

Parametric Models Fitting with Automatic Differentiation

Gabriel Peyré



www.numerical-tours.com





Mathematical Coffees

Huawei-FSMP joint seminars

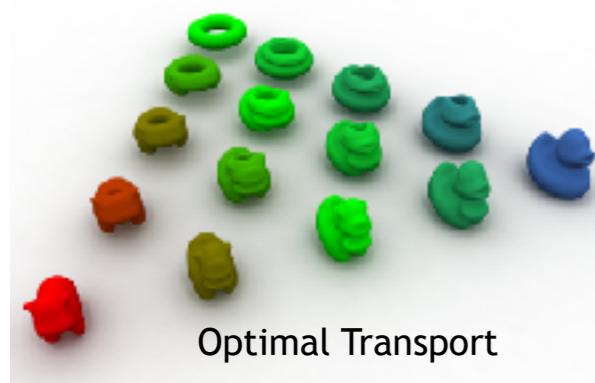
<https://mathematical-coffees.github.io>



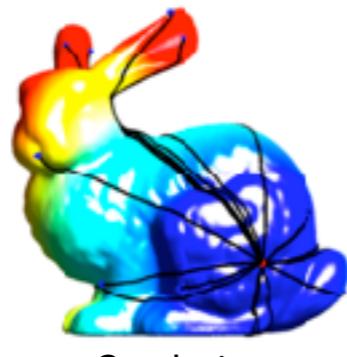
FSMP

Fondation Sciences
Mathématiques de Paris

Organized by: Mérouane Debbah & Gabriel Peyré



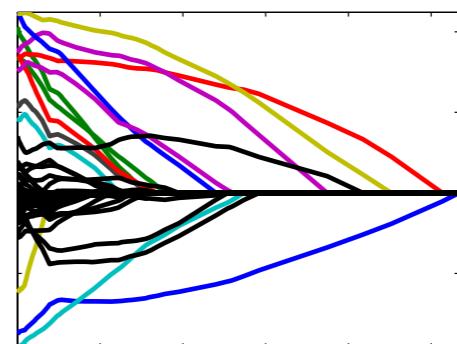
Optimal Transport



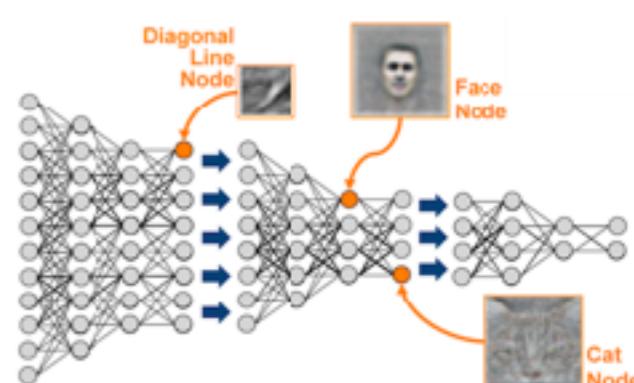
Geodesics



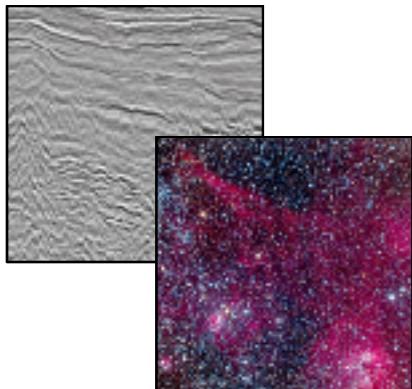
Mesches



Optimization



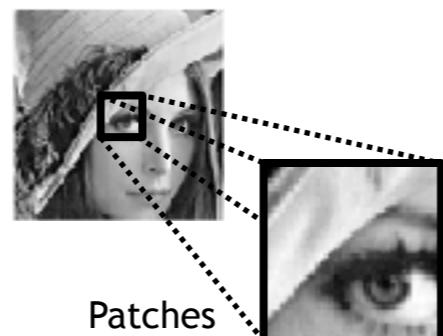
Deep Learning



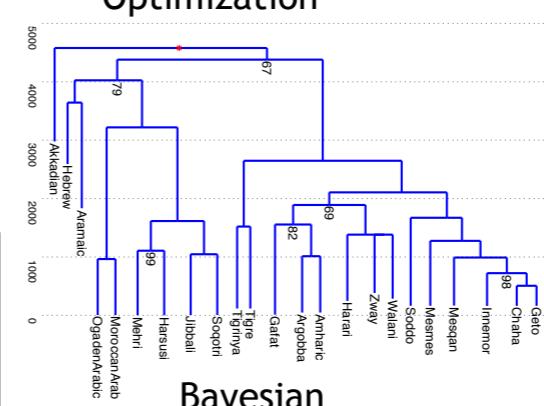
Sparsity



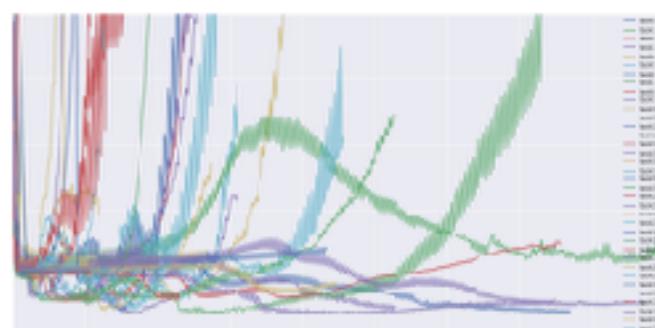
Neuro-imaging



Patches



Bayesian



Parallel/Stochastic

Alexandre Allauzen, Paris-Sud.
Pierre Alliez, INRIA.
Guillaume Charpiat, INRIA.
Emilie Chouzenoux, Paris-Est.

Nicolas Courty, IRISA.
Laurent Cohen, CNRS Dauphine.
Marco Cuturi, ENSAE.
Julie Delon, Paris 5.

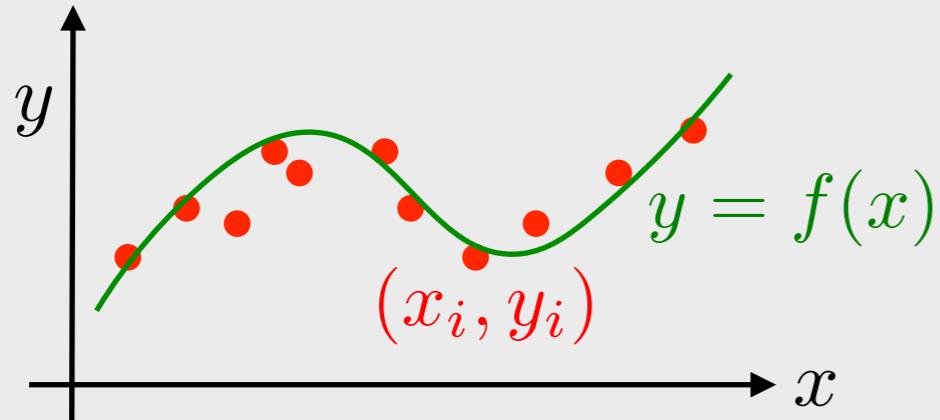
Fabian Pedregosa, INRIA.
Julien Tierny, CNRS and P6.
Robin Ryder, Paris-Dauphine.
Gael Varoquaux, INRIA.

Jalal Fadili, ENSICAEN.
Alexandre Gramfort, INRIA.
Matthieu Kowalski, Supelec.
Jean-Marie Mirebeau, CNRS, P-Sud.

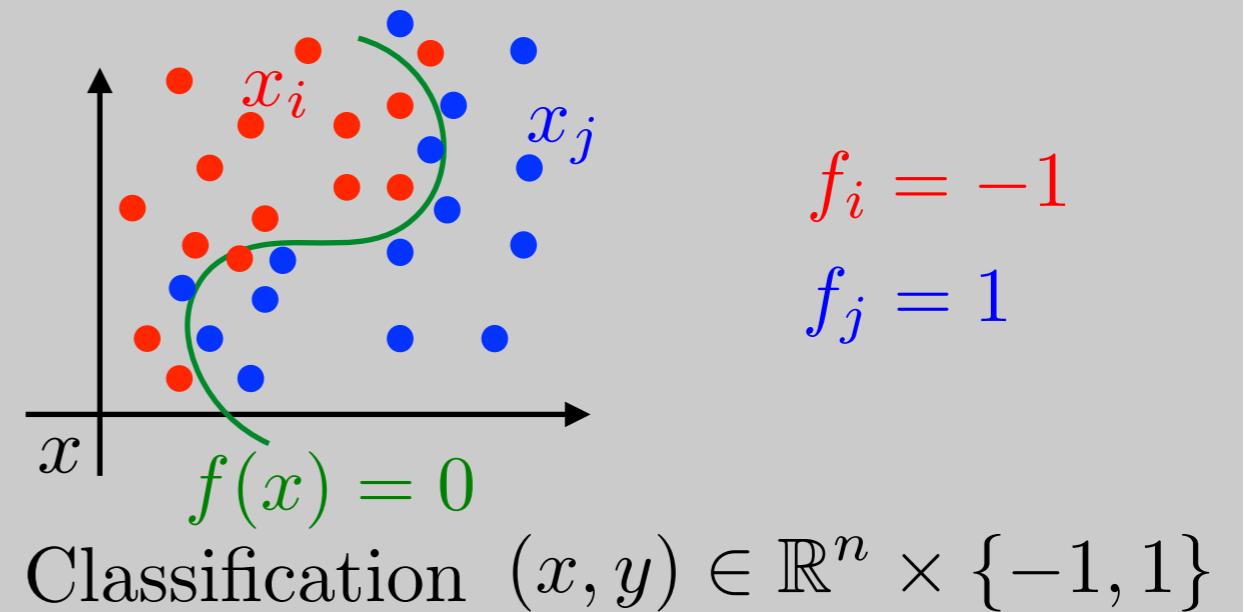


Parametric Models

(Noisy) observations (x_i, y_j) , try to infer $y = f(x)$.



Regression $(x, y) \in \mathbb{R}^n \times \mathbb{R}^p$

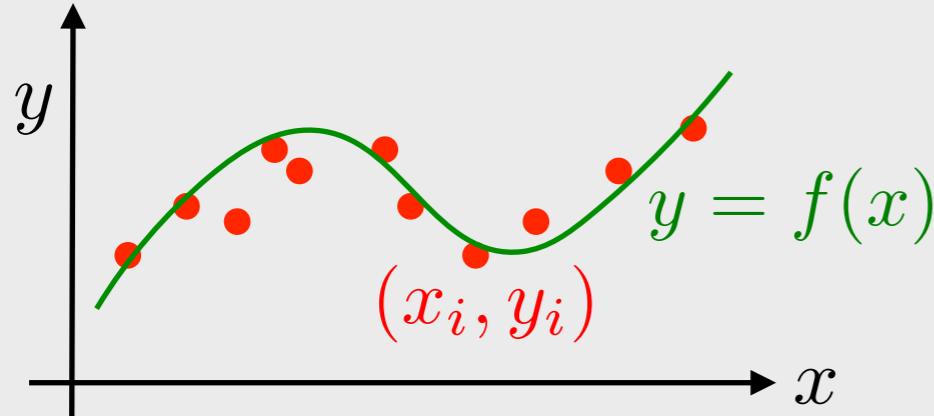


Classification $(x, y) \in \mathbb{R}^n \times \{-1, 1\}$

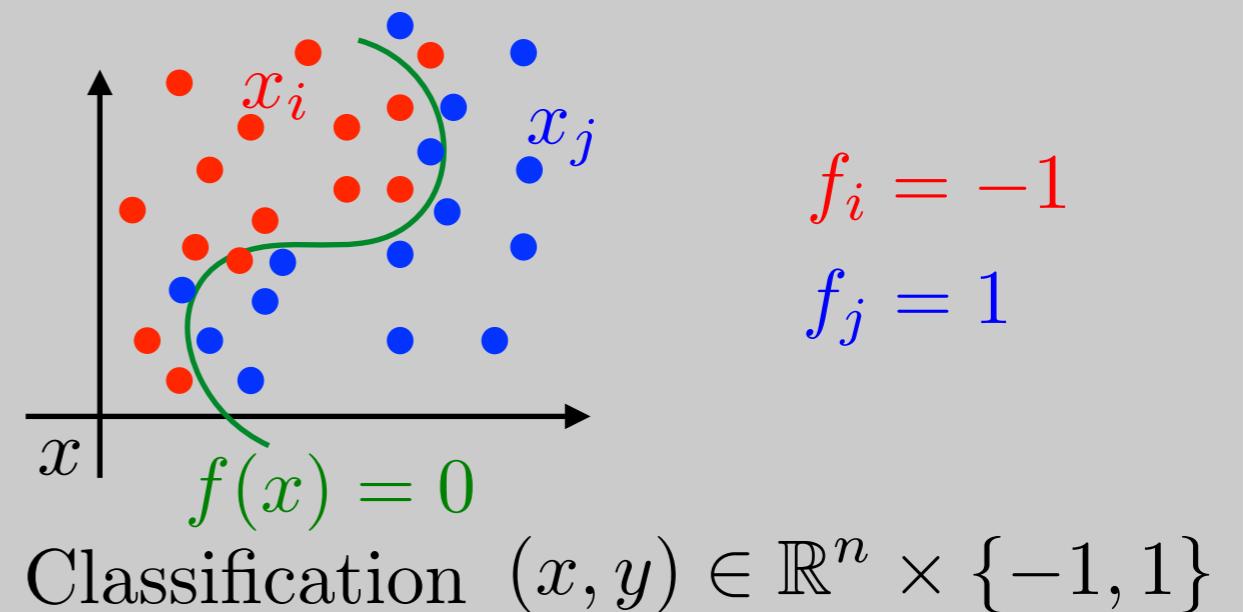
$$\begin{aligned}f_i &= -1 \\f_j &= 1\end{aligned}$$

Parametric Models

(Noisy) observations (x_i, y_j) , try to infer $y = f(x)$.

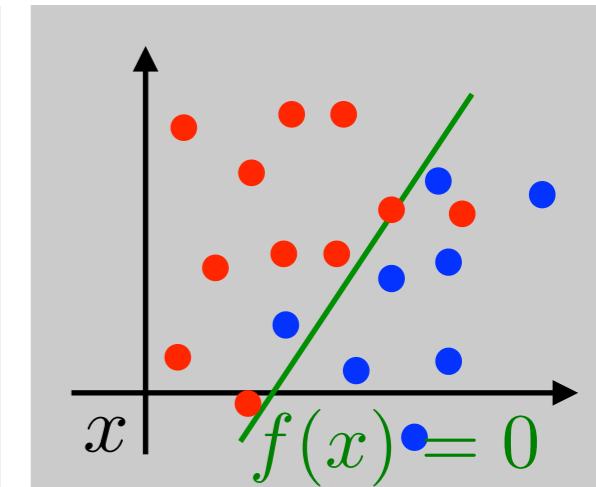
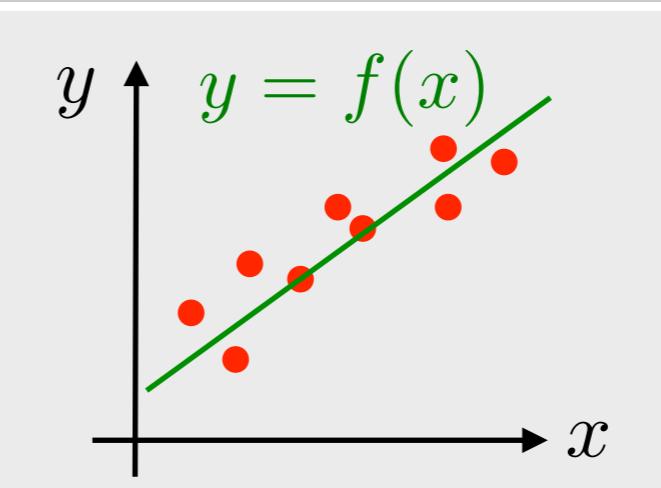


Regression $(x, y) \in \mathbb{R}^n \times \mathbb{R}^p$



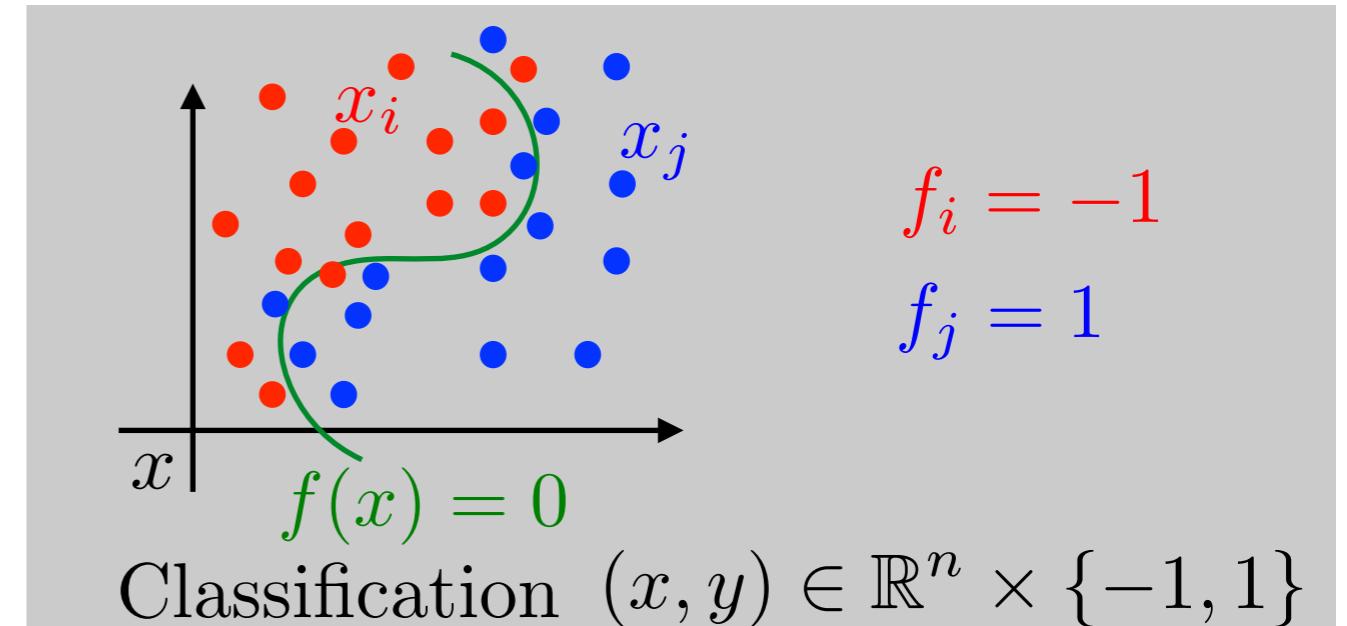
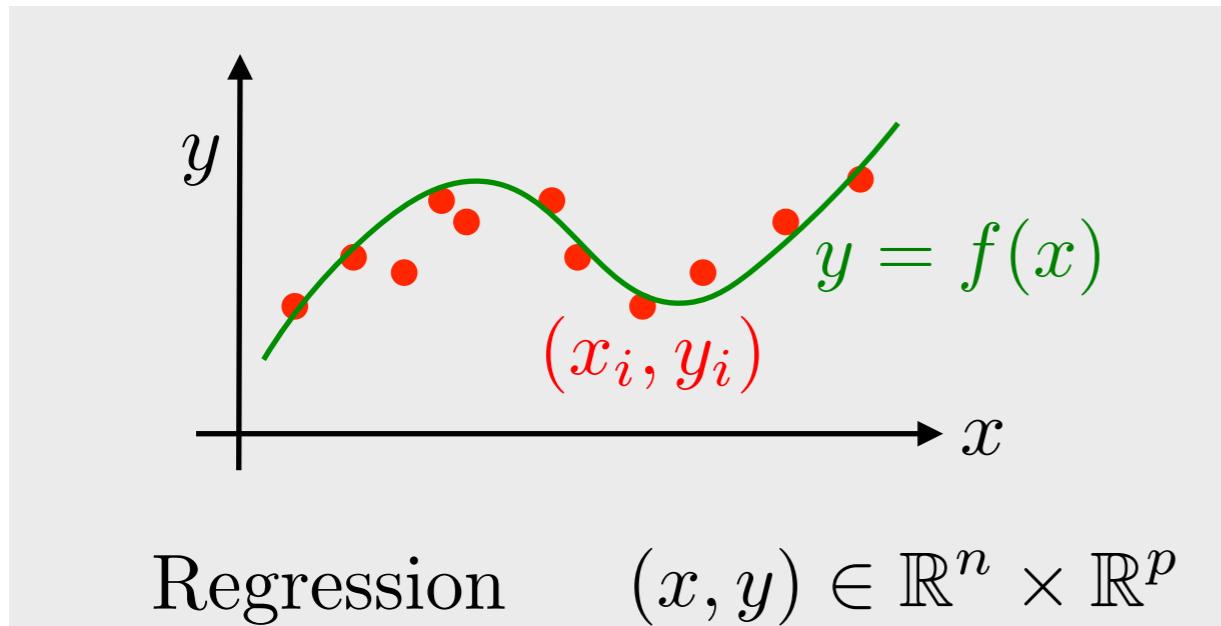
Parametric model: $y = f(x, \theta)$, find θ .

Linear model: $f(x, \theta) = \langle x, \theta \rangle$.



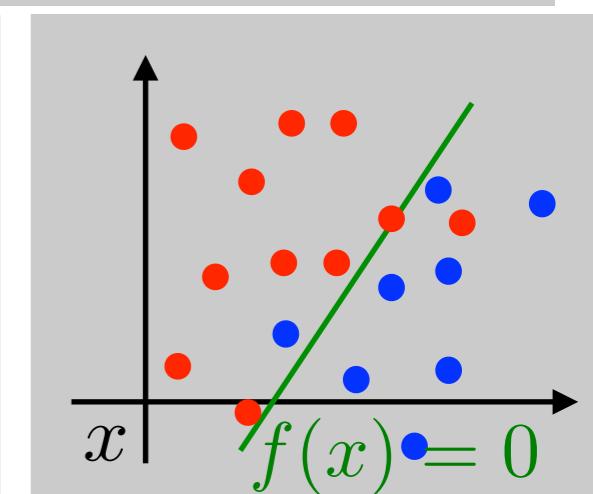
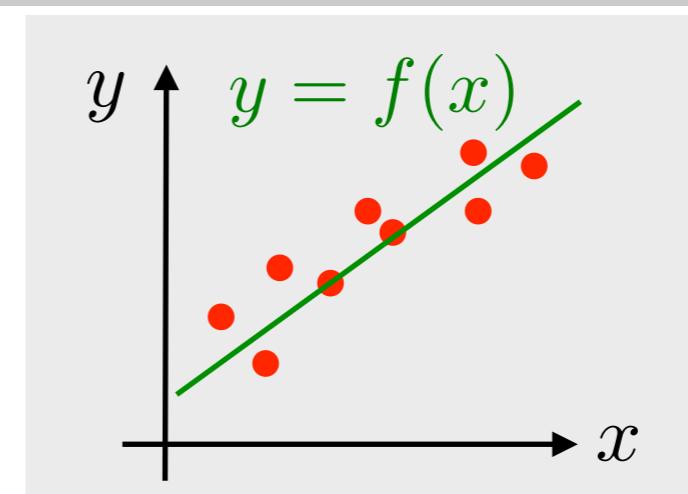
Parametric Models

(Noisy) observations (x_i, y_j) , try to infer $y = f(x)$.



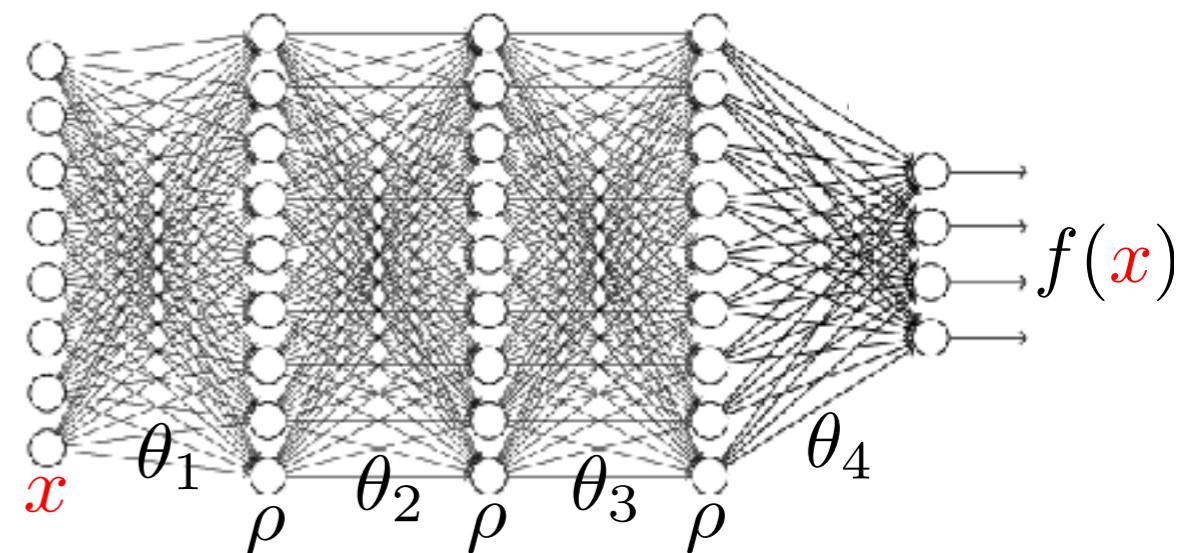
Parametric model: $y = f(x, \theta)$, find θ .

Linear model: $f(x, \theta) = \langle x, \theta \rangle$.



Deep network:

$$f(x, \theta) = \theta_K(\dots \rho(\theta_2(\rho(\theta_1(x) \dots)))$$



Empirical Loss Minimization

Regression: $(y, y') \in \mathbb{R}^d \times \mathbb{R}^d$, $L(y, y') = \|y - y'\|^2$

Classification: $(y, y') \in \mathbb{R}^d \times \{-1, 1\}$, $L(y, y') = \log(\exp(-y'y) + 1)$

Loss minimization:

$$\min_{\theta} \sum_i L(f(x_i, \theta), y_i)$$

$$\min_{\theta} \mathbb{E}_{(X,Y)}(L(f(X, \theta), Y))$$

Empirical Loss Minimization

Regression: $(y, y') \in \mathbb{R}^d \times \mathbb{R}^d, \quad L(y, y') = \|y - y'\|^2$

Classification: $(y, y') \in \mathbb{R}^d \times \{-1, 1\}, \quad L(y, y') = \log(\exp(-y'y) + 1)$

Loss minimization:

$$\min_{\theta} \sum_i L(f(x_i, \theta), y_i)$$

$$\min_{\theta} \mathbb{E}_{(X,Y)}(L(f(X, \theta), Y))$$

Stochastic gradient descent:

– Sample:

$$(x, y) \in \{(x_i, y_i)\}_i$$

$$(x, y) \sim (X, Y)$$

– Update: $\theta^{(\ell+1)} \stackrel{\text{def.}}{=} \theta^{(\ell)} - \tau_\ell \nabla_{\theta} \ell_{x,y}(\theta)$

where $\ell_{x,y}(\theta) \stackrel{\text{def.}}{=} L(f(x, \theta), y)$

Automatic Differentiation

How to compute $\nabla \ell_{x,y}(\theta)$? $\ell_{x,y}(\theta) \stackrel{\text{def.}}{=} L(f(x, \theta), y)$

Chain rule: $\nabla \ell_{x,y}(\theta) = [\partial f(x, \theta)]^\top (\nabla L(f(x, \theta), y))$

Linear $f(x, \theta) = \theta \times x$: $\partial f(x, \theta) = \theta$.

Non-linear $f(x, \theta)$: painful ... but $\ell_{x,y}$ it is just a computer program.

Automatic Differentiation

How to compute $\nabla \ell_{x,y}(\theta)$?

$$\ell_{x,y}(\theta) \stackrel{\text{def.}}{=} L(f(x, \theta), y)$$

Chain rule: $\nabla \ell_{x,y}(\theta) = [\partial f(x, \theta)]^\top (\nabla L(f(x, \theta), y))$

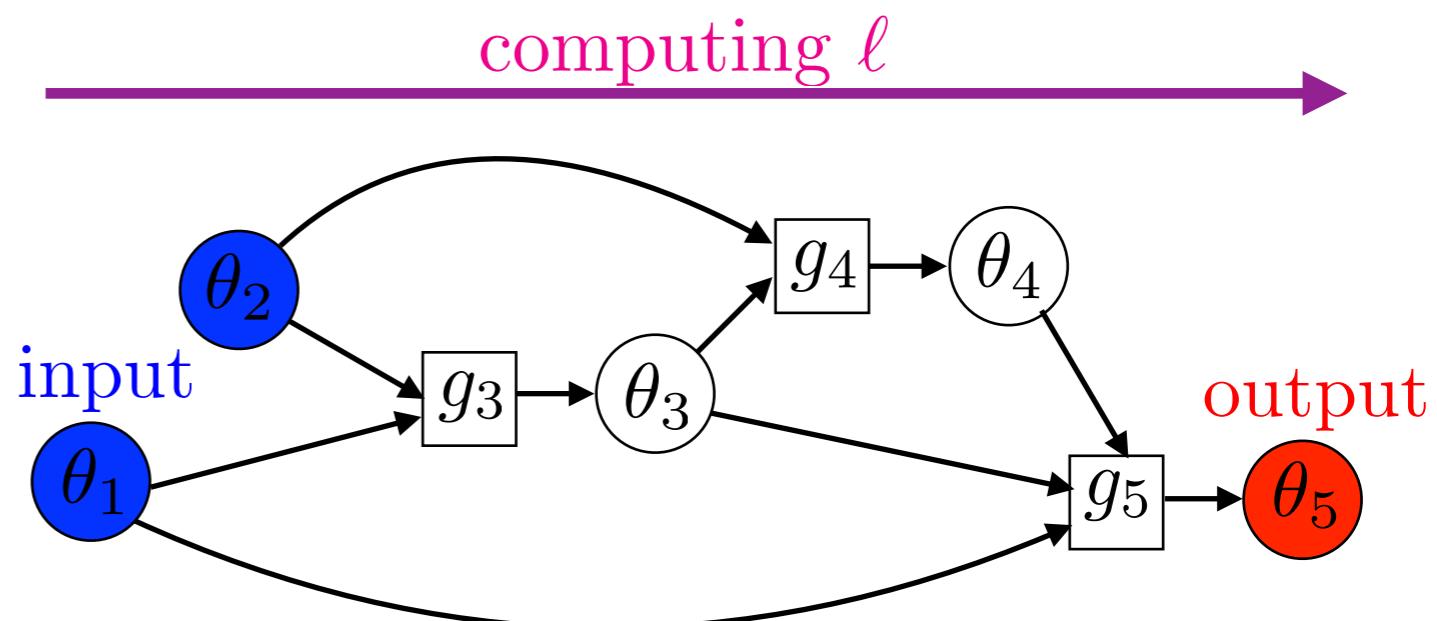
Linear $f(x, \theta) = \theta \times x$: $\partial f(x, \theta) = \theta$.

Non-linear $f(x, \theta)$: painful ... but $\ell_{x,y}$ it is just a computer program.

Computer program \Leftrightarrow directed acyclic graph \Leftrightarrow linear ordering of nodes $(\theta_r)_r$

forward

```
function  $\ell(\theta_1, \dots, \theta_M)$ 
  for  $r = M + 1, \dots, R$ 
    |  $\theta_r = g_r(\theta_{\text{Parents}(r)})$ 
  return  $\theta_R$ 
```



Automatic Differentiation

How to compute $\nabla \ell_{x,y}(\theta)$?

$$\ell_{x,y}(\theta) \stackrel{\text{def.}}{=} L(f(x, \theta), y)$$

Chain rule: $\nabla \ell_{x,y}(\theta) = [\partial f(x, \theta)]^\top (\nabla L(f(x, \theta), y))$

Linear $f(x, \theta) = \theta \times x$: $\partial f(x, \theta) = \theta$.

Non-linear $f(x, \theta)$: painful ... but $\ell_{x,y}$ it is just a computer program.

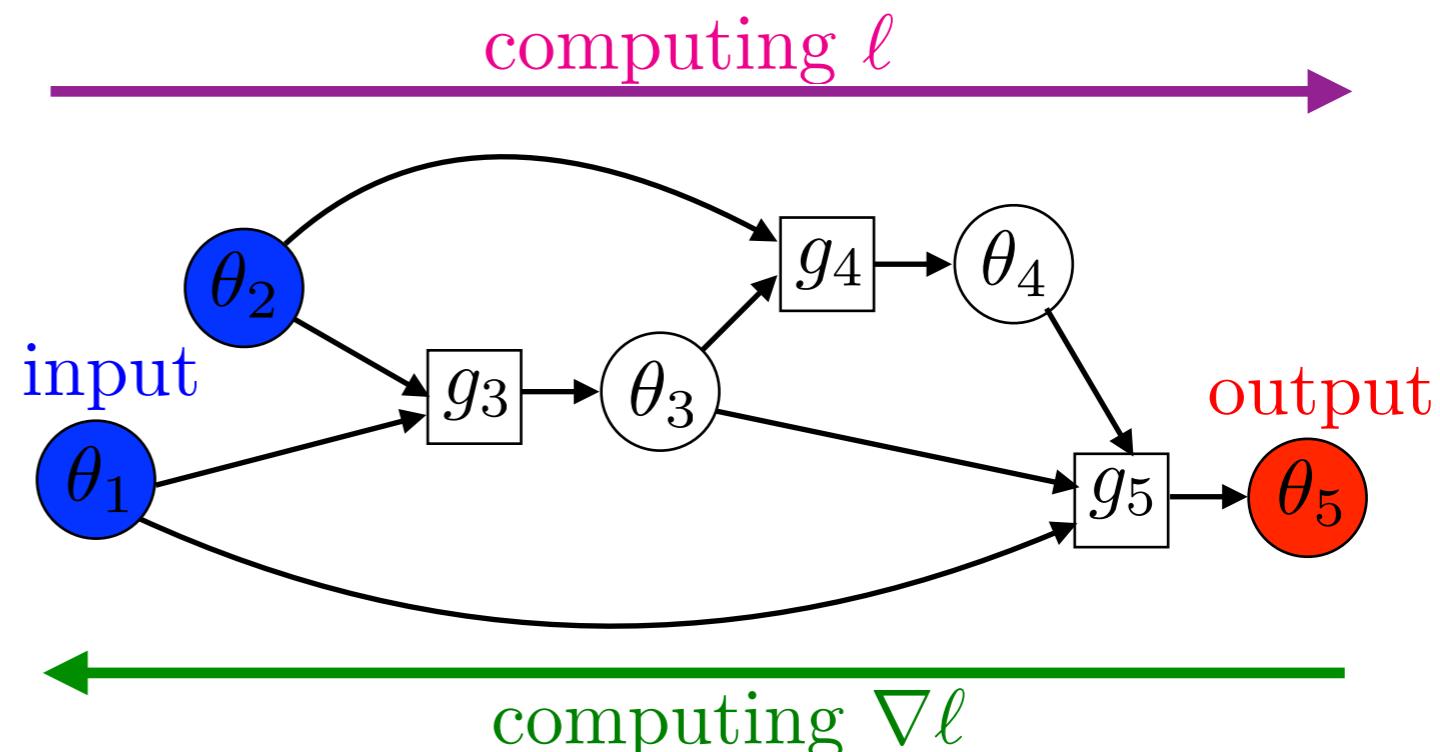
Computer program \Leftrightarrow directed acyclic graph \Leftrightarrow linear ordering of nodes $(\theta_r)_r$

forward

```
function  $\ell(\theta_1, \dots, \theta_M)$ 
  for  $r = M + 1, \dots, R$ 
    |  $\theta_r = g_r(\theta_{\text{Parents}(r)})$ 
  return  $\theta_R$ 
```

backward

```
function  $\nabla \ell(\theta_1, \dots, \theta_M)$ 
   $\nabla_R \ell = 1$ 
  for  $r = R - 1, \dots, 1$ 
    |  $\nabla_r \ell = \sum_{s \in \text{Child}(r)} \partial_r g_s(\theta) \nabla_s \ell$ 
  return  $(\nabla_1 \ell, \dots, \nabla_M \ell)$ 
```

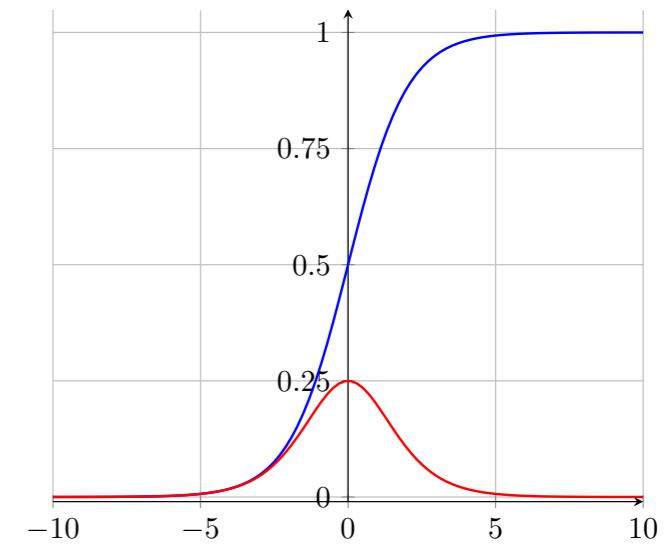
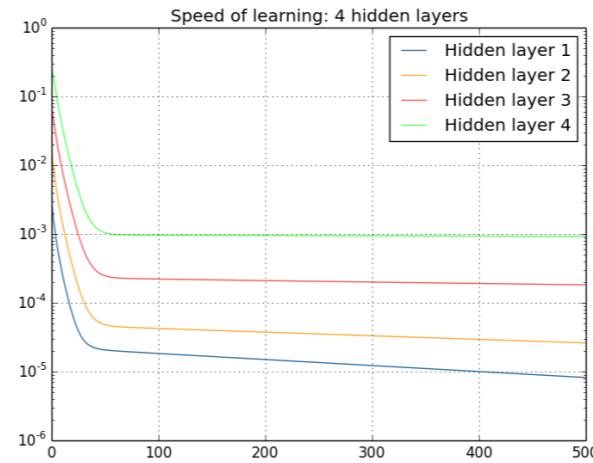
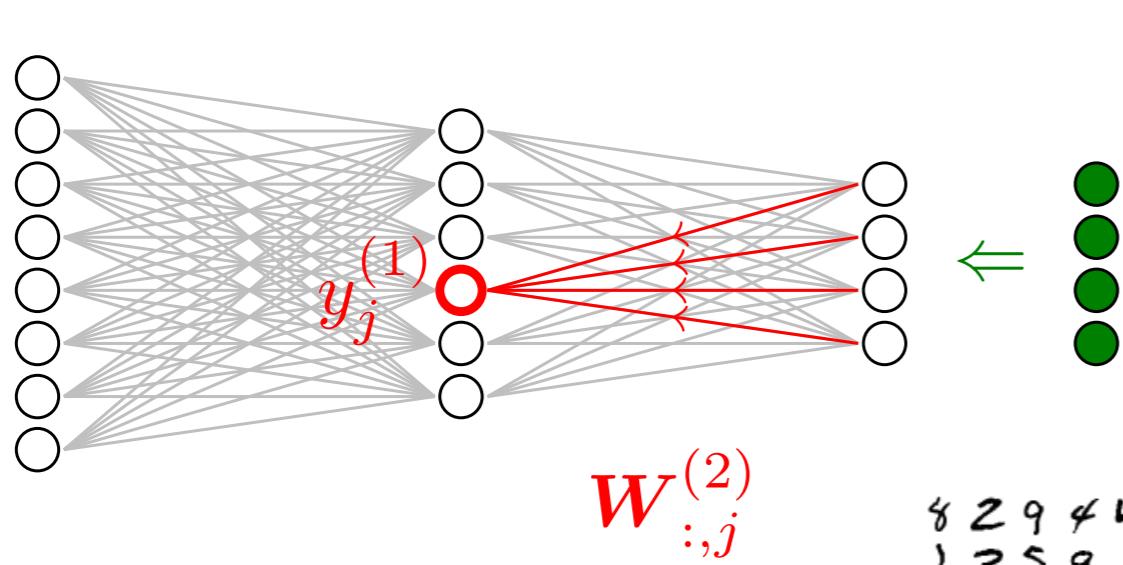


$$\frac{\partial \ell}{\partial \theta_r} = \sum_{s \in \text{Child}(r)} \frac{\partial \ell}{\partial \theta_s} \frac{\partial \theta_s}{\partial \theta_r}$$

$\nabla_r \ell(\theta)$ $\nabla_s \ell(\theta)$ $\partial_r g_s(\theta)$

What's next

Alexandre Allauzen: deep neural networks training.



8 2 9 4 4 6 4 9 7 0 9 2 9 5 1 5 9 1 0 3
2 3 5 9 1 7 6 2 8 2 2 5 0 7 4 9 7 8 3 2
1 1 8 3 6 1 0 3 1 0 0 1 1 2 7 3 0 4 6 5
2 6 4 7 1 8 9 9 3 0 7 1 0 2 0 3 5 4 6 5

Guillaume Charpiat: architecture of deep neural networks.

