

Bases de données multimédia

IX – Machine à vecteurs de support et CNN

ENSIMAG

2014-2015

Matthijs Douze & Karteek Alahari

Inria



Plan

- Classification
- Machine à vecteurs de support
- Reconnaissance d'actions
- CNN
- Quelques perspectives

Inria



Classification

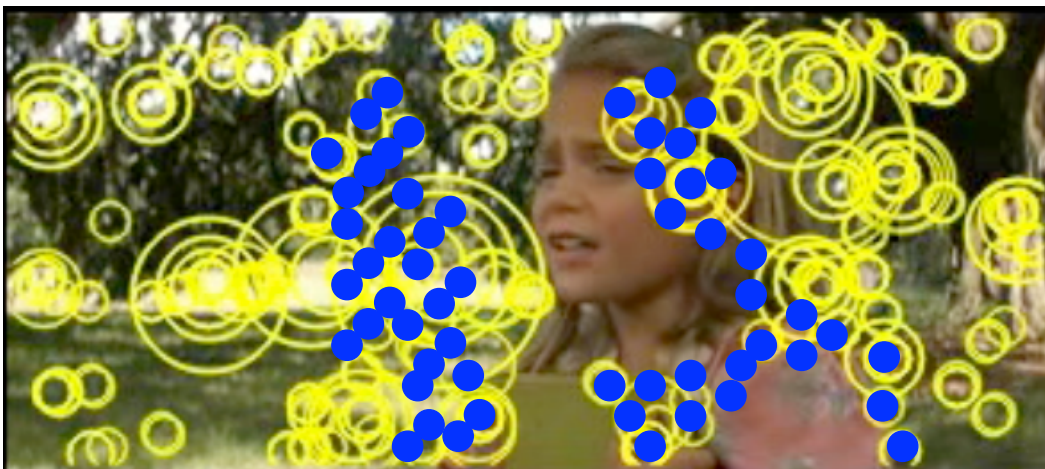
- What can we do, given all these features ?



Inria

Classification

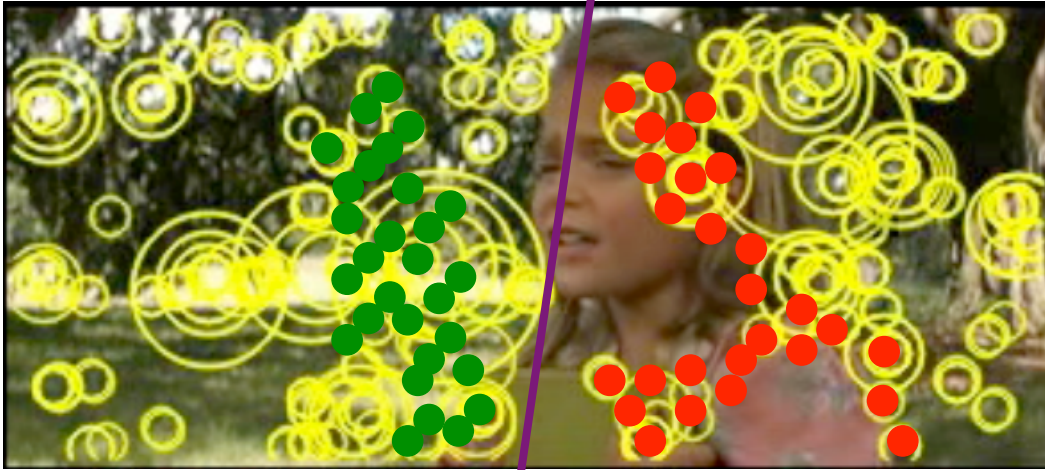
- What can we do, given all these features ?



Inria

Classification

- What can we do, given all these features ?

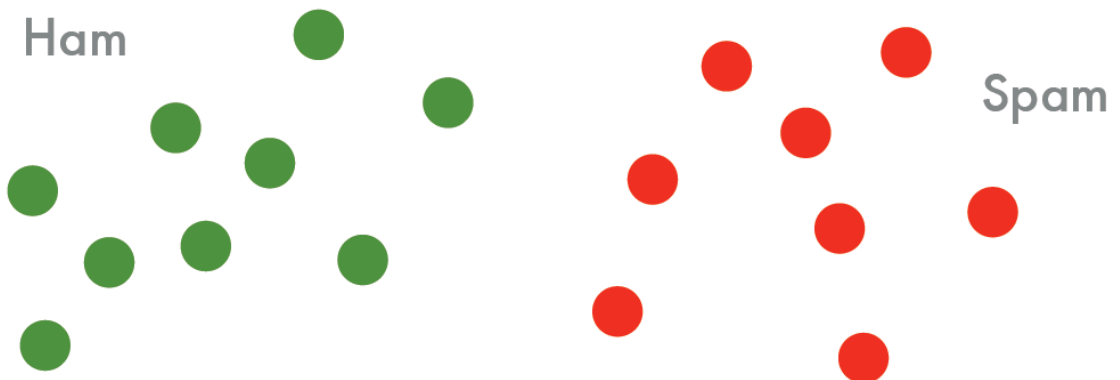


Inria

?

Grenoble INP
ensimag

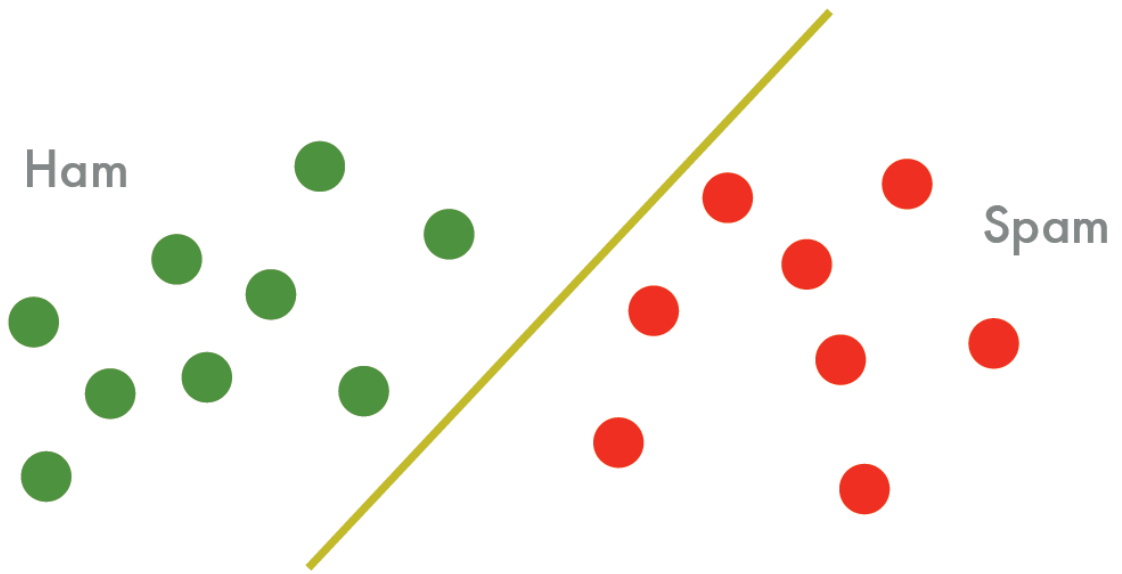
Classification



Inria

Grenoble INP
ensimag
Slides courtesy Alex Smola

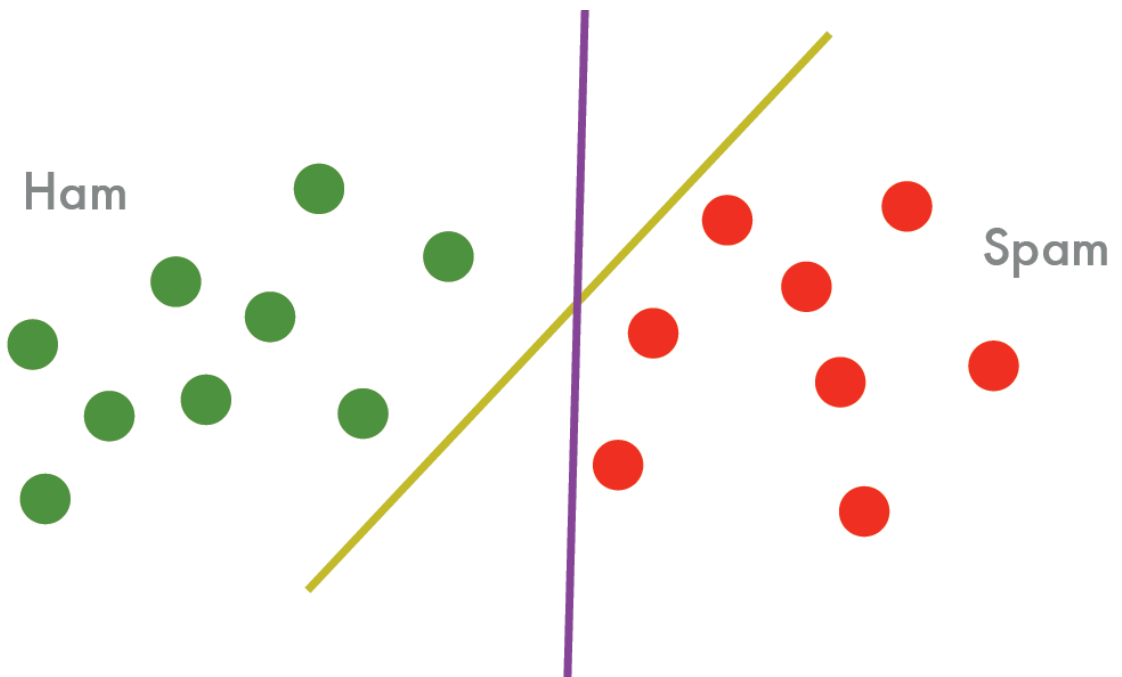
Classification



Inria



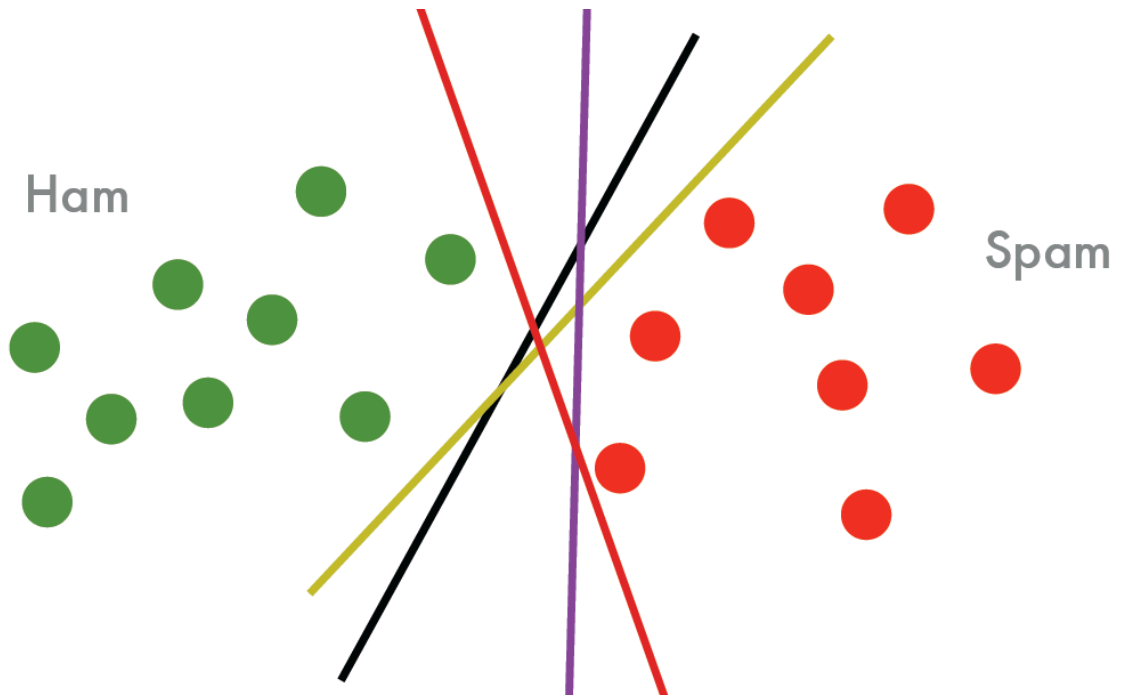
Classification



Inria



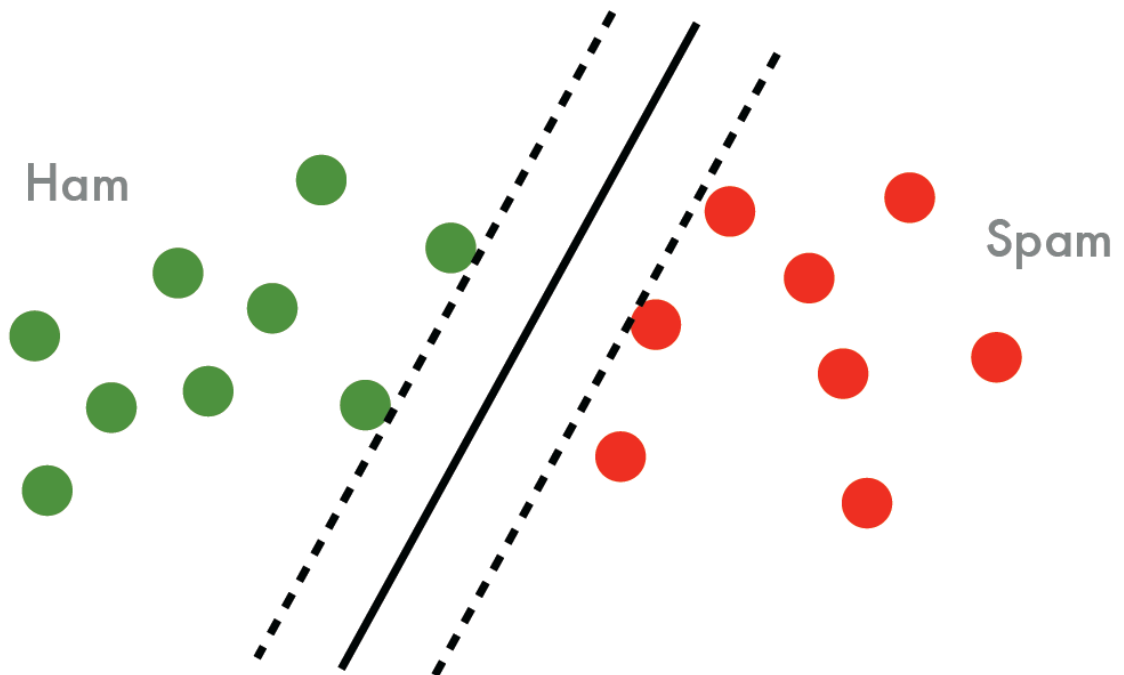
Classification



Inria

Grenoble INP
ensimag

Classification

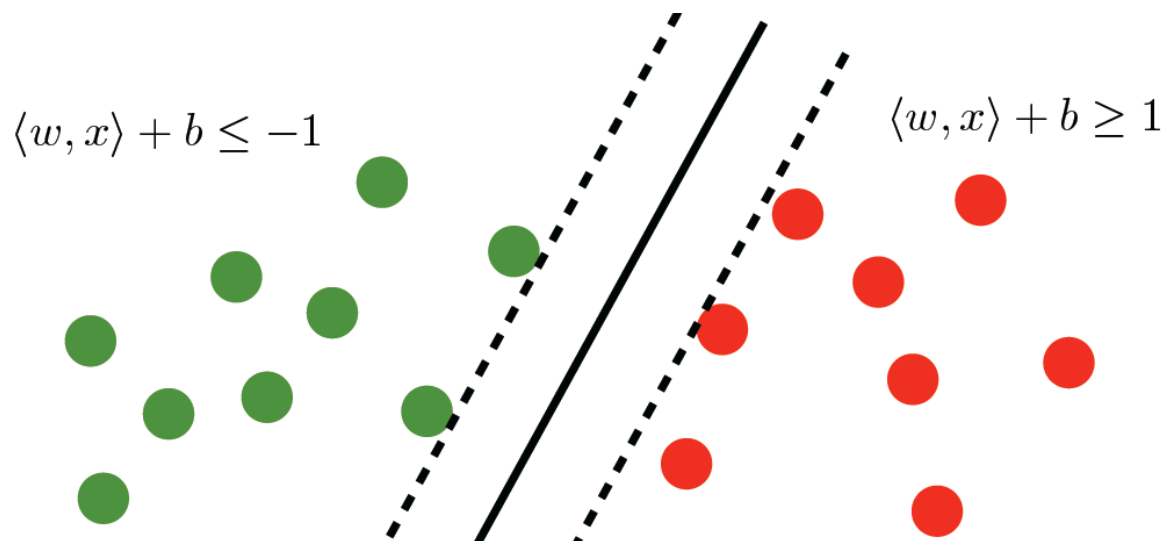


Large margin classifier

Inria

Grenoble INP
ensimag

Classification

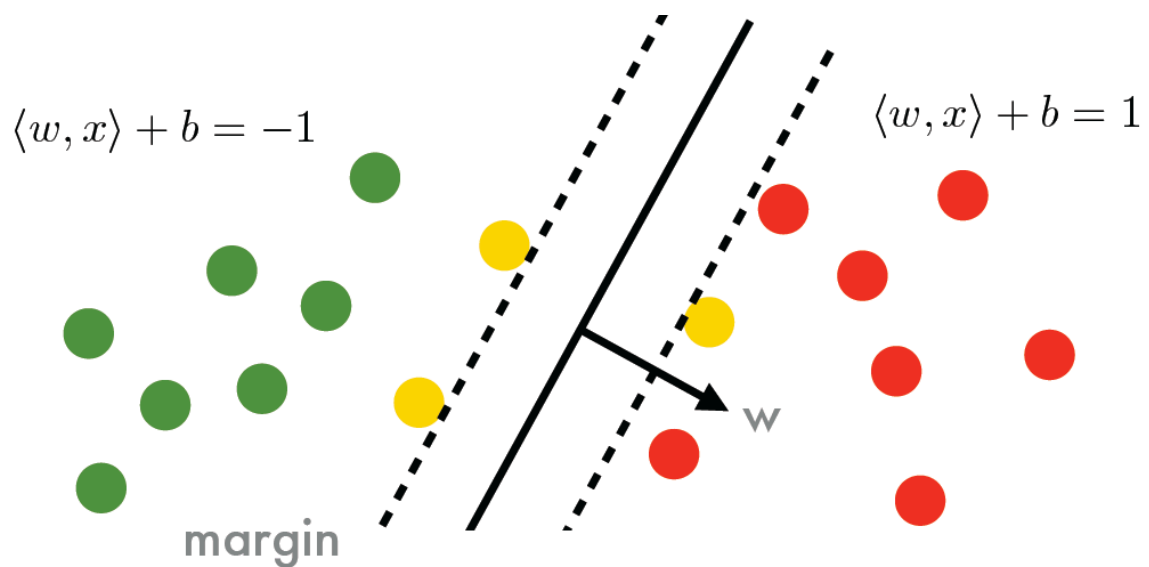


linear function
 $f(x) = \langle w, x \rangle + b$

Inria



Classification

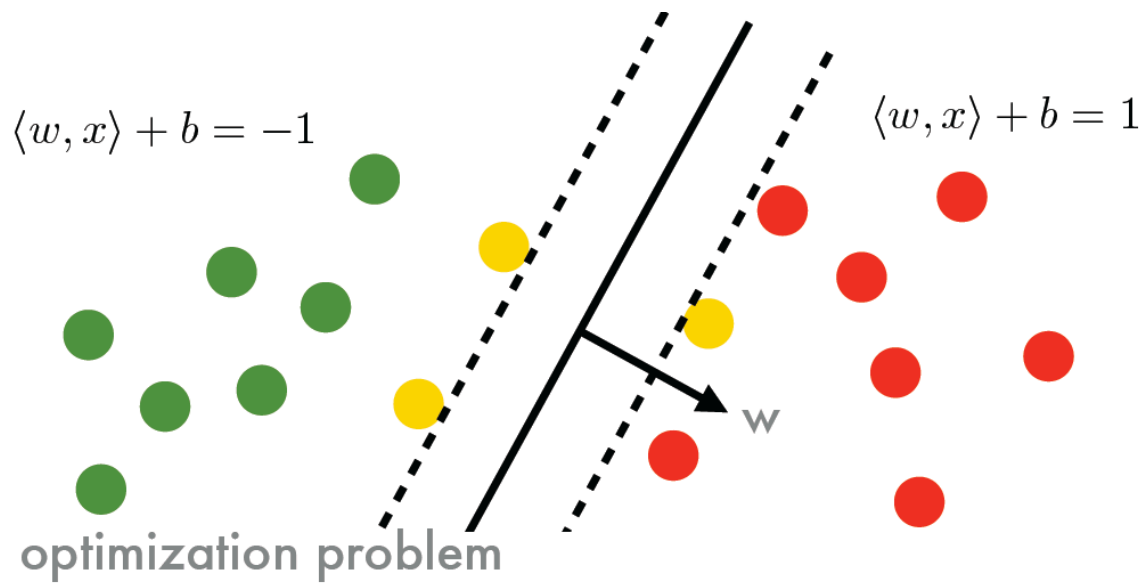


$$\frac{\langle x_+ - x_-, w \rangle}{2 \|w\|} = \frac{1}{2 \|w\|} [[\langle x_+, w \rangle + b] - [\langle x_-, w \rangle + b]] = \frac{1}{\|w\|}$$

Inria



Classification



$$\text{maximize}_{w,b} \frac{1}{\|w\|} \text{ subject to } y_i [\langle x_i, w \rangle + b] \geq 1$$

Inria



Support Vector Machines

$$\text{minimize}_{w,b} \frac{1}{2} \|w\|^2 \text{ subject to } y_i [\langle x_i, w \rangle + b] \geq 1$$

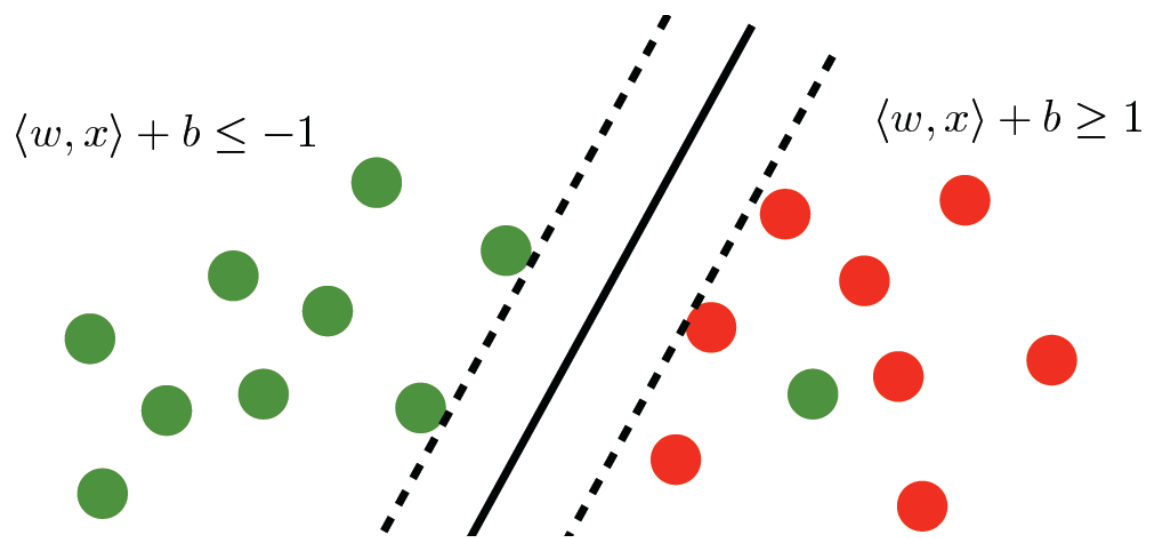
- Many optimization techniques to solve this
 - e.g., Stochastic Gradient Descent (SGD)
- Implementations available
 - SVM^{light} (Thorsten Joachims)
 - SGD-SVM (Léon Bottou)

Inria



Support Vector Machines

- What about linearly inseparable cases ?



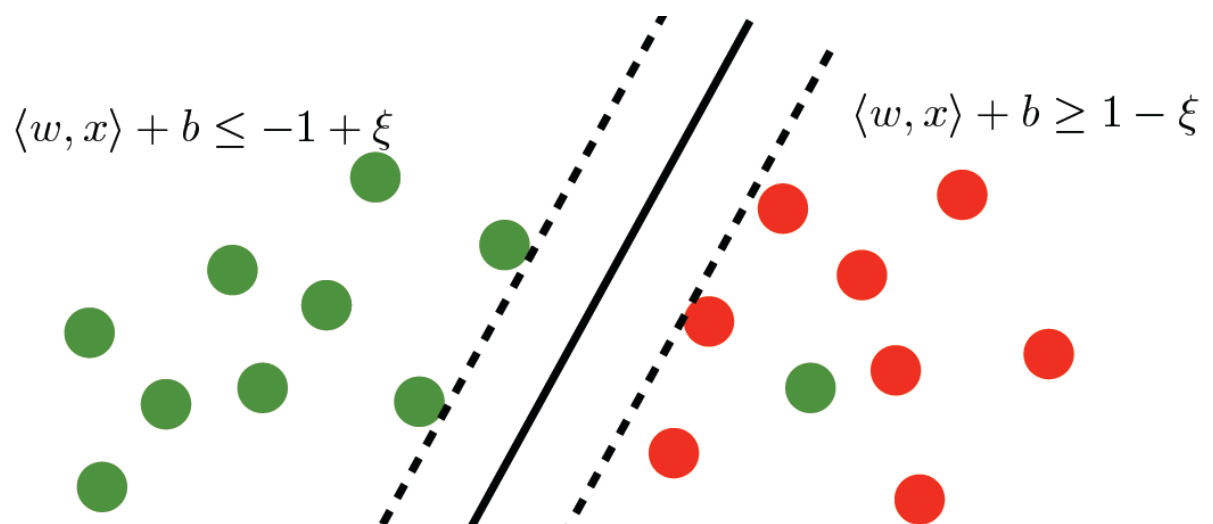
linear function
 $f(x) = \langle w, x \rangle + b$

Inria

Grenoble INP
ensimag

Support Vector Machines

- What about linearly inseparable cases ?



Inria

Grenoble INP
ensimag

Action recognition in realistic videos

Challenges

- Severe camera motion
- Variation in human appearance and pose
- Cluttered background and occlusion
- Viewpoint and illumination changes

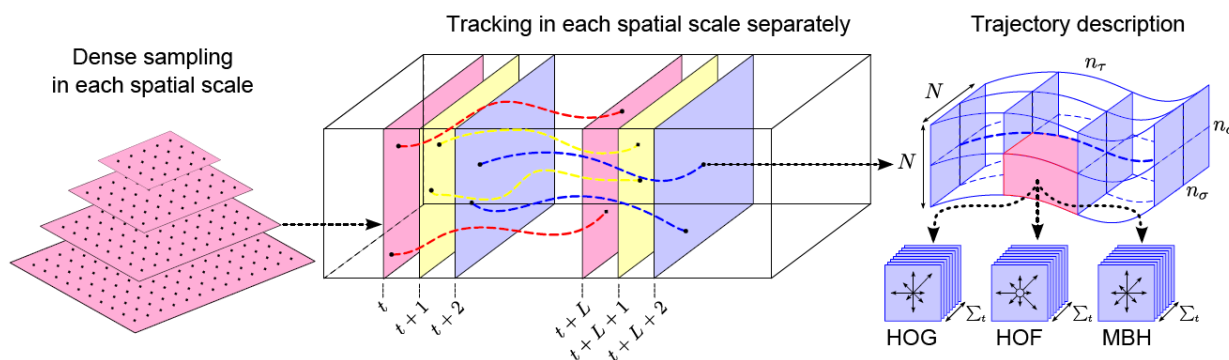
Current state of the art

- Local space-time features + bag-of-features model
- Dense trajectories performs the best on a large variety of datasets [Wang et.al. IJCV'13]



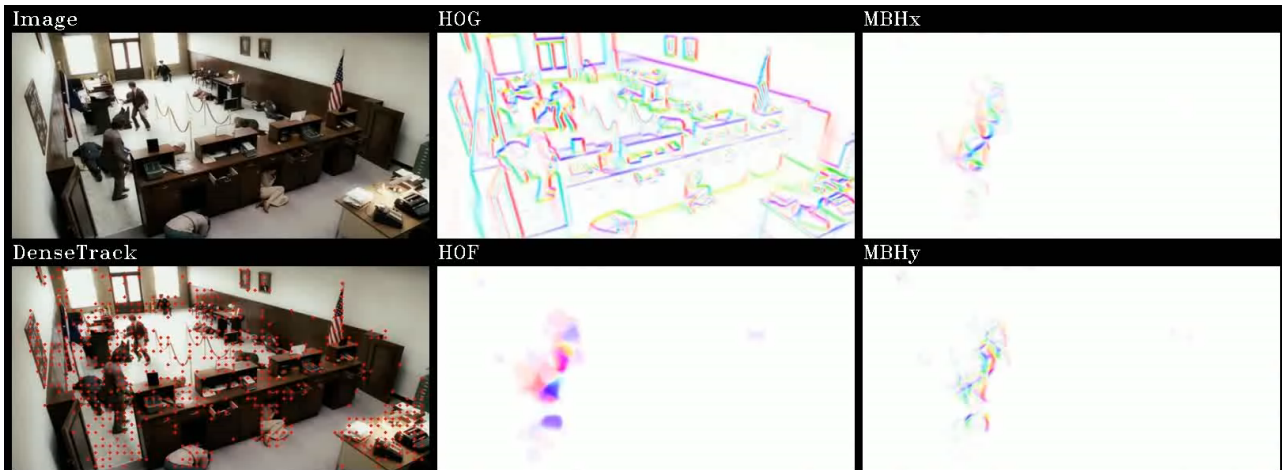
Dense trajectories

- Three major steps:
 - Dense sampling
 - Feature tracking
 - Trajectory-aligned descriptors



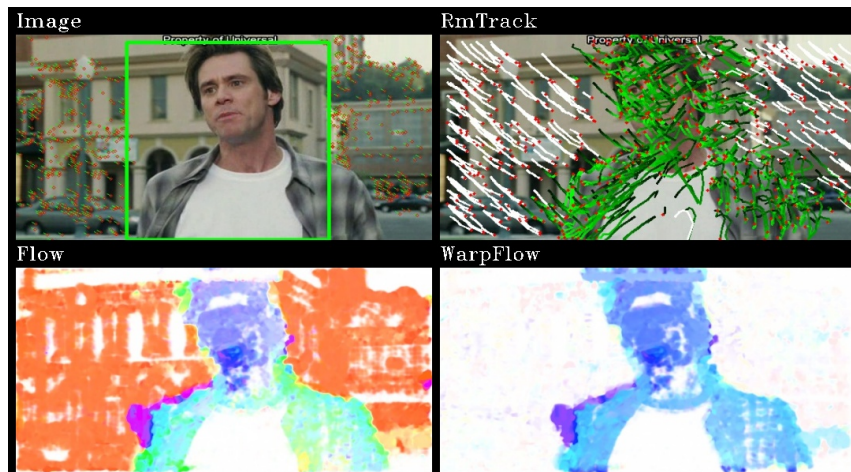
Dense trajectories

- Advantages:
 - Capture
 - MBH is robust to camera motion
- Disadvantages:
 - Generate irrelevant trajectories in background due to camera motion
 - Motion descriptors are corrupted due to camera motion, e.g., HOF, MBH



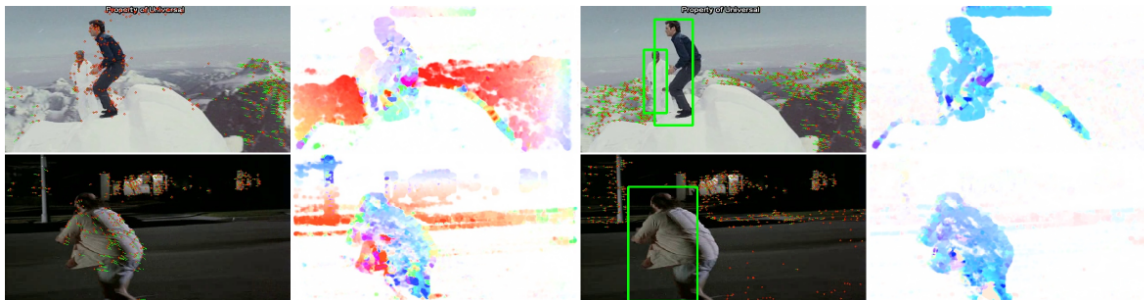
Improved dense trajectories

- Improve dense trajectories by explicit camera motion estimation
- Detect humans to remove outlier matches for homography estimation
- Stabilize optical flow to eliminate camera motion
- Remove trajectories caused by camera motion



Remove inconsistent matches due to humans

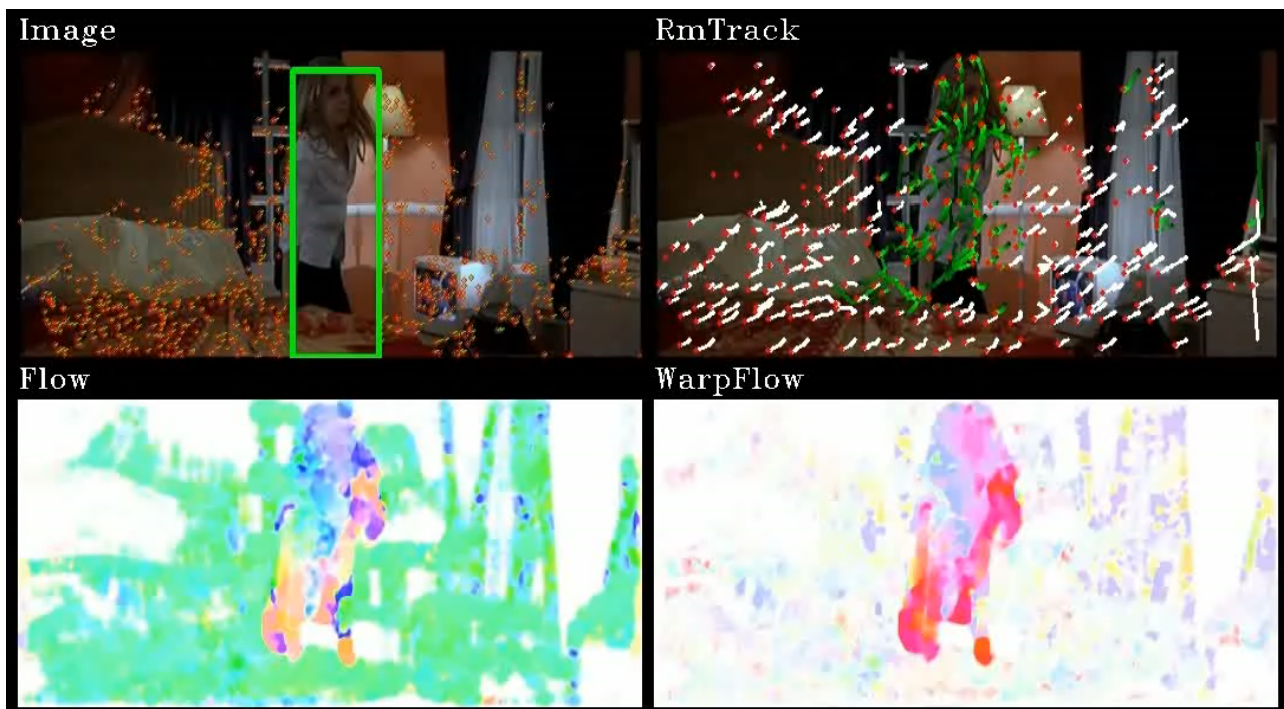
- Human motion: not constrained by camera motion → generates outlier matches
- Apply a human detector in each frame, and track the human bounding box
- Remove feature matches inside the human bounding box during estimation



Inria

Grenoble INP
ensimag

Example of the refined flow



Source code: http://lear.inrialpes.fr/~wang/improved_trajectories

Inria

Grenoble INP
ensimag

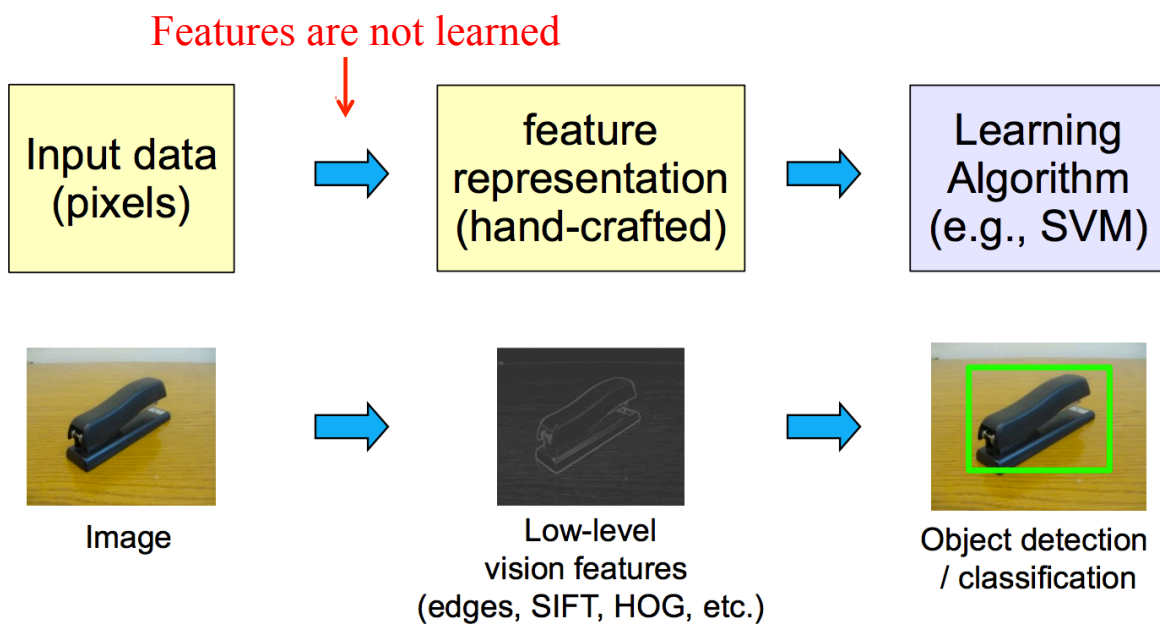
Plan

- Classification
- Machine à vecteurs de support
- Reconnaissance d'actions
- **CNN**
- Quelques perspectives

Inria



Traditional Approaches for Recognition



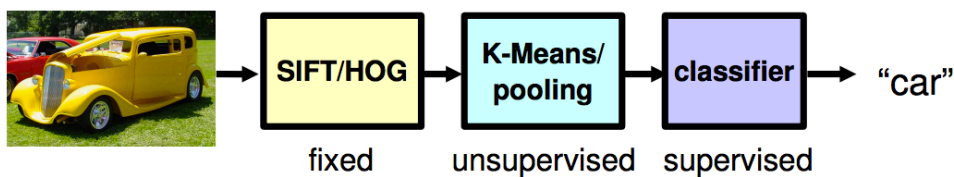
Inria

CNN Slides courtesy: M. Renzato, R. Fergus

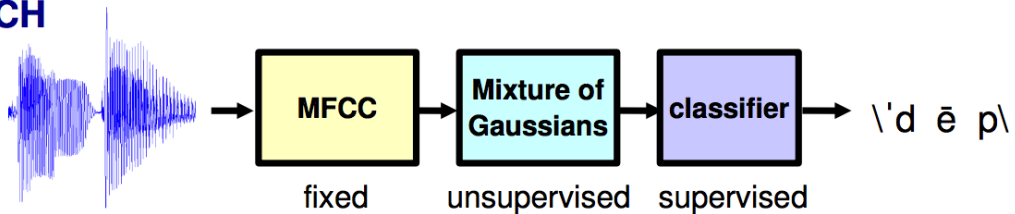


Traditional Approaches for Recognition

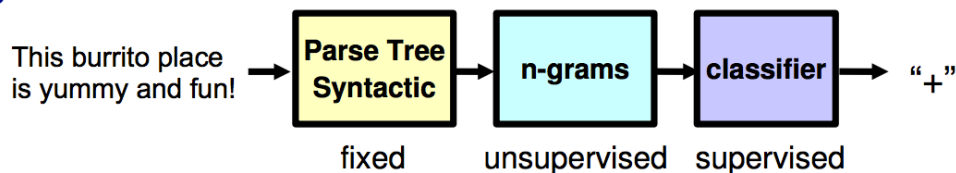
VISION



SPEECH



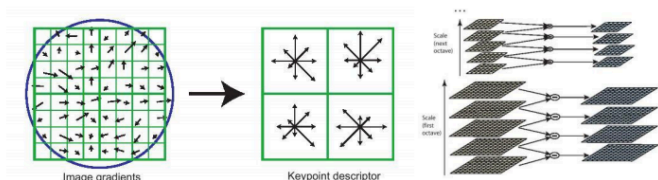
NLP



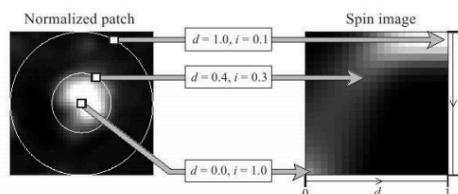
Inria

Grenoble INP
ensimag

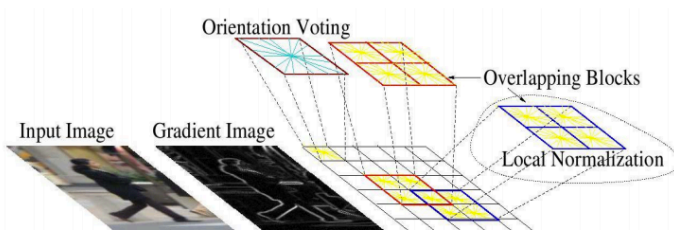
Computer Vision Features



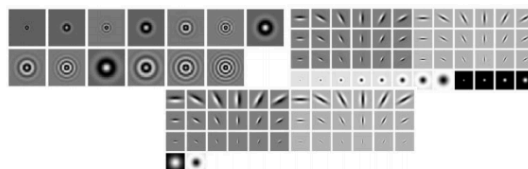
SIFT



Spin image



HoG



Textons

and many others:

SURF, MSER, LBP, Color-SIFT, Color histogram, GLOH,

Inria

Grenoble INP
ensimag

Computer Vision Features

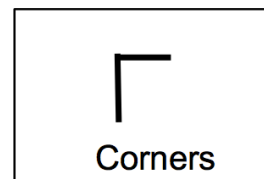
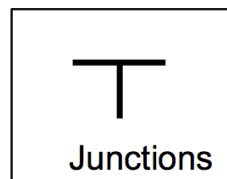
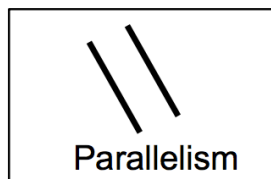
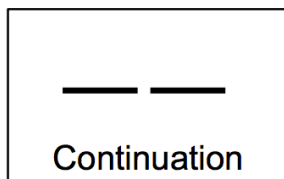
- Features are key to progress
- Have led to impressive results in various competitions (e.g., PASCAL VOC)
- Where do we go from here? Better features? Better classifiers?

Inria



Mid-level Representations

- Mid-level cues



“Tokens” from Vision by D.Marr:



- Object parts:



Inria



Mid-level Representations

VISION

pixels → edge → texture → motif → part → object

SPEECH

sample → spectral band → formant → motif → phone → word

NLP

character → word → NP/VP/.. → clause → sentence → story

Difficult to hand-engineer → What about learning them?

Inria



Learning Feature Hierarchy

- Learn hierarchy
- All the way from pixels → classifier
- One layer extracts features from output of previous layer



- Train all layers jointly

Inria



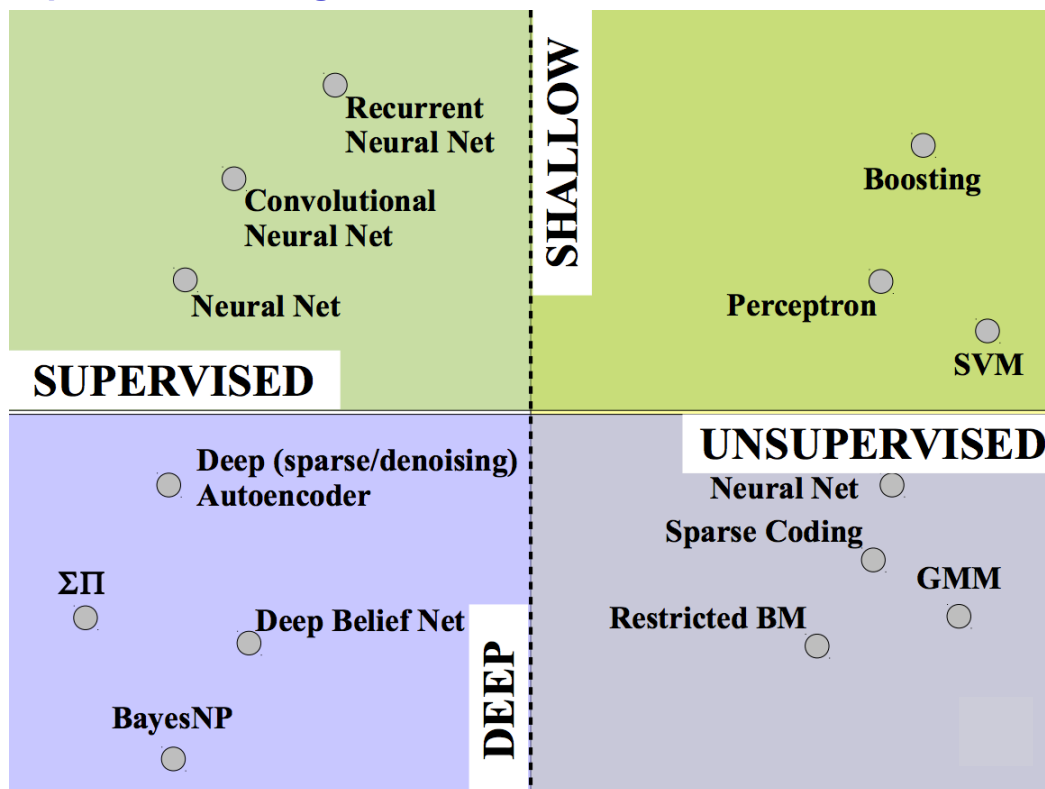
Learning Feature Hierarchy

- Supervised Learning
 - End-to-end learning of deep architectures (e.g., deep neural networks) with back-propagation
 - Works well when the amounts of labels is large
 - Structure of the model is important (e.g. convolutional structure)
- Unsupervised Learning
 - Learn statistical structure or dependencies of the data from unlabeled data
 - Layer-wise training
 - Useful when the amount of labels is not large

Inria



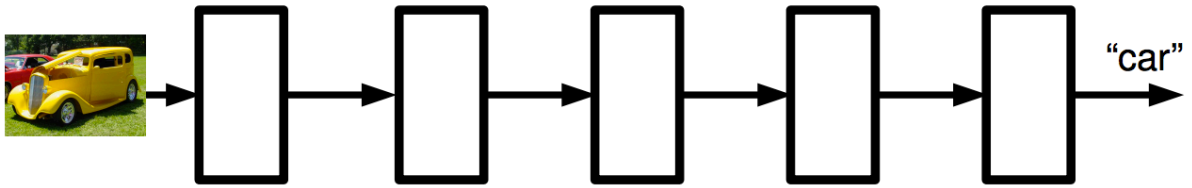
The Space of Learning Methods



Inria



Deep Learning



What is Deep Learning

- Cascade of non-linear transformations
- End to end learning
- General framework (any hierarchical model is deep)

Inria



But, before that, some basics (... on the board)

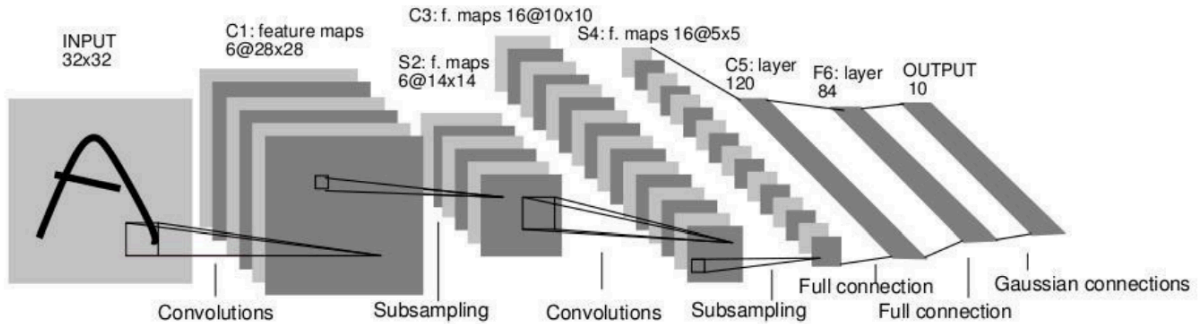
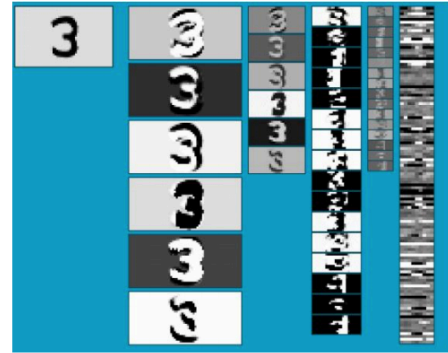
- Perceptron
 - Single-layer
 - Multi-layer
- Non-linearity
 - Sigmoid
- Backpropagation algorithm
- Neural networks

Inria



Example: Convolutional Neural Networks (CNN)

- LeCun et al. 1989
- Neural network with specialized connectivity structure

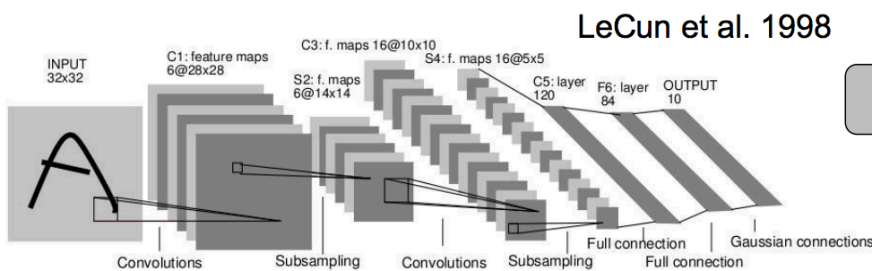
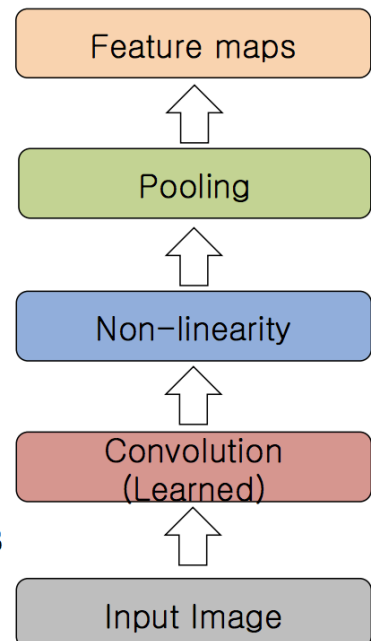


Inria



CNN

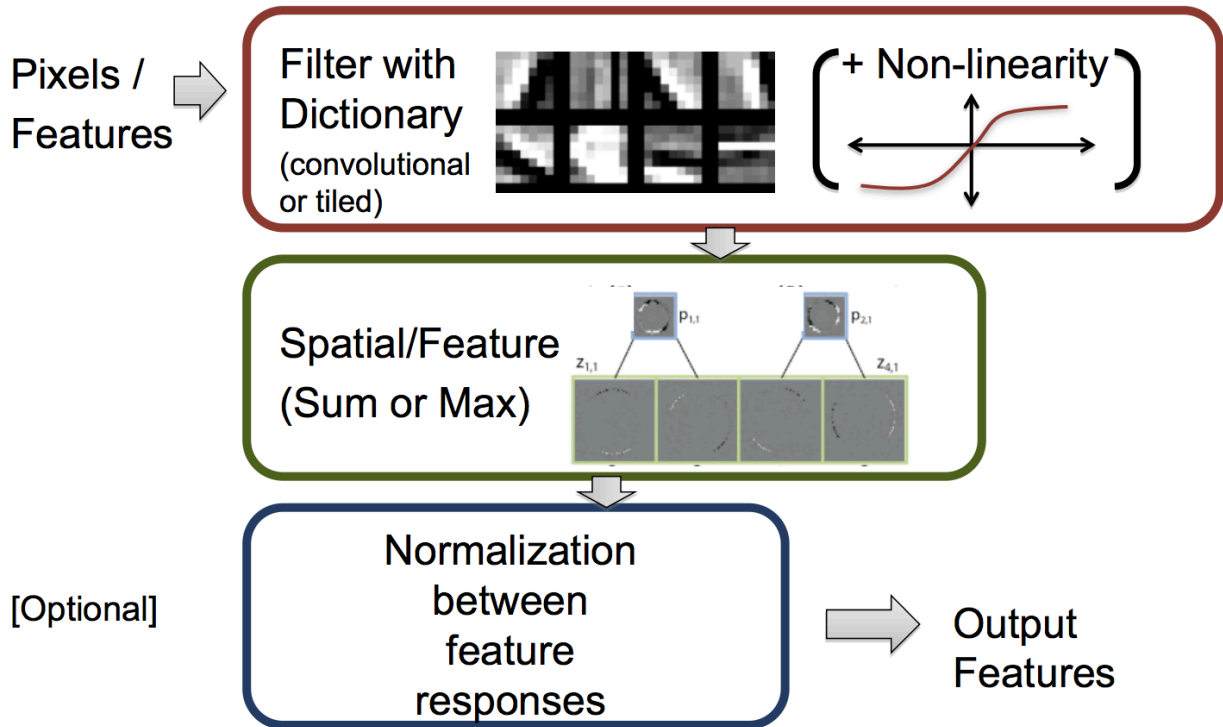
- Feed-forward:
 - Convolve input
 - Non-linearity (rectified linear)
 - Pooling (local max)
- Supervised
- Train convolutional filters by back-propagating classification error



Inria



Components of each CNN Layer

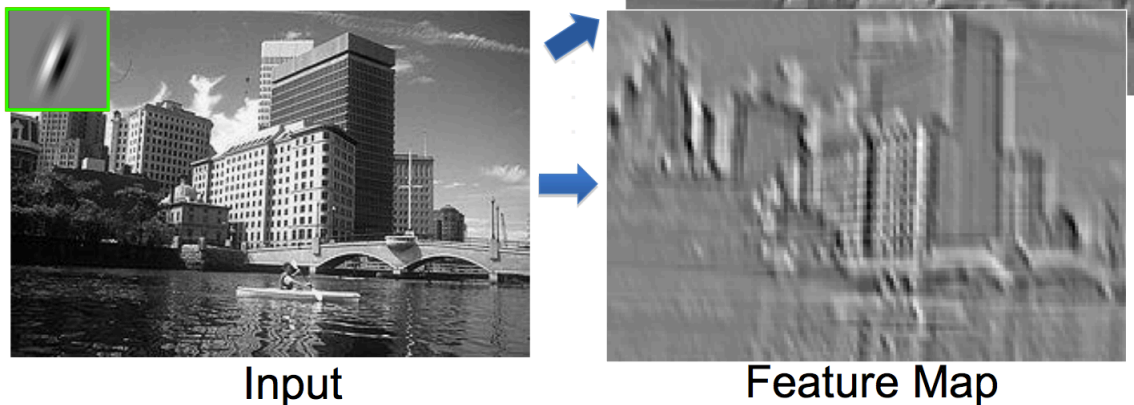


Inria

Grenoble INP
ensimag

Filtering: Convolutional

- Convolutional
 - Dependencies are local
 - Translation equivariance
 - Tied filter weights (few params)
 - Stride 1,2,... (faster, less mem.)

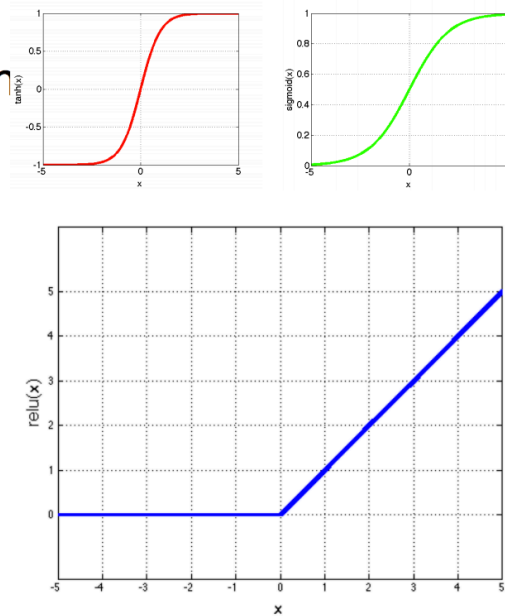


Inria

Grenoble INP
ensimag

Filtering: Non-linearity

- Non-linearity
 - Per-element (independent)
 - **Tanh**
 - **Sigmoid**: $1/(1+\exp(-x))$
 - **Rectified linear**
 - Simplifies backprop
 - Makes learning faster
 - Avoids saturation issues
- Preferred option

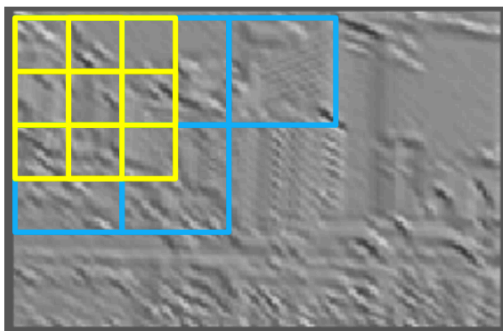


Inria

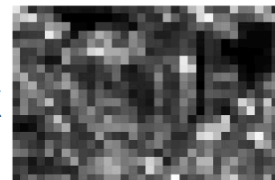
Grenoble INP
ensimag

Pooling

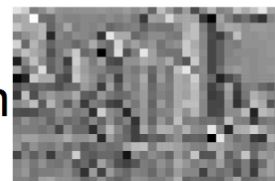
- Spatial Pooling
 - Non-overlapping / overlapping regions
 - Sum or max
 - Boureau et al. ICML'10 for theoretical analysis



Max



Sum

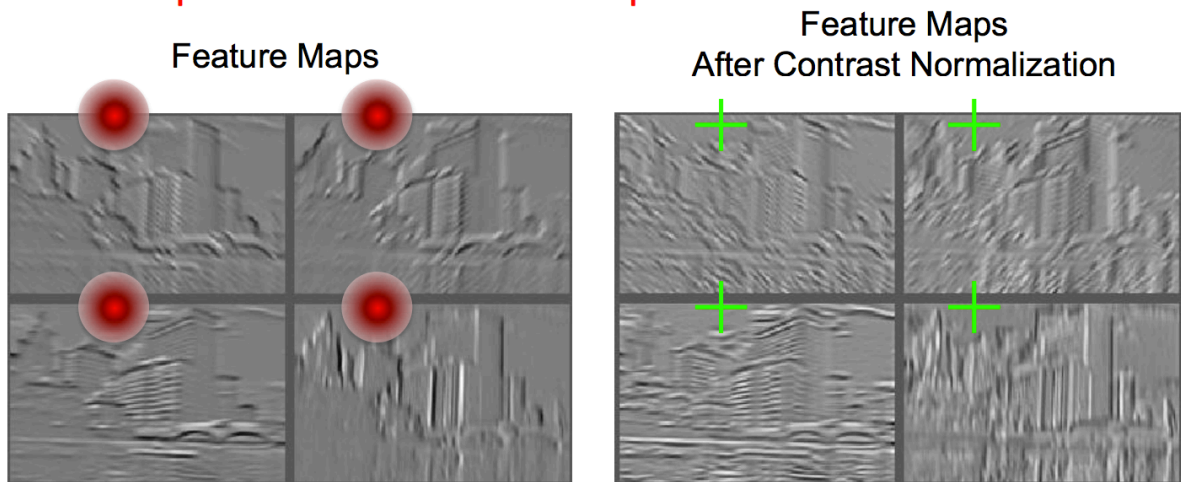


Inria

Grenoble INP
ensimag

Normalization

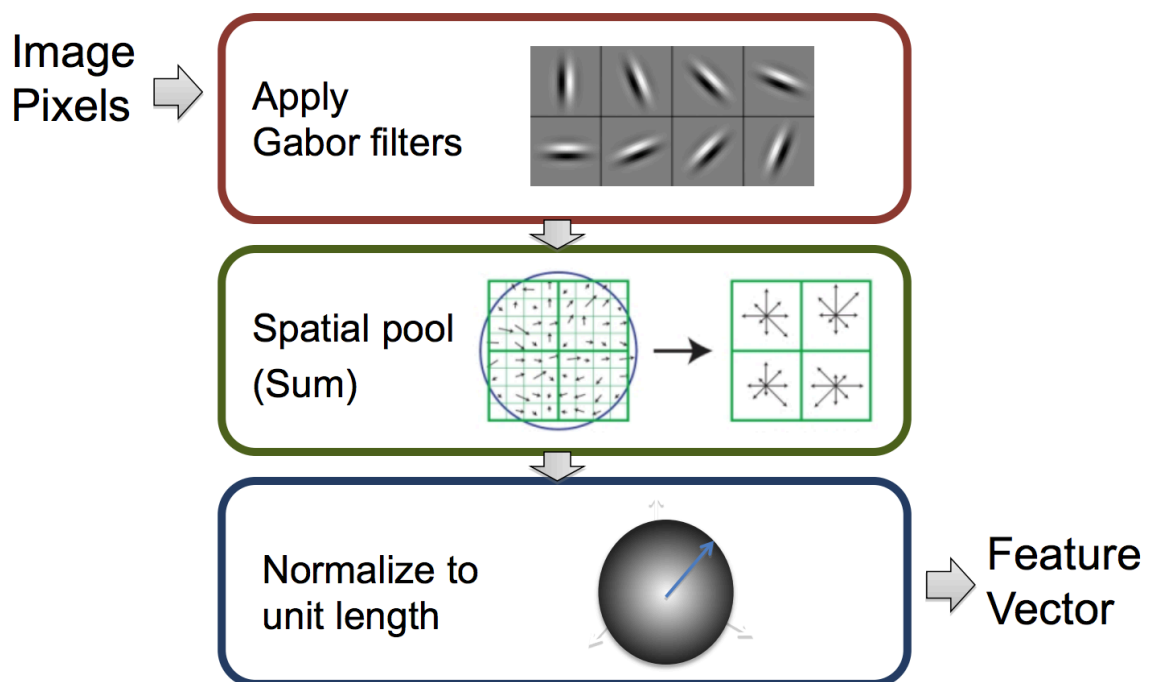
- Contrast normalization (across feature maps)
 - Local mean = 0, local std. = 1, “Local” \rightarrow 7x7 Gaussian
 - Equalizes the features maps



Inria

Grenoble INP
ensimag

Comparison with SIFT

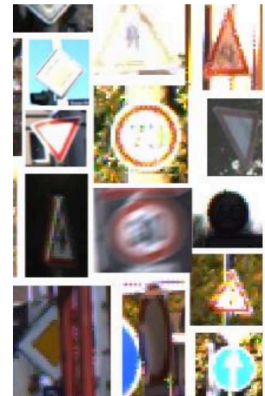


Inria

Grenoble INP
ensimag

CNN: Applications

- Handwritten text/digits
 - MNIST (0.17% error [Ciresan et al. 2011])
 - Arabic & Chinese [Ciresan et al. 2012]
- Simpler recognition benchmarks
 - CIFAR-10 (9.3% error [Wan et al. 2013])
 - Traffic sign recognition
 - 0.56% error vs 1.16% for humans [Ciresan et al. 2011]

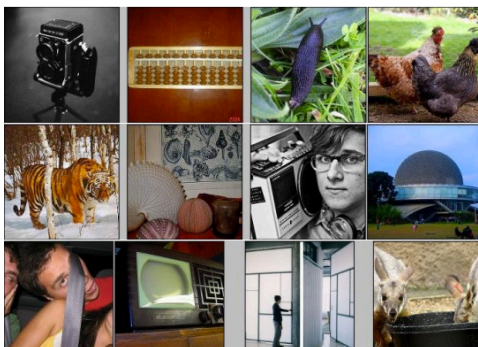


Inria

Grenoble INP
ensimag

CNN: Applications

IMAGENET



[Deng et al. CVPR 2009]

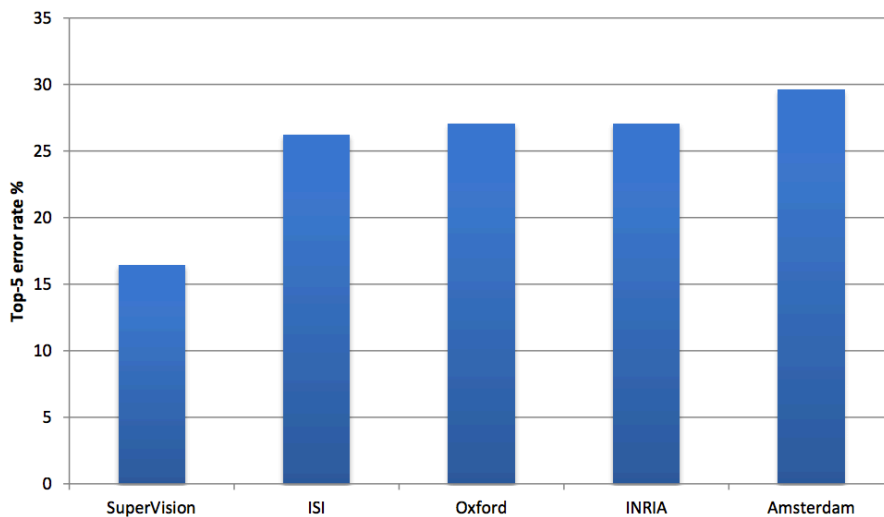
- ~14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon Turk

Inria

Grenoble INP
ensimag

ImageNet 2012

- Krizhevsky et al. -- 16.4% error (top-5)
- Next best (non-convnet) – 26.2% error

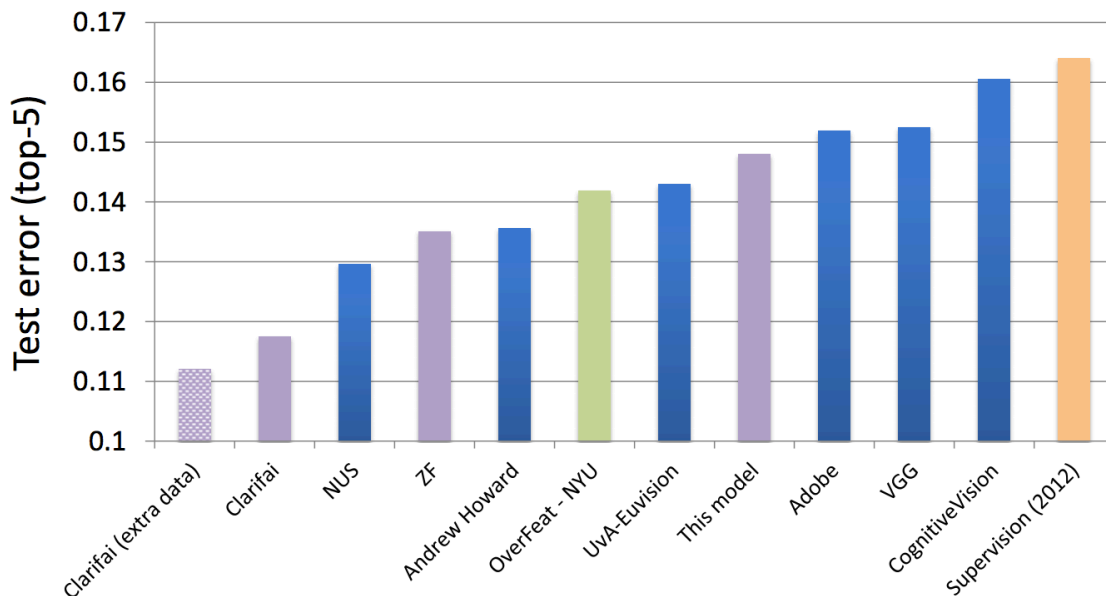


Inria

Grenoble INP
ensimag

ImageNet 2013

- <http://www.image-net.org/challenges/LSVRC/2013/results.php>



Inria

Grenoble INP
ensimag

Plan

- Classification
- Machine à vecteurs de support
- Reconnaissance d'actions
- CNN
- Quelques perspectives

Inria



Problèmes d'intérêt et perspectives

- En conclusion
 - ▶ domaine de recherche très actif car récent
 - les capacités de calcul des ordinateurs permettent depuis très peu de temps de traiter des grands volumes de données
 - ▶ Compétitions pour comparer les méthodes : ImageNet, Trecvid, VideOlympics

Inria

