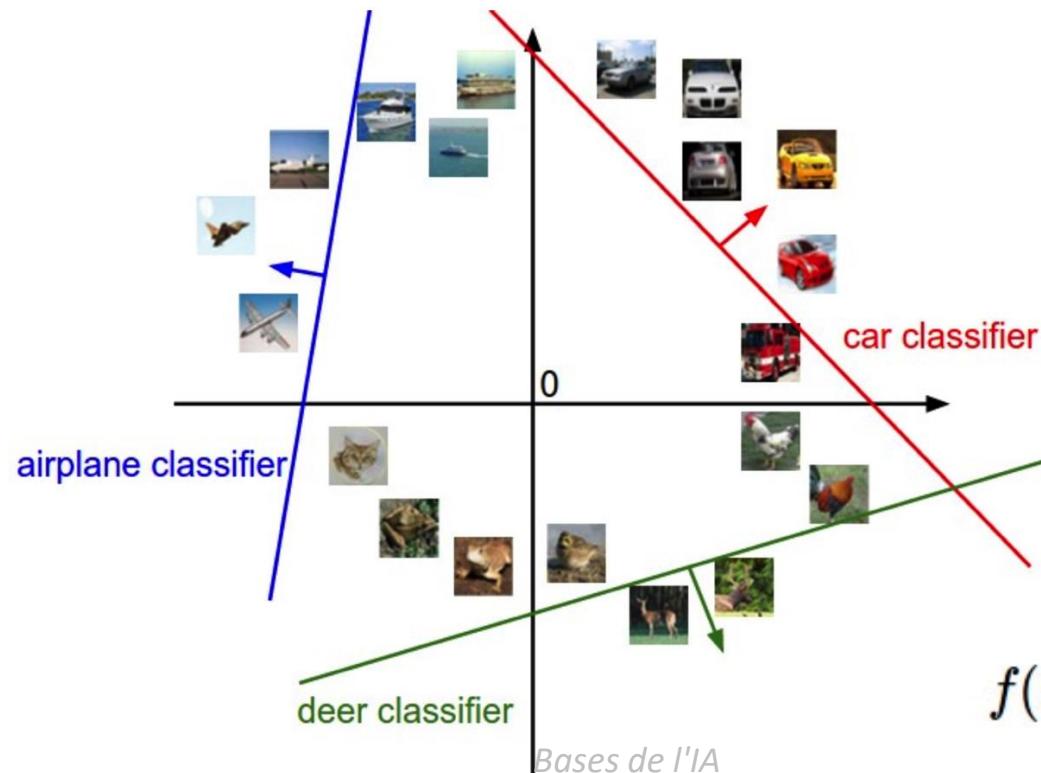
$f(\text{cow})$ = “cow”

Récapitulation : classification linéaire

- **Fonction de score et fonction de perte**
 - La fonction de score associe une instance de données d'entrée (par exemple, une image) à un vecteur de scores, un pour chaque catégorie.
 - La dernière fois, notre fonction de score était basée sur un classificateur linéaire.
$$f(x, W) = Wx + b$$

f : fonction de score
x : instance d'entrée
W, b : paramètres d'une fonction linéaire
- Trouver **W** et **b** pour minimiser une perte, par exemple la perte d'entropie croisée.

Les classificateurs linéaires séparent l'espace des features en demi-espaces



Réseaux de neurones

(Before) Linear score function: $f = Wx$

Réseaux de neurones

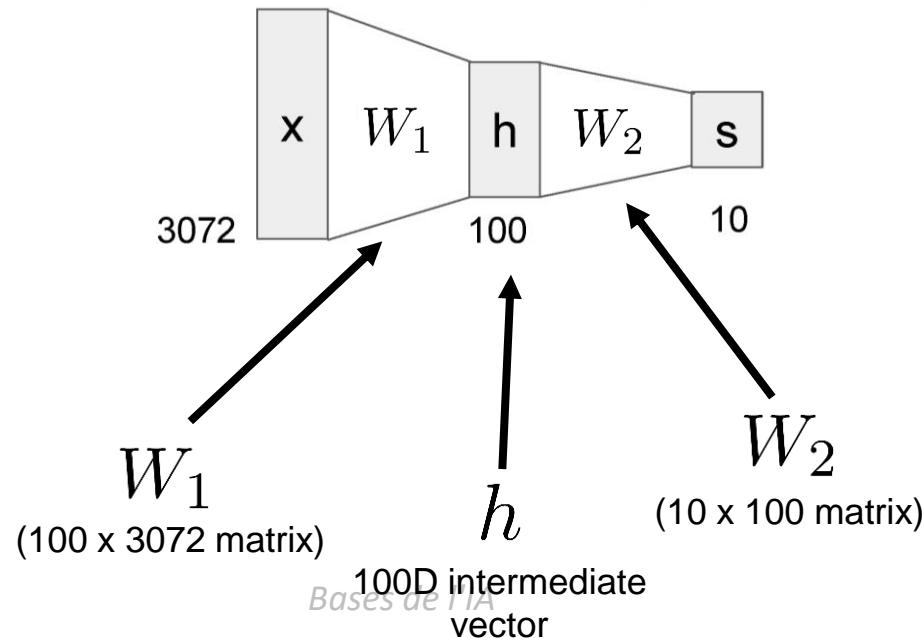
(Before) Linear score function: $f = Wx$

(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$

Réseaux de neurones

(Before) Linear score function: $f = Wx$

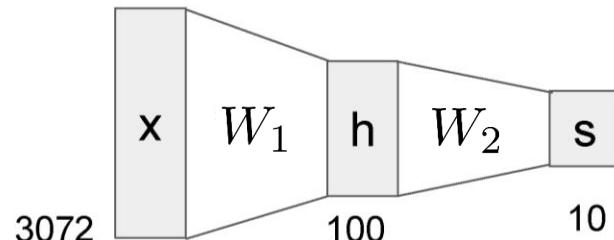
(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$



Réseaux de neurones

(Before) Linear score function: $f = Wx$

(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$



Nombre total de poids à apprendre :

$$3\,072 \times 100 + 100 \times 10 = 308\,200$$

Réseaux de neurones

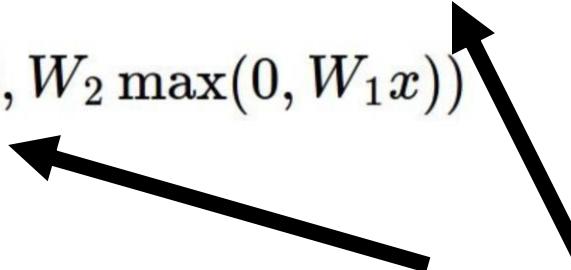
(Before) Linear score function:

$$f = Wx$$

(Now) 2-layer Neural Network
or 3-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$



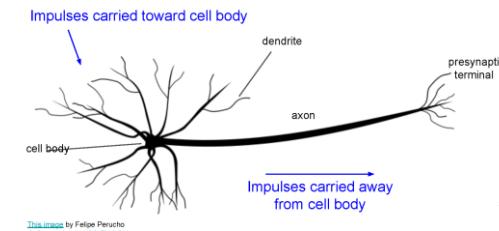
également appelés
"perceptrons
multicouches" (MLP)

Réseaux de neurones

- Généralisation des réseaux neuronaux :
- Fonctions linéaires enchaînées et séparées par des fonctions non linéaires (fonctions d'activation), par exemple "max".

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

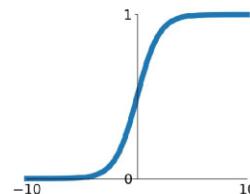
- Pourquoi séparer les fonctions linéaires par des fonctions non linéaires ?
- Très vaguement inspiré par les vrais neurones



Activation functions

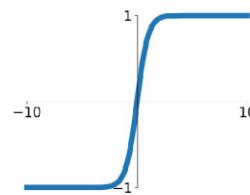
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



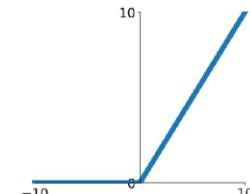
tanh

$$\tanh(x)$$



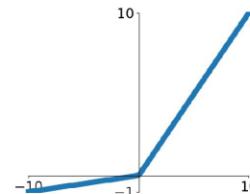
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

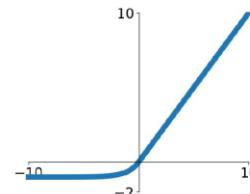


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

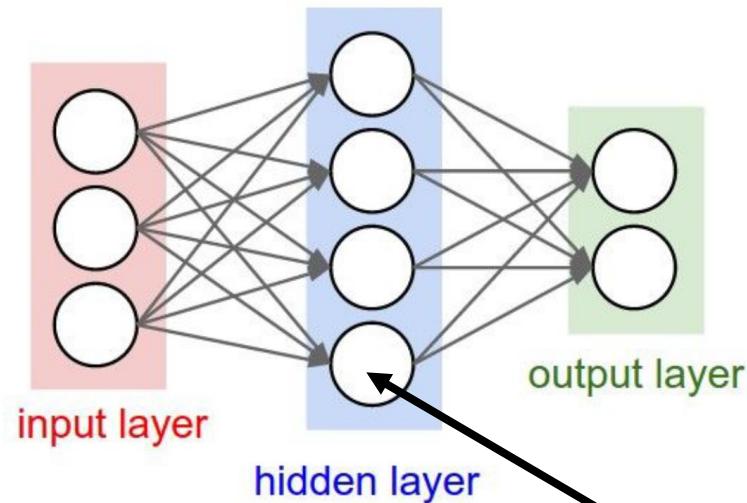
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



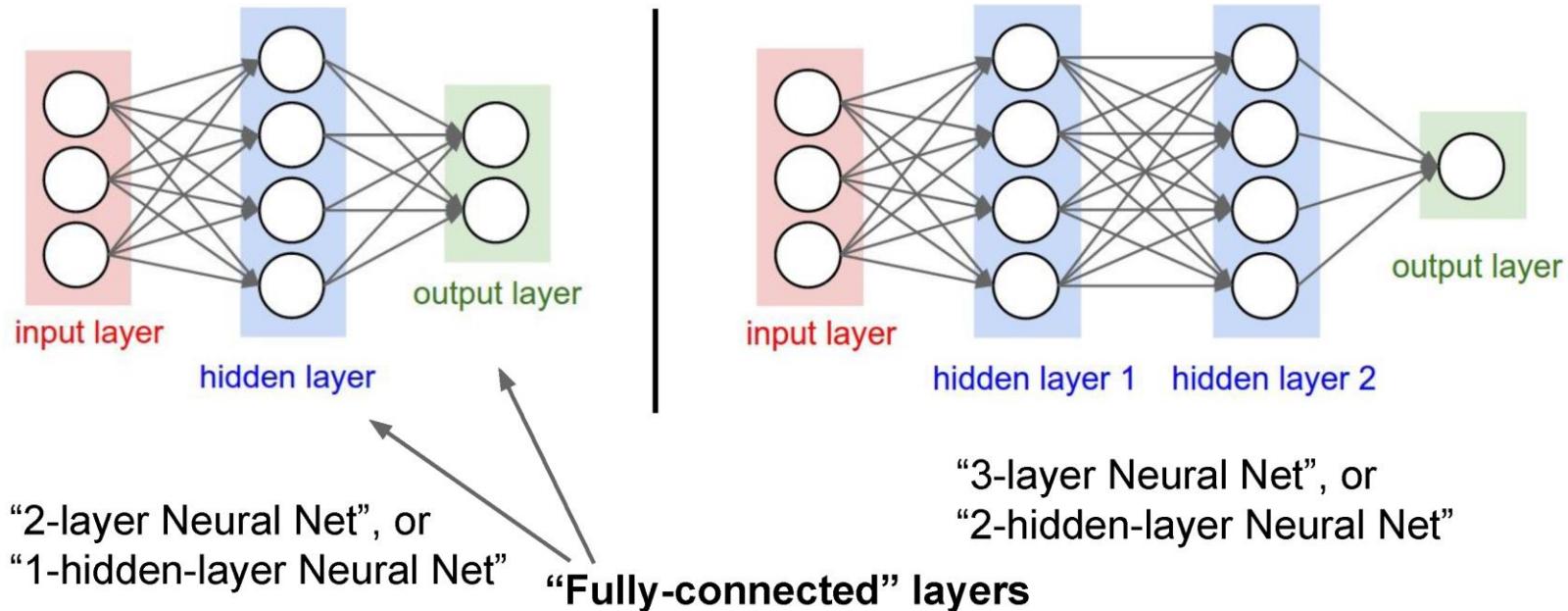
Architecture du réseau de neurones

Réseau de neurones à 2 couches



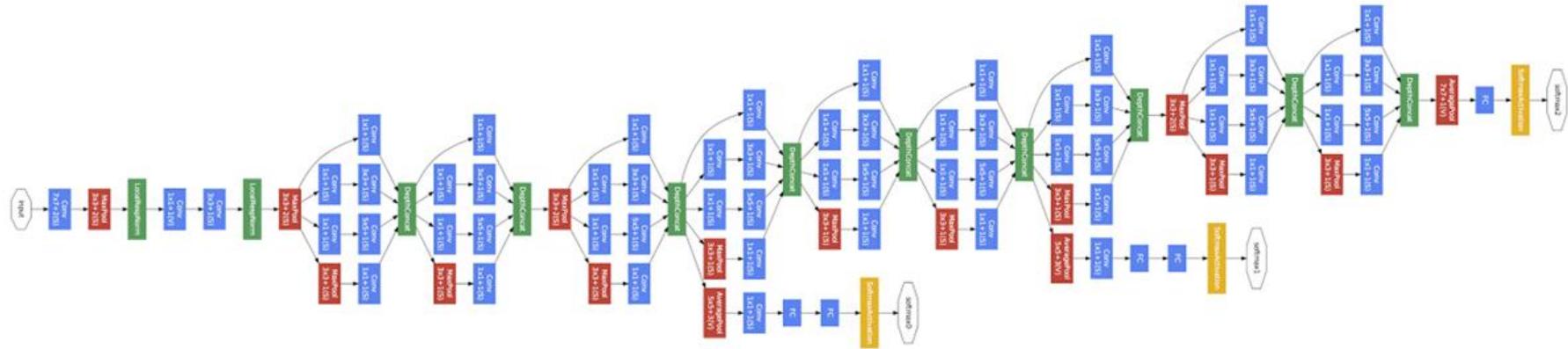
Neurone ou unité

Architecture du réseau de neurones



Les réseaux profonds comportent généralement de nombreuses couches et potentiellement des millions de paramètres

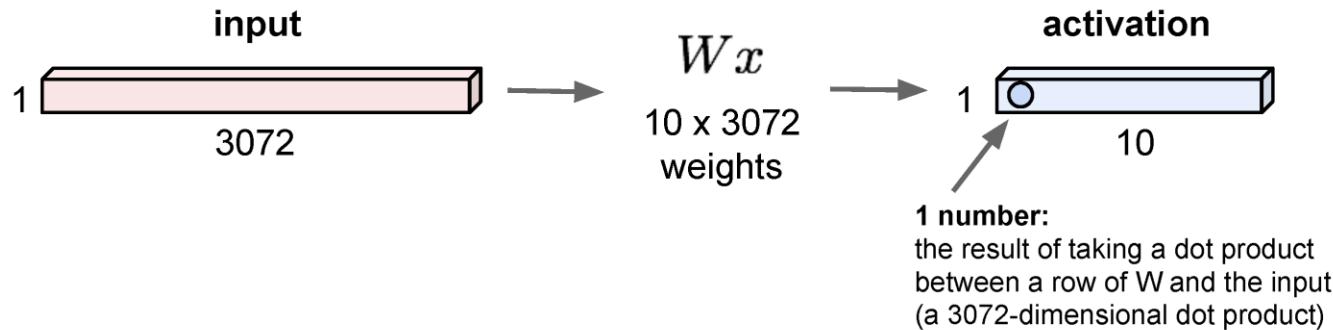
Réseau neuronal profond



- *Inception network* (Szegedy et al, 2015)
- 22 couches

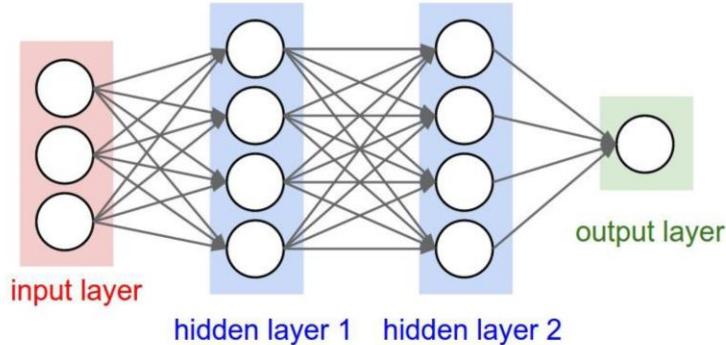
Couches entièrement connectées

32x32x3 image -> stretch to 3072 x 1



Tout comme un classificateur linéaire, mais dans ce cas, il s'agit d'une seule couche d'un réseau plus grand.

Exemple de calcul feed-forward d'un réseau neuronal



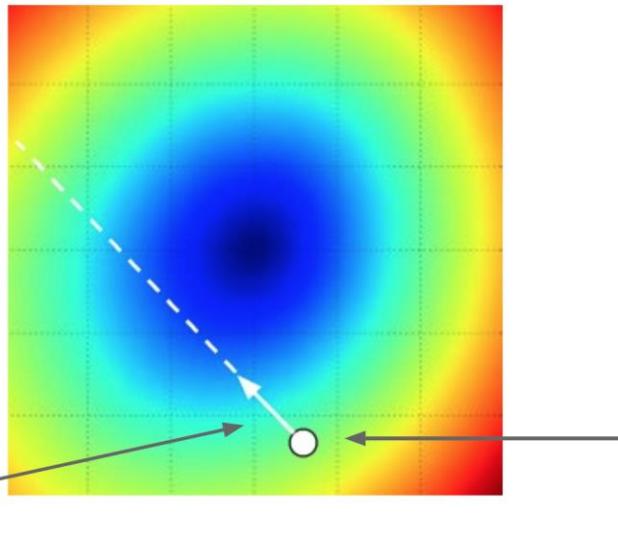
```
# forward-pass of a 3-layer neural network:  
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)  
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)  
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)  
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)  
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

Résumé

- Nous organisons les neurones en **couches entièrement connectées**
- L'abstraction d'une couche a l'avantage de nous permettre d'utiliser un **code vectoriel efficace** (par exemple, les multiplications de matrices).

Optimisation des paramètres avec la descente de gradient

- Comment trouver les meilleurs paramètres W et b ?
- En général : descente de gradient
 - Commencer par une estimation d'un bon W et d'un bon b (ou les initialiser de manière aléatoire)
 - Calculer la fonction de perte pour cette hypothèse initiale et le gradient de la fonction de perte
 - Progresser sur une certaine distance dans la direction du gradient négatif (direction de la descente la plus abrupte)
 - Répéter les étapes 2 et 3
- Remarque : l'exécution efficace de l'étape 2 pour les réseaux profonds est appelée **rétropropagation**.



negative gradient direction

Descente de gradient : marcher dans la direction opposée au gradient

Q : Jusqu'où ?

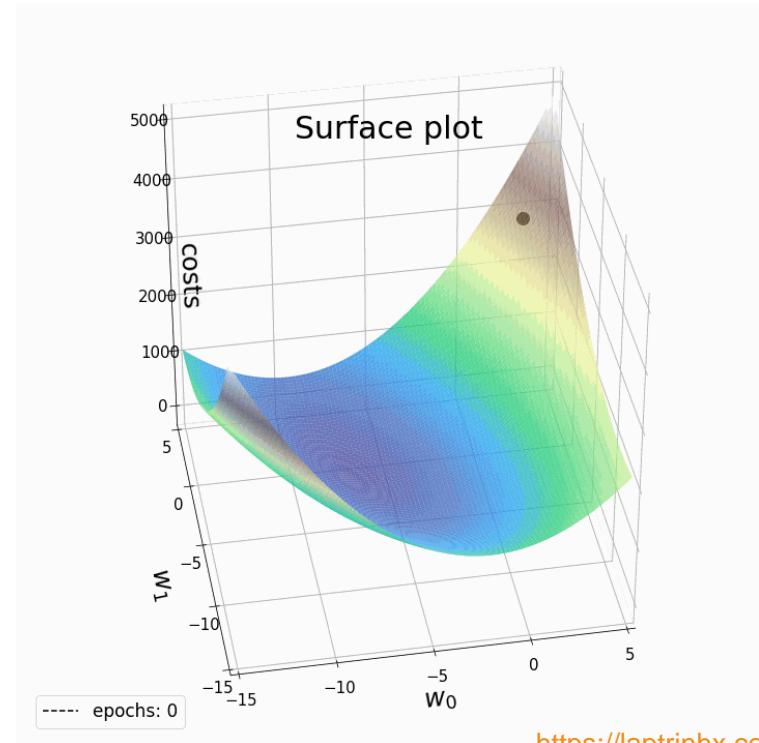
R : Taille du pas : taux d'apprentissage

Trop grand : manque le minimum

Trop petit : convergence lente

Bases de l'IA

Exemple 2D de descente de gradient



En réalité, dans l'apprentissage profond, nous optimisons une fonction de perte très complexe avec des millions de variables (ou plus)

<https://laptrinhx.com/gradient-descent-animation-2-multiple-linear-regression-3070246823/>

Exemple en 2D : TensorFlow

Tinker With a **Neural Network** Right Here in Your Browser.
Don't Worry, You Can't Break It. We Promise.

Epoch 000,000 Learning rate 0.03 Activation Tanh Regularization None Regularization rate 0 Problem type Classification

DATA
Which dataset do you want to use?

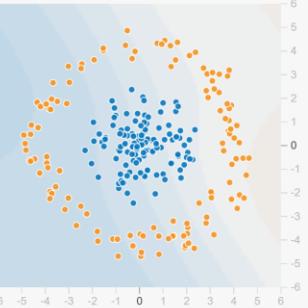


Ratio of training to test data: 50%
Noise: 0
Batch size: 10

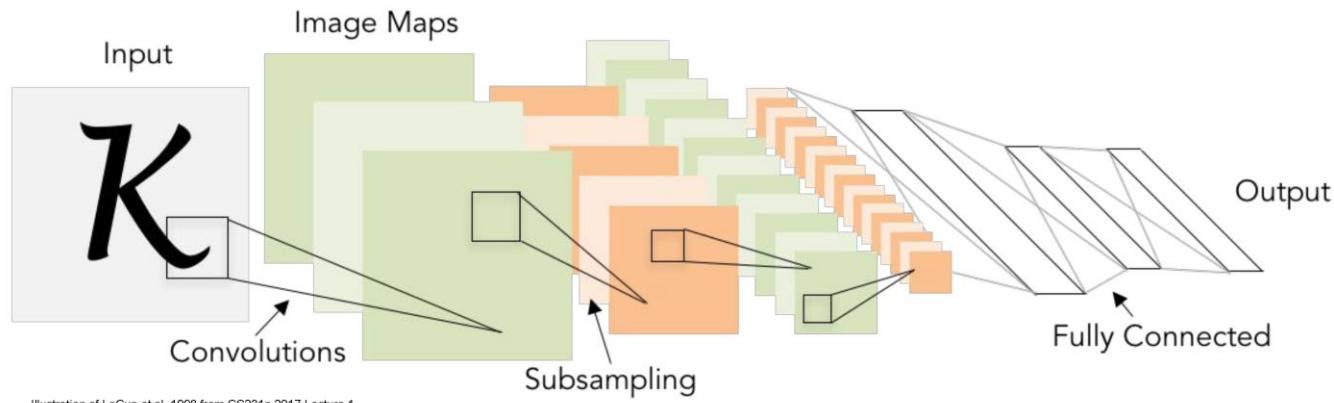
FEATURES
Which properties do you want to feed in?
 X_1 X_2 X_1^2 X_2^2 $X_1 X_2$

2 HIDDEN LAYERS
4 neurons 2 neurons
The outputs are mixed with varying weights, shown by the thickness of the lines.
This is the output from one neuron. Hover to see it larger.

OUTPUT
Test loss 0.505
Training loss 0.502



Réseaux de neurones convolutifs



Les réseaux convolutifs sont omniprésents

Classification



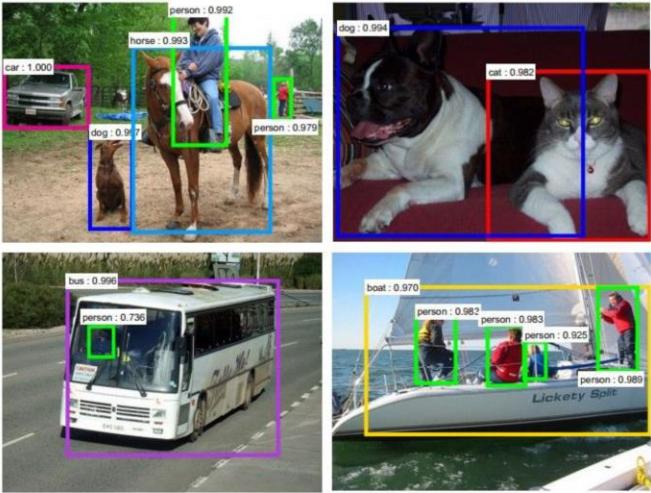
Retrieval



Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

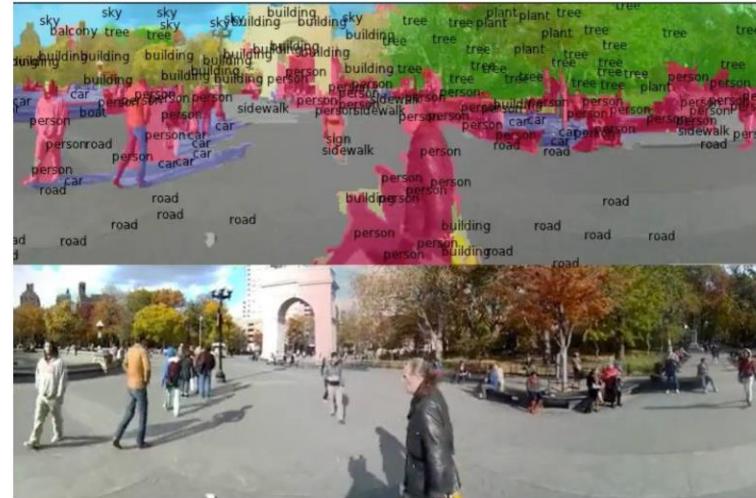
Les réseaux convolutifs sont omniprésents

Detection



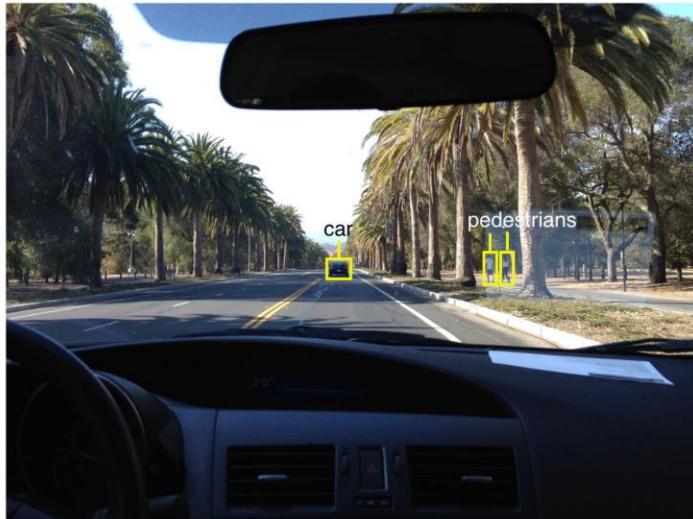
[Faster R-CNN: Ren, He, Girshick, Sun 2015]

Segmentation



[Farabet et al., 2012]

Les réseaux convolutifs sont omniprésents



self-driving cars

Photo by Lane McIntosh. Copyright CS231n 2017.



[This image](#) by GBPublic_PR is
licensed under [CC-BY 2.0](#)

NVIDIA Tesla line

(these are the GPUs on rye01.stanford.edu)

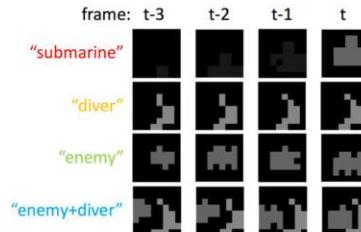
Note that for embedded systems a typical setup would involve NVIDIA Tegras, with integrated GPU and ARM-based CPU cores.

Les réseaux convolutifs sont omniprésents



Images are examples of pose estimation, not actually from Toshev & Szegedy 2014. Copyright Lane McIntosh.

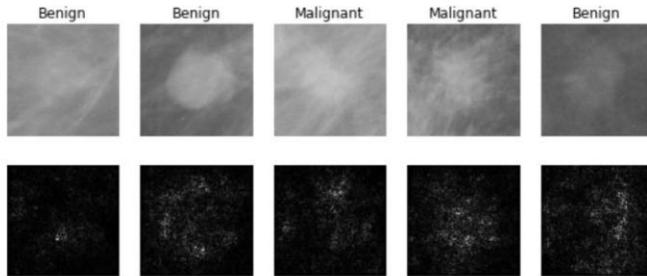
[Toshev, Szegedy 2014]



[Guo et al. 2014]

Figures copyright Xiaoxiao Guo, Satinder Singh, Honglak Lee, Richard Lewis, and Xiaoshi Wang, 2014. Reproduced with permission.

Les réseaux convolutifs sont omniprésents



[Levy et al. 2016]

Figure copyright Levy et al. 2016.
Reproduced with permission.



[Dieleman et al. 2014]

From left to right: [public domain by NASA](#), usage [permitted](#) by
ESA/Hubble, [public domain by NASA](#), and [public domain](#).



[Sermanet et al. 2011]
[Ciresan et al.]

Photos by Lane McIntosh.
Copyright CS231n 2017.

No errors



A white teddy bear sitting in the grass

Minor errors



A man in a baseball uniform throwing a ball

Somewhat related



A woman is holding a cat in her hand

Image Captioning

[Vinyals et al., 2015]
[Karpathy and Fei-Fei, 2015]



A man riding a wave on top of a surfboard



A cat sitting on a suitcase on the floor



A woman standing on a beach holding a surfboard

All images are CC0 Public domain:
<https://pixabay.com/en/luggage-antique-cat-1642010/>
<https://pixabay.com/en/teddy-plush-bears-cute-teddy-bear-1623436/>
<https://pixabay.com/en/surf-wave-summer-sport-literal-1668716/>
<https://pixabay.com/en/woman-female-model-portrait-adult-983967/>
<https://pixabay.com/en/handsstand-lake-meditation-498009/>
<https://pixabay.com/en/baseball-player-shortstop-infield-1045263/>

Captions generated by Justin Johnson using [NeuralTalk2](#)

De la légende au texte

TEXT PROMPT

an illustration of a baby daikon radish in a tutu walking a dog



AI-GENERATED IMAGES



[Edit prompt or view more images↓](#)

TEXT PROMPT

an armchair in the shape of an avocado [...]

AI-GENERATED IMAGES



[Edit prompt or view more images↓](#)

DALL·E: Creating Images from Text,
OpenAI

<https://openai.com/blog/dall-e/>

Bases de l'IA

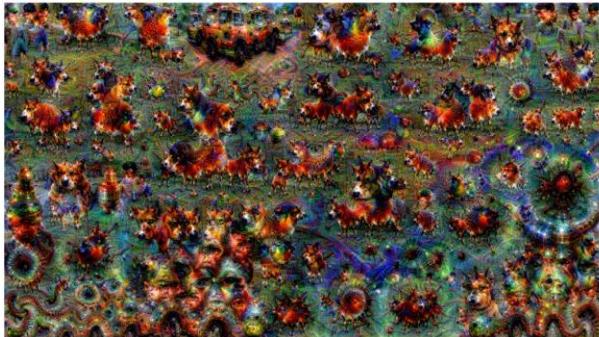
TEXT PROMPT

a store front that has the word 'openai' written on it [...]

AI-GENERATED IMAGES



5/4/2023



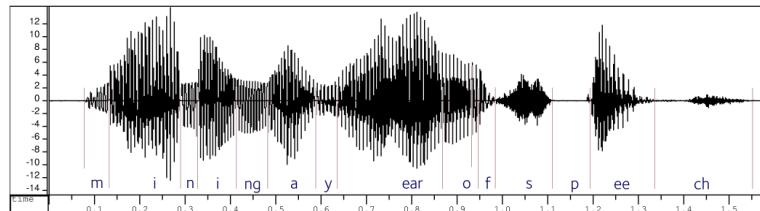
Figures copyright Justin Johnson, 2015. Reproduced with permission. Generated using the Inceptionism approach from a [blue post](#) by Google Research.

Original Image is CC0 public domain
Starry Night and *Tree Roots* by Van Gogh are in the public domain
Bokeh image is in the public domain
 Gatys et al, "Image Style Transfer using Convolutional Neural Networks", CVPR 2016
 Gatys et al, "Controlling Perceptual Factors in Neural Style Transfer", CVPR 2017

Réseaux de neurones convolutifs

Version des réseaux de neurones profonds conçue pour les signaux

- Signaux 1D (par exemple, formes d'ondes vocales)



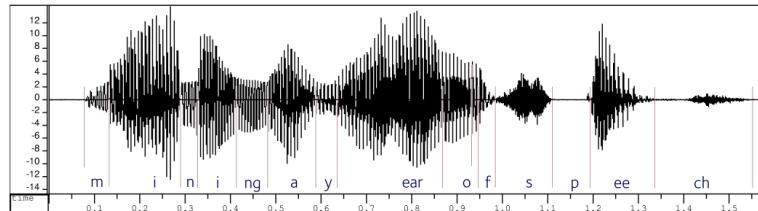
- Signaux 2D (par exemple, images)



Réseaux de neurones convolutifs

Version des réseaux de neurones profonds conçue pour les signaux

- Signaux 1D (par exemple, formes d'ondes vocales)



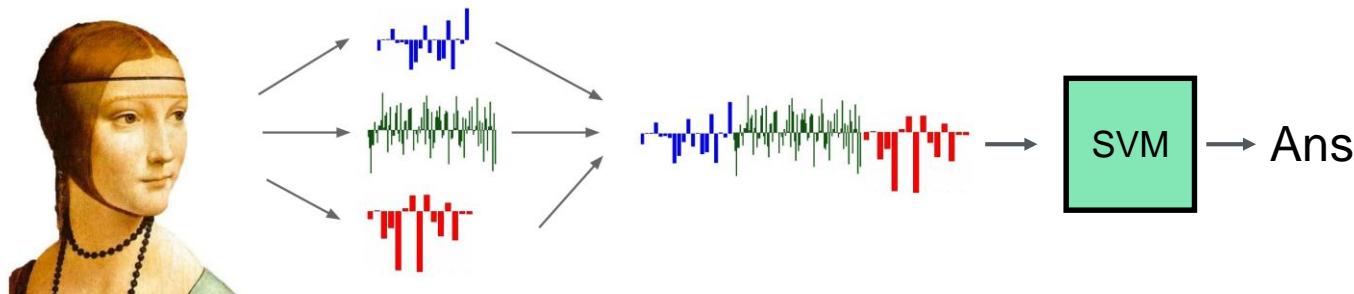
- Signaux 2D (par exemple, images)



Motivation:
Apprentissage des features

Bases de l'IA

La vie avant l'apprentissage profond...



*Pixels
d'entrée*

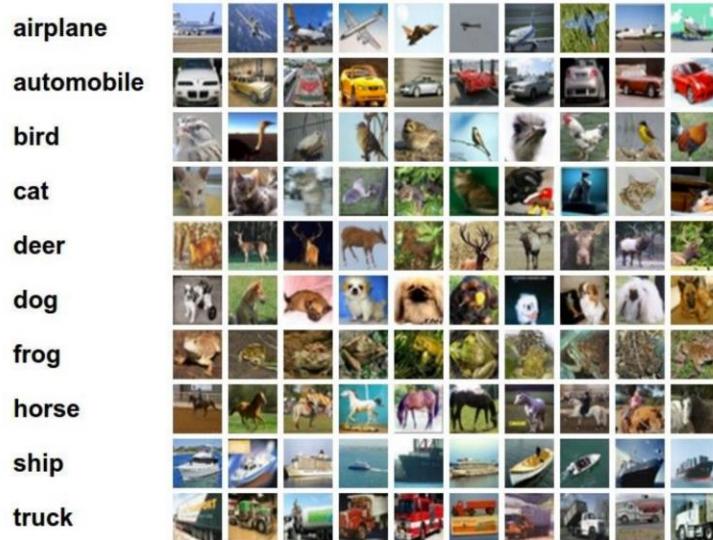
*Extraction de
features*

*Concaténer dans
un vecteur x*

Classificateur linéaire

Pourquoi utiliser des features ?

Pourquoi pas des pixels ?

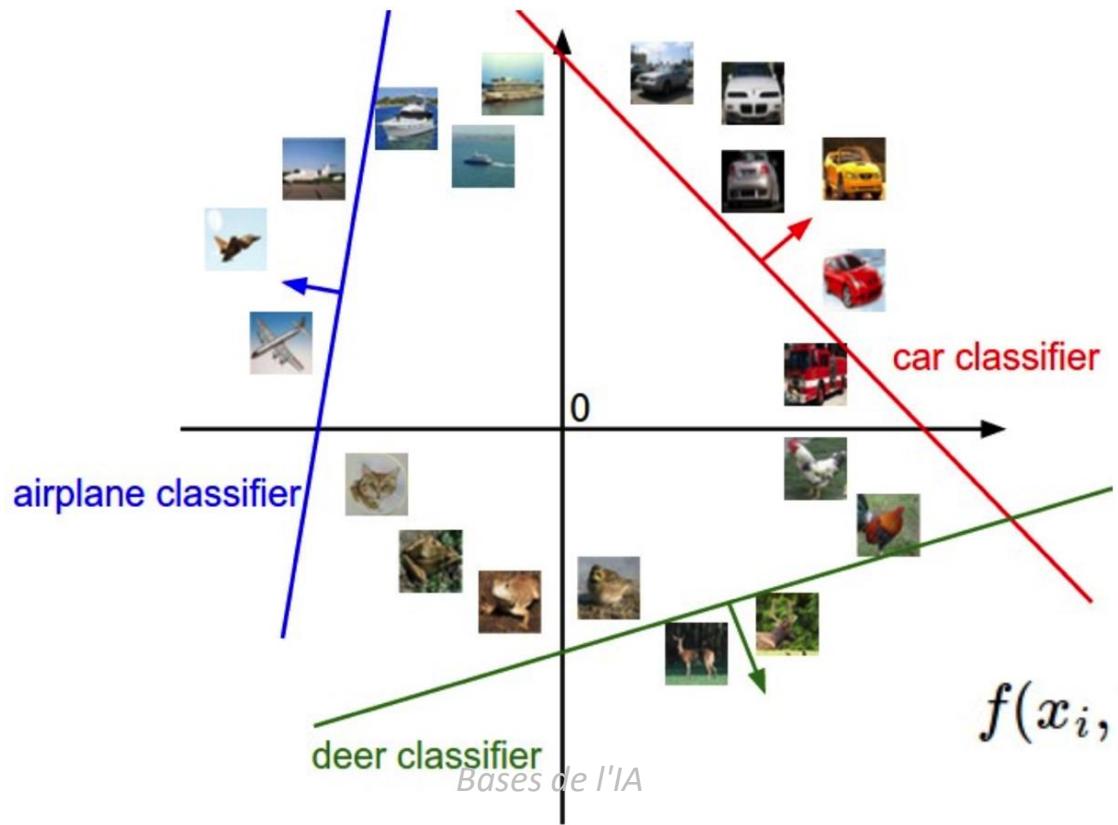


$$f(x_i, W, b) = Wx_i + b$$

Q : Quel serait un ensemble de classes très difficile à distinguer pour un classificateur linéaire ?

(en supposant que $x = \text{pixels}$)

Classes séparables linéairement



$$f(x_i, W, b) = Wx_i + b$$

Features de l'image

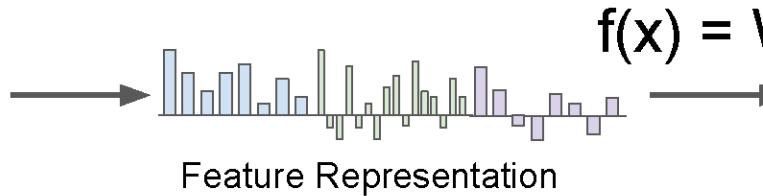


$$f(x) = Wx$$

Class
scores



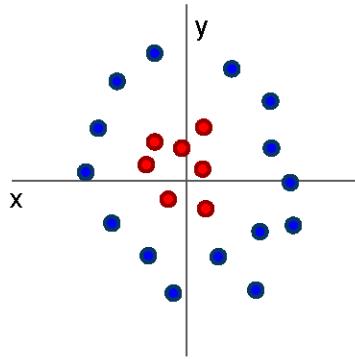
Features de l'image



$$f(x) = Wx$$

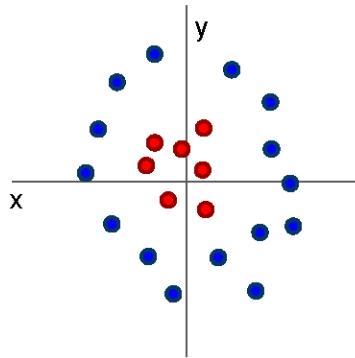
Class
scores

Features de l'image: motivation



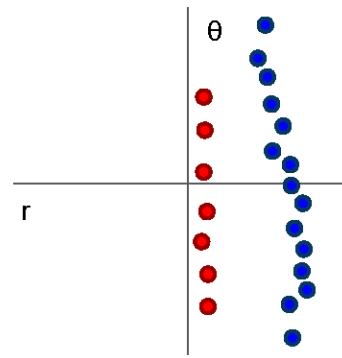
ne peut pas séparer les
points rouges et bleus
avec un classificateur
linéaire

Features de l'image: motivation



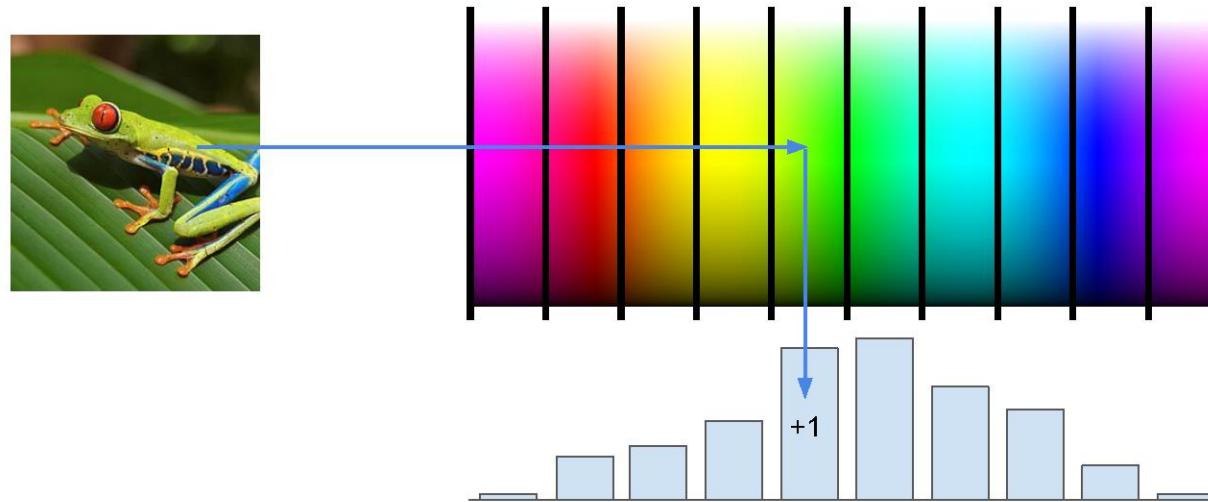
ne peut pas séparer les points rouges et bleus avec un classificateur linéaire

$$f(x, y) = (r(x, y), \theta(x, y))$$



après application de la transformation des features, les points peuvent être séparés par un classificateur linéaire

Exemple : histogramme des couleurs

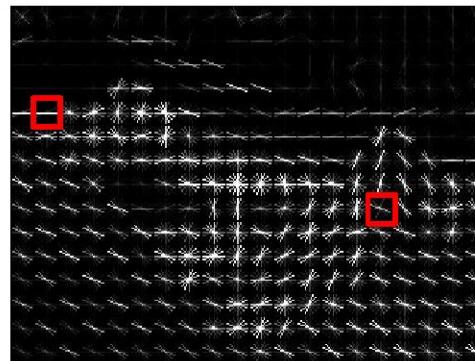


Exemple : histogramme des gradients orientés



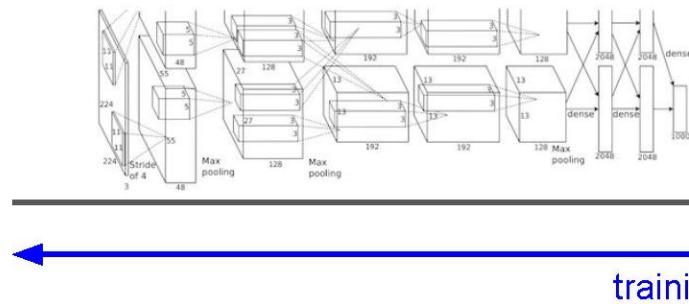
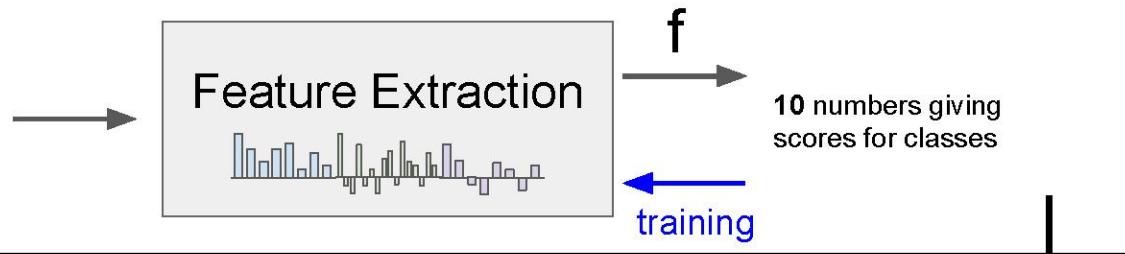
Diviser l'image en régions de 8x8 pixels
Dans chaque région, quantifier la direction des bords en 9 cases

Lowe, "Object recognition from local scale-invariant features", ICCV 1998
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005



Exemple : L'image 320x240 est divisée en 40x30 cases ; dans chaque case, il y a 9 nombres ; le vecteur de caractéristiques comprend donc $30 \times 40 \times 9 = 10800$ nombres.

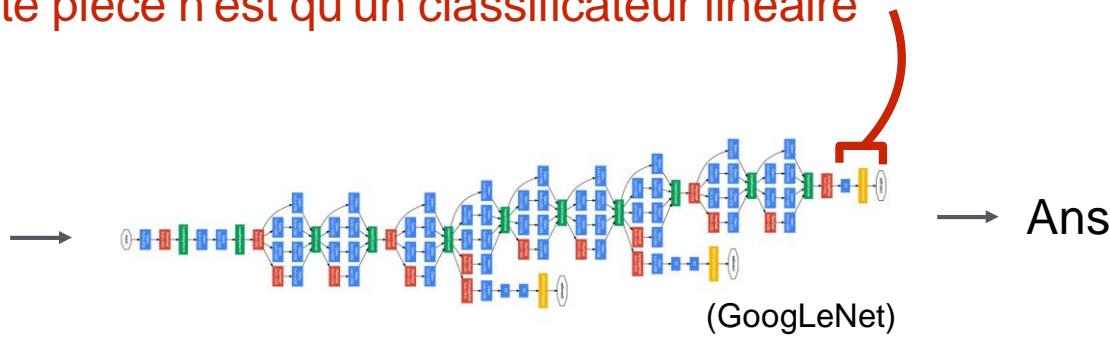
Features de l'image vs ConvNets



La dernière couche de la plupart des CNN est un classificateur linéaire



Cette pièce n'est qu'un classificateur linéaire



Pixels
d'entrée

Tout exécuter à l'aide d'un grand réseau de neurones, entraîné de bout en bout

Clé : effectuer suffisamment de traitements pour que les classes soient linéairement séparables lorsque l'on arrive à la fin du réseau.

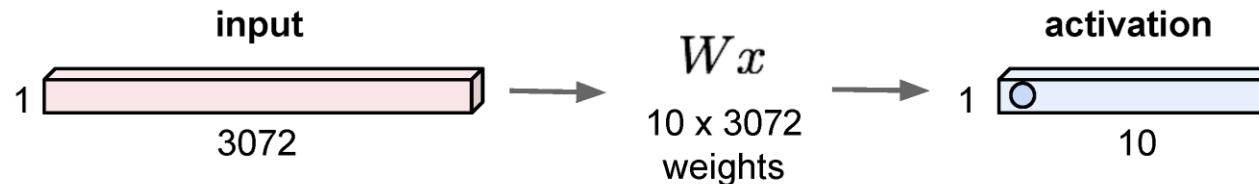
Réseaux de neurones convolutifs

Types de couches :

- Couche entièrement connectée
- Couche convulsive
- Pooling layer

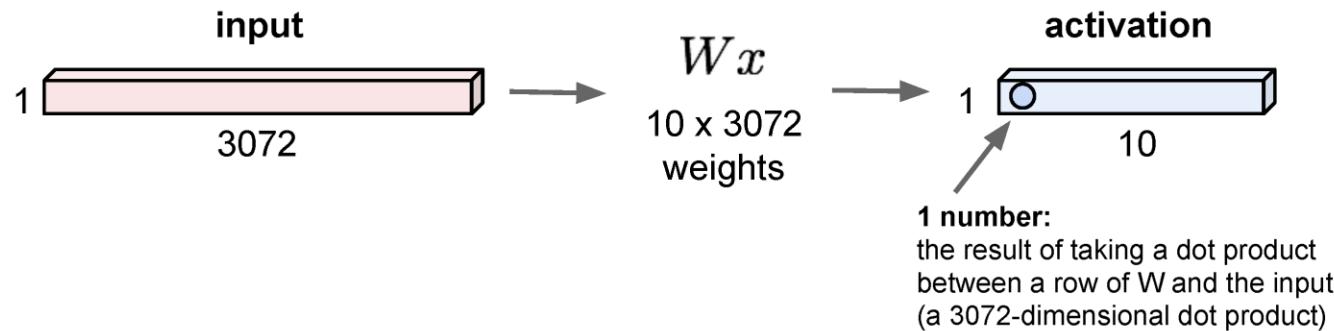
Couche entièrement connectée

32x32x3 image -> stretch to 3072 x 1



Couche entièrement connectée

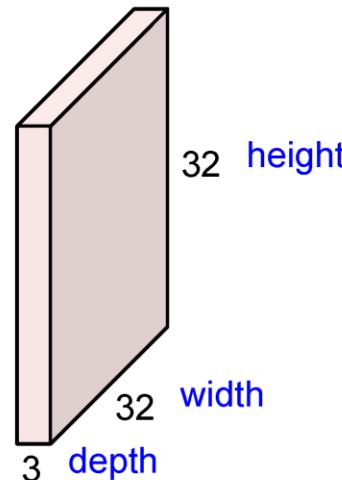
32x32x3 image -> stretch to 3072 x 1



Identique à un classificateur linéaire!

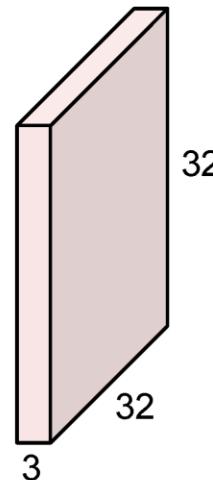
Couche convulsive

32x32x3 image -> preserve spatial structure



Couche convulsive

32x32x3 image

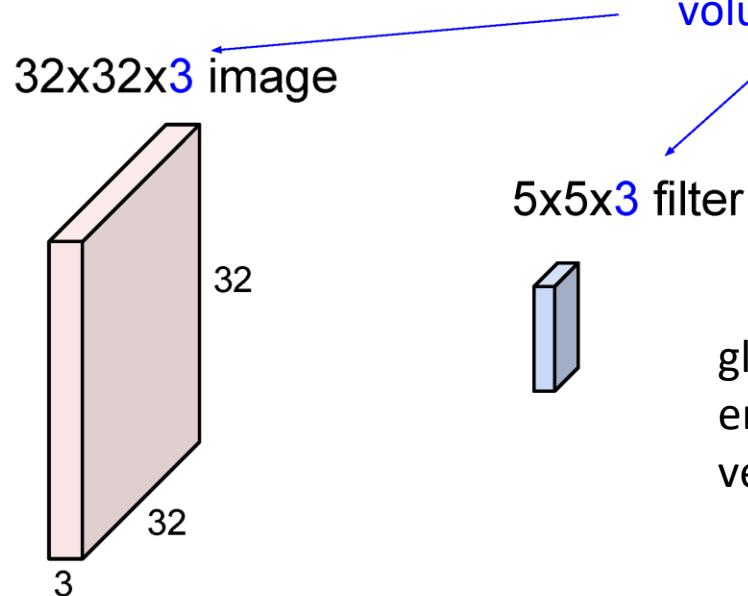


5x5x3 filter



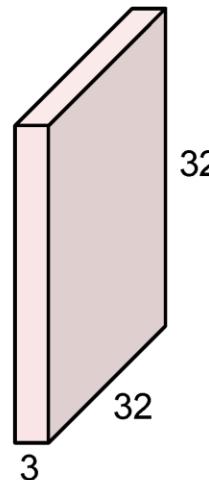
glisser le filtre sur l'image dans l'espace,
en calculant les produits scalaire de deux
vecteurs

Couche convulsive



Couche convulsive

32x32x3 image



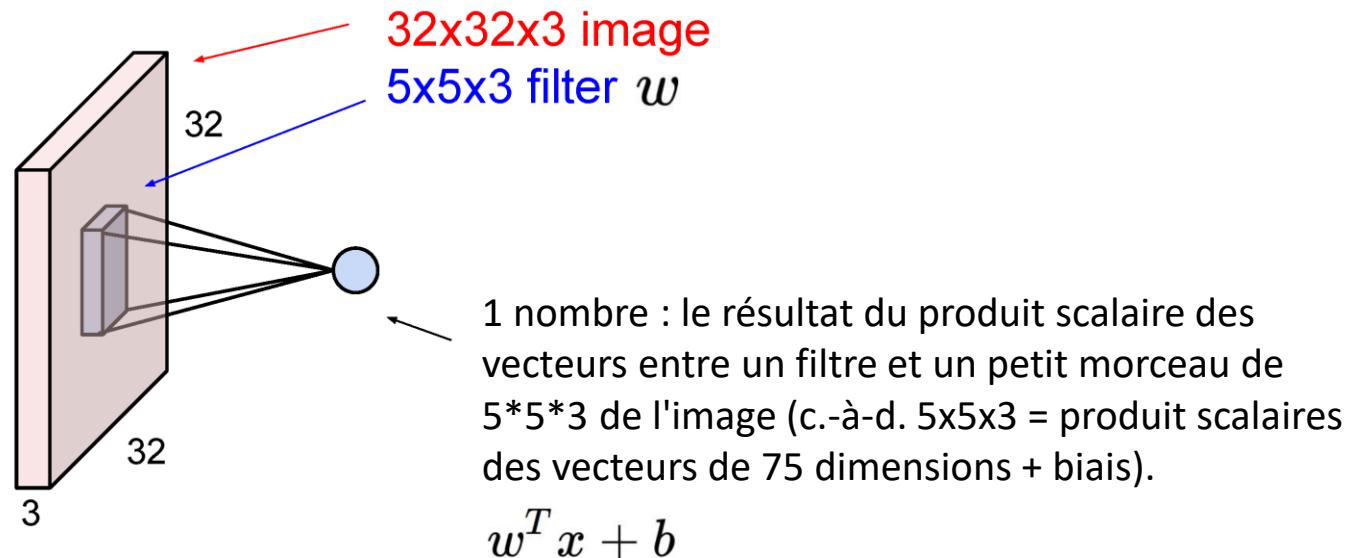
5x5x3 filter



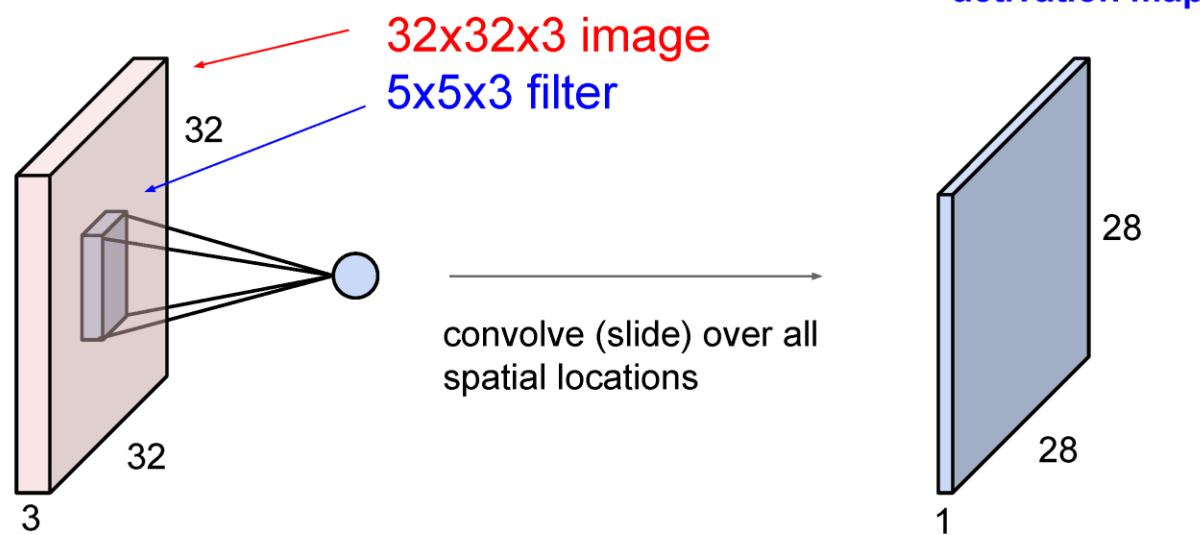
glisser le filtre sur l'image dans l'espace,
en calculant les produits scalaire de deux
vecteurs

Nombre de poids : $5 \times 5 \times 3 + 1 =$
76 (contre 3072 pour une couche
entièrement connectée) (+1 pour le biais)

Couche convulsive



Couche convulsive

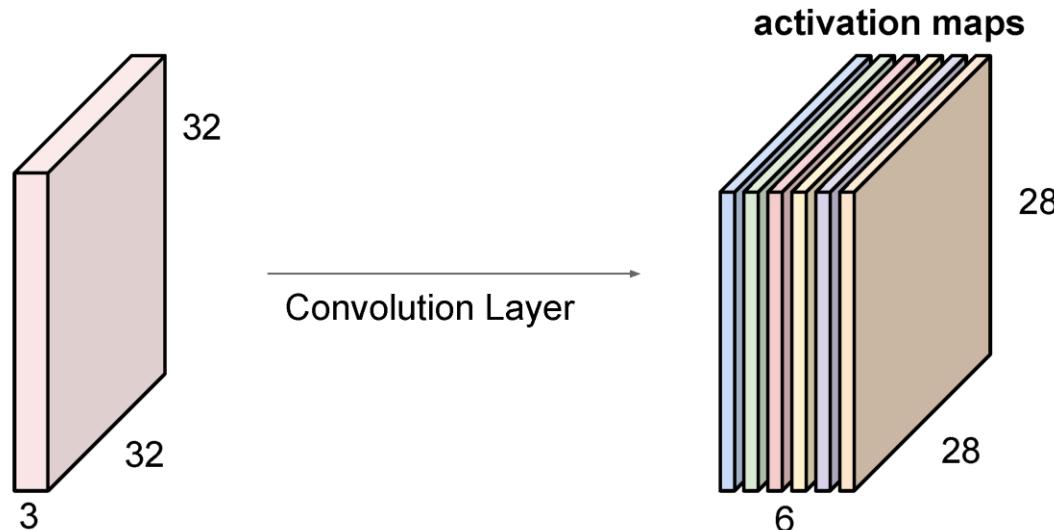


Couche convulsive

consider a second, green filter



For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

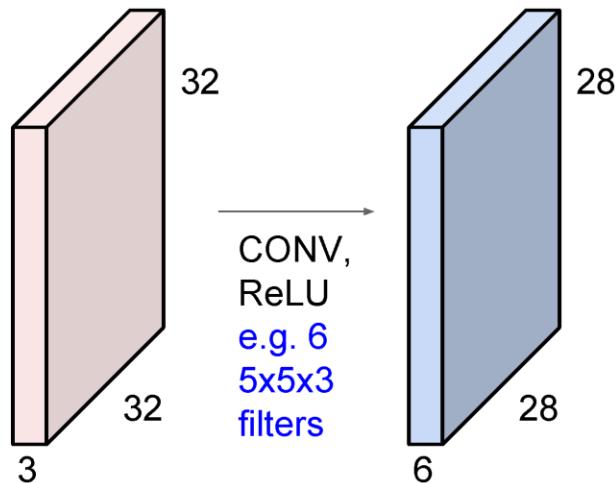


We stack these up to get a “new image” of size 28x28x6!

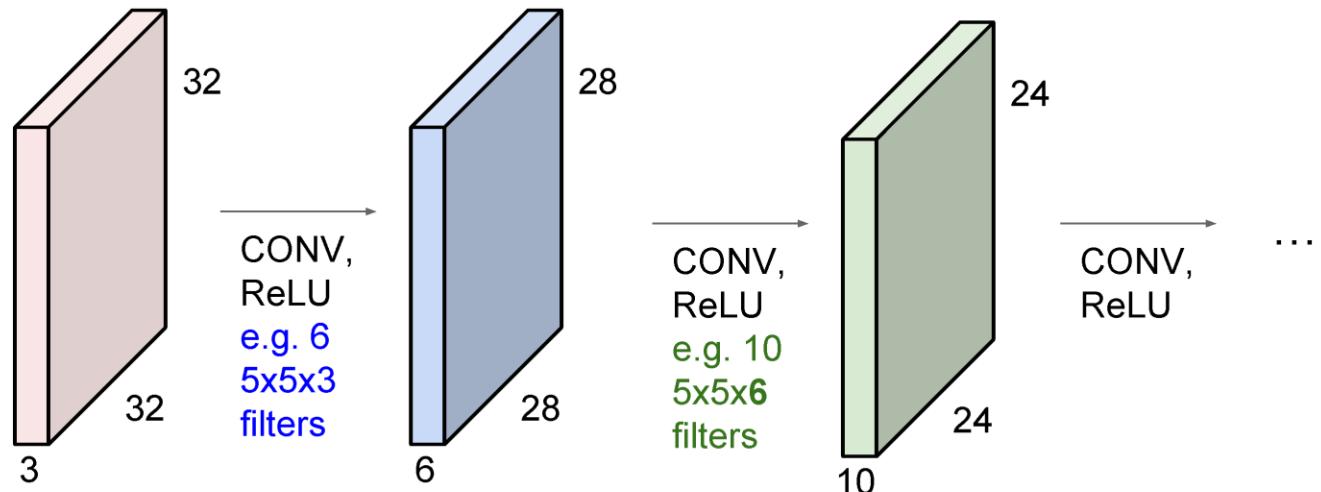
(nombre total de paramètres : $6 \times (75 + 1) = 456$)

Bases de l'IA

ConvNet est une séquence de couches de convolution, entrecoupées de fonctions d'activation.



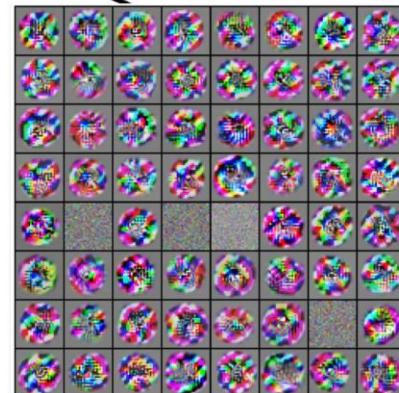
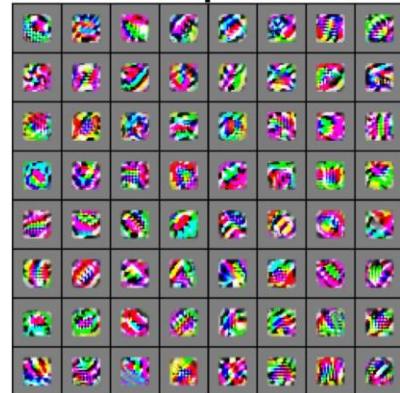
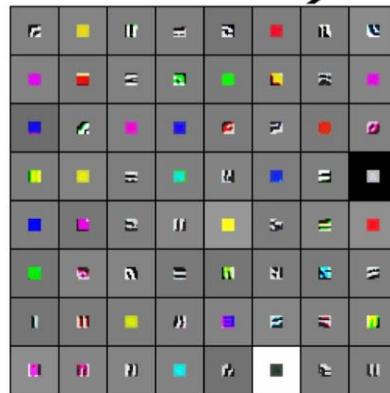
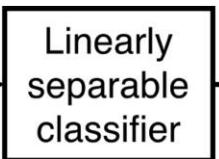
ConvNet est une séquence de couches de convolution, entrecoupées de fonctions d'activation.



Preview

[Zeiler and Fergus 2013]

Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].



Propriétés de la couche convolutive

- Nombre réduit de paramètres à apprendre par rapport à une couche entièrement connectée
- Préserve la structure spatiale - la sortie d'une couche convolutive a la forme d'une image.

[ConvNetJS demo: training on CIFAR-10]

ConvNetJS CIFAR-10 demo

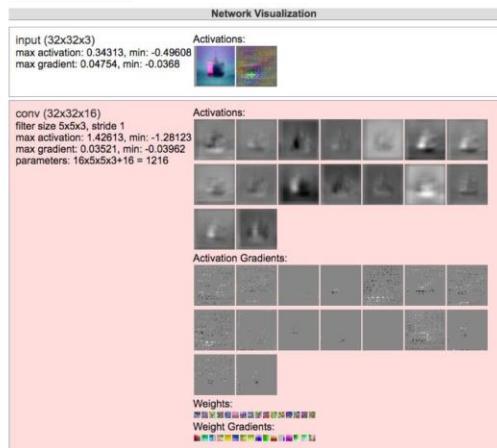
Description

This demo trains a Convolutional Neural Network on the [CIFAR-10 dataset](#) in your browser, with nothing but Javascript. The state of the art on this dataset is about 90% accuracy and human performance is at about 94% (not perfect as the dataset can be a bit ambiguous). I used [this python script](#) to parse the [original files](#) (python version) into batches of images that can be easily loaded into page DOM with img tags.

This dataset is more difficult and it takes longer to train a network. Data augmentation includes random flipping and random image shifts by up to 2px horizontally and vertically.

By default, in this demo we're using Adadelta which is one of per-parameter adaptive step size methods, so we don't have to worry about changing learning rates or momentum over time. However, I still included the text fields for changing these if you'd like to play around with SGD+Momentum trainer.

Report questions/bugs/suggestions to [@karpathy](#).



<https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

Vue d'ensemble

- Un réseau neuronal convolutif peut être considéré comme une fonction allant des images aux scores des classes
- Avec des millions de poids réglables...
- ... conduisant à un mappage très non linéaire entre les images et les caractéristiques / scores de classe.
- Nous réglerons ces poids en fonction de la précision de la classification sur les données d'entraînement...
- ... et nous espérons que notre réseau se généralisera à de nouvelles images au moment du test.

Vue d'ensemble

- Un réseau neuronal fonction allant des i
- Avec des millions de
- ... conduisant à un r les caractéristiques
- Nous réglerons ces classification sur les
- ... et nous espérons nouvelles images au

