Numerical Probability: Stochastic Approximation from Finance to Data Science

Gilles Pagès

LPSM-Sorbonne-Université

(Labo. Proba., Stat. et Modélisation)



12th European Summer School in Financial Mathematics

Univ. of Padova September 2-8, 2019

Gilles PAGÈS (LPSM)

Stochastic approximation I

LPSM-Sorbonne Université 1 / 94

Optimization (deterministic, the origins)

Examples from Finance

- Implicitation
- Minimization

Learning procedures

- Abstract Learning
- Supervised Learning
- Unsupervised Learning (clustering)

Stochastic algorithms/Approximation

- From Robbins-Monro to Robbins-Siegmund
- Stochastic Gradient Descent (SGD) and pseudo-SGD

Examples revisited by SFD

- Numerical probability
- Learning (supervised and unsupervised)

Application to Neural Networks and deep learning

- Linear neural network
- One hidden layer feedforward perceptron
- Toward deep learning
- Multilayer feedforward perceptron and Backpropagation



Optimization (deterministic, the origins)

- **Examples from Finance**
- Implicitation
- Minimization
- ³ Learning procedures
 - Abstract Learning
 - Supervised Learning
 - Unsupervised Learning (clustering)
- 4 Stochastic algorithms/Approximation
 - From Robbins-Monro to Robbins-Siegmund
 - Stochastic Gradient Descent (SGD) and pseudo-SGD
- 5 Examples revisited by SFD
 - Numerical probability
 - Learning (supervised and unsupervised)
- 6 Application to Neural Networks and deep learning
 - Linear neural network
 - One hidden layer feedforward perceptron
 - Toward deep learning
 - Multilayer feedforward perceptron and Backpropagation

Deterministic zero search and optimization

 Zero search: One aims at finding a zero θ* of a function h : ℝ^d → ℝ^d. In view of generic notations in stochastic approximation, we will denote

$$h(\theta), \ \theta \in \mathbb{R}^d$$



(d = 1 is mandatory just for graphs).

- Various methods:
 - Local recursive zero search (standard): θ₀ be fixed and let γ > 0 be small enough. Set

$$\theta_{n+1} = \theta_n - \gamma h(\theta_n), \quad n \ge 0$$

- Various methods (the sequel):
 - Local recursive zero search if $h C^1$ (Newton-Raphson "false position" algoritm)

 $\theta_{n+1} = \theta_n - [J_h(\theta_n)]^{-1}h(\theta_n), \quad n \ge 0,$

where $J_h(\theta)$ denotes the Jacobian of *h* at θ .

Idea: The tangent hyperplane is the best approximation of h (by an hyperplane)

$$h(\theta) \simeq h(\theta_n) + J_h(\theta_n)(\theta_n - \theta)$$

so θ_{n+1} is solution to $h(\theta_n) + J_h(\theta_n)(\theta_n - \theta) = 0$.



Very fast but also very unstable, especially when $J_h(\theta^*)$ is "small".

Yet another local recursive zero search if h C¹ (Levenberg-Marquardt algorithm): Let λ_n > 0, n ≥ 1,

$$\theta_{n+1} = \theta_n - \left[J_h(\theta_n) + \lambda_{n+1}I_d\right]^{-1}h(\theta_n), \quad n \ge 0.$$

turns out to be more stable.

Gilles PAGÈS (LPSM)

4 / 94

• Global recursive zero search:

- Idea: make the step decrease (not too fast) to "enlarge" in an adaptive way the convergence area of the algorithm...

– Let
$$\gamma_n$$
, $n \ge 1$ satisfy

$$\sum_{n} \gamma_{n} = +\infty \text{ and } \sum_{n} \gamma_{n}^{2} < +\infty.$$
$$\theta_{n+1} = \theta_{n} - \gamma_{n+1} h(\theta_{n}), \ n \ge 0$$

• Etc.

- Set

• WARNING! All these methods require that

h to be computed at a reasonable cost.

Minimizing a (potential function)

• Gradient descent (GD):

Let $V : \mathbb{R}^d \to \mathbb{R}_+$, \mathcal{C}^1 with $\lim_{|x| \to +\infty} V(x) = +\infty$ so that $\operatorname{argmin} V \neq \emptyset$.

How to compute $\operatorname{argmin} \& \min_{\mathbb{R}^d} V$???

• If V is convex, then

$$\operatorname{argmin} V = \{\nabla V = 0\}$$
 (is a convex set)

- Solution: set $h = \nabla V$,

– If ∇V Lipschitz, then (exercise)

$$heta_n o heta^* \in \{
abla V = 0 \} = \operatorname{argmin} V \quad ext{ as } \quad n o +\infty.$$

• If V is not convex only

$$\operatorname{argmin} V \subsetneq \{\nabla V = 0\}.$$

Still set $h = \nabla V$ (what else?)





7 / 94

• Pseudo-gradient (back to zero search!):

The function h is often given (model) and (hopefully) there exists a Lyapunov function V s.t. $(h|\nabla V) \ge 0$ and

$$\{h = 0\} \simeq \{(h|\nabla V) = 0\}.$$

If $(d = 2)$, $\mathcal{H}(V)(x) = \begin{pmatrix} -\partial_{x_2} V \\ \partial_{x_1} V \end{pmatrix}$ (Hamiltonian of $\nabla V(x)$) and
 $h(x) = \lambda \nabla V(x) + \mu \mathcal{H}(V)(x)$

then, the above conditions are satisfied and $|h|^2$ has V-linear growth so that $\theta_n \to C(0; 1)$ (if $\theta_0 \neq 0$) but does not converge "pointwise".



However, on this example, $V(\theta_n) \rightarrow \operatorname{argmin} V$.

• It may happen that $\{h = 0\} \neq \{(h | \nabla V) = 0\} \neq \{\nabla V = 0\} \neq \operatorname{argmin} V !!$.



ptimization (deterministic, the origins)

- Examples from Finance
- Implicitation
- Minimization
- Learning procedures
 - Abstract Learning
 - Supervised Learning
 - Unsupervised Learning (clustering)
- 4 Stochastic algorithms/Approximation
 - From Robbins-Monro to Robbins-Siegmund
 - Stochastic Gradient Descent (SGD) and pseudo-SGD
- 5 Examples revisited by SFD
 - Numerical probability
 - Learning (supervised and unsupervised)
- 6 Application to Neural Networks and deep learning
 - Linear neural network
 - One hidden layer feedforward perceptron
 - Toward deep learning
 - Multilayer feedforward perceptron and Backpropagation



Implicitation: Volatility

- Black-Scholes model: traded asset $X_t = x_0 e^{(r \frac{\sigma^2}{2})t + \sigma W_t}$, x_0 , volatility $\sigma > 0$, interest rate r, W standard Brownian motion.
- Call payoff (X_τ − K)₊ = max(X_τ − K, 0) with strike price K and maturity T.
- Mark-to-Market quoted price: $\operatorname{Call}_{M2Mkt} \in (0, x_0)$.
- Black-Scholes price at time 0

$$\begin{aligned} \operatorname{Call}_{BS}(x_0, K, r, \sigma) &= e^{-rT} \mathbb{E} \left(X_{\tau} - K \right)_+ \\ &= x_0 \Phi_0(d_1) - K e^{-rt} \Phi_0(d_2) \\ d_1 &= \frac{\log(\frac{x_0}{K}) + \left(r + \frac{\sigma^2}{2}\right)T}{\sigma\sqrt{T}}, \quad d_2 = d_1 - \sigma\sqrt{T}. \end{aligned}$$

• Implicitation of the volatility: solve in σ the inverse problem

$$\operatorname{Call}_{BS}(\ldots,\sigma,\ldots)-\operatorname{Call}_{M2Mkt}=0.$$





• The function is even in σ and the equation has two opposite solutions.

• Algo₁:

$$\sigma_{n+1} = \sigma_n - \underbrace{\gamma_{n+1}(\operatorname{Call}_{BS}(\sigma_n) - \operatorname{Call}_{M2Mkt})}_{=:h(\sigma_n)}, \quad \sigma_0 > 0.$$

with $\gamma_n = \gamma > 0$ or decreasing assumption.

- Algo₂:
 - The Vega:

$$\operatorname{Vega}_{BS}(\sigma) = x_0 \operatorname{sign}(\sigma) \sqrt{T} \frac{e^{-\frac{d_1(\sigma)^2}{2}}}{\sqrt{2\pi}}$$

• Implicit volatility search reads:

$$\sigma_{n+1} = \sigma_n - \underbrace{\frac{\operatorname{Call}_{BS}(x_0, K, r, \sigma_n) - \operatorname{Call}_{M2Mkt}}_{\operatorname{Vega}_{BS}(\sigma_n)}}_{=:h(\sigma_n)}, \quad \sigma_0 > 0.$$

[This is the actual algorithm with a "good choice" of σ_0]

Implicitation: Correlation I

• 2-dim (correlated) Black-Scholes model:

$$X_t^i = x_0^i e^{(r - \frac{\sigma_i^2}{2})t + \sigma_i W_t^i}, \ x_0^i, \ \sigma_i > 0, i = 1, 2$$

with $\langle W^1, W^2 \rangle_t = \rho t$.

• Best-of-Call Payoff:

$$\left(\max(X_{\tau}^1,X_{\tau}^1)-K\right)_+$$

• Premium at time 0

$$\mathsf{Best-of-Call}_{BS}(\dots,\rho,\dots) = e^{-rT}\mathbb{E}\left(\,\mathsf{max}(X^1_\tau,X^1_\tau) - \mathcal{K}\right)_+$$

- Organized markets on such options are market of the correlation ρ .
- The volatilities σ_i , i = 1, 2, are known from vanilla option markets on X^1 and X^2 .

How to "extract" the correlation ρ ?

• Deterministic algo(s):

$$\rho_{n+1} = \rho_n - \gamma_{n+1} \underbrace{\frac{\mathsf{Best-of-Call}_{BS}(\rho_n) - \mathsf{Best-of-Call}_{M2Mkt}}{\partial_{\rho}\mathsf{Best-of-Call}_{BS}(\rho_n) + \lambda_n}_{=:h(\rho_n)}}_{=:h(\rho_n)}.$$

- Except that we have no (simple) closed form for the B-S price and its ρ -derivative
- The correlation $ho \in [-1,1]$. Projections are possible but....
- What to do?



Optimization (deterministic, the origins)

Examples from Finance

- Implicitation
- Minimization
- Learning procedure
 - Abstract Learning
 - Supervised Learning
 - Unsupervised Learning (clustering)
- 4 Stochastic algorithms/Approximation
 - From Robbins-Monro to Robbins-Siegmund
 - Stochastic Gradient Descent (SGD) and pseudo-SGD
- 5 Examples revisited by SFD
 - Numerical probability
 - Learning (supervised and unsupervised)
- 6 Application to Neural Networks and deep learning
 - Linear neural network
 - One hidden layer feedforward perceptron
 - Toward deep learning
 - Multilayer feedforward perceptron and Backpropagation

Minimization: Value-at-risk/Conditional Value-at-risk/I

Let X = φ(Z), Z : (Ω, A, P) → R^q be an integrable random variable representative of a loss and let α ∈ (0, 1), α ≃ 1.

 $\mathsf{Value-at-Risk}_{\alpha}(X) = \alpha \text{-quantile} = \inf \left\{ \xi : \mathbb{P}(X \leq \xi) \geq \alpha \right\}.$

• For simplicity, assume X has a density $f_X > 0$ on \mathbb{R} . Then $\xi_{\alpha} = \operatorname{VaR}_{\alpha}(X)$ is the unique solution to

$$\mathbb{P}(X \leq \xi_{\boldsymbol{\alpha}}) = \boldsymbol{\alpha} \Longleftrightarrow \mathbb{P}(X > \xi_{\boldsymbol{\alpha}}) = 1 - \boldsymbol{\alpha}.$$

• The conditional Value-at-Risk is defined by

$$\mathsf{CVaR}_{\boldsymbol{\alpha}}(X) = \mathbb{E}(X \mid X \ge \mathsf{VaR}_{\boldsymbol{\alpha}}(X)).$$

• Rockafellar-Uryasev Potential (¹):

$$V(\xi)=\xi+rac{1}{1-lpha}\mathbb{E}\,(X-\xi)_+,\quad \xi\in\mathbb{R}.$$

Gilles PAGÈS (LPSM)

17 / 94

¹ R.T. Rockafellar, S. Uryasev (2000). Optimization of Conditional Value-At-Risk, The Journal of Risk, 2(3):21-41. www.ise.ufl.edu/uryasev.

• The function V is convex and $\lim_{|\xi| \to +\infty} V(\xi) = +\infty$ since

$$V(\xi) \geq \xi$$
 so that $\lim_{\xi o +\infty} V(\xi) = +\infty$

and

$$V(\xi) \ge \xi + \frac{1}{1-\alpha} (\mathbb{E} X - \xi)_+$$

= $\xi + \frac{1}{1-\alpha} (\mathbb{E} X - \xi)$ for ξ low enough
= $-\frac{\alpha}{1-\alpha} \xi + \frac{1}{1-\alpha} \mathbb{E} X \to +\infty$ as $\xi \to -\infty$.

 $\bullet\,$ By differentiation under the $\mathbb E,$ we get

$$V'(\xi)=1-rac{1}{1-lpha}\mathbb{P}(X>\xi).$$

• $V'(\xi) = 0$ iff $\mathbb{P}(X > \xi) = 1 - \alpha$ iff $\xi = \xi_{\alpha}$

Moreover

$$V(\xi_{\alpha}) = \frac{\xi_{\alpha} + \mathbb{E} (X - \xi_{\alpha})_{+}}{\mathbb{P}(X > \xi_{\alpha})} = \frac{\mathbb{E} X \mathbf{1}_{\{X > \xi_{\alpha}\}}}{\mathbb{P}(X \ge \xi_{\alpha})}$$
$$= \mathbb{E} (X \mid X \ge \mathsf{VaR}_{\alpha}(X)) = \mathsf{CVaR}_{\alpha}(X)$$

• (GD) pour la Va $\mathsf{R}_{lpha}(X)$: $h(\xi) = V'(\xi)$. Let $\xi_0 \in \mathbb{R}$,

$$\begin{aligned} \xi_{n+1} &= \xi_n - \gamma_{n+1} \big(1 - \frac{1}{1 - \alpha} \big(1 - F_x(\xi_n) \big) \big) \\ &= \xi_n - \frac{\gamma_{n+1}}{1 - \alpha} \big(F_x(\xi_n) - \alpha \big), \quad n \ge 0 \end{aligned}$$

• Newton/Levenberg-Marquardt algo: $\xi_0 \in \mathbb{R}$,

$$\xi_{n+1} = \xi_n - \frac{F_x(\xi_n) - \alpha}{(1-\alpha)f_x(\xi_n)}, \quad n \ge 0.$$

 Why not ! But X = φ(Z) (the whole portfolio of a CIB Bank!) ⇒ q large and no closed form for the c.d.f. F_x(ξ) = ℙ(X ≤ ξ) of X.

- Optimization (deterministic, the origins)
- Examples from Finance
- Implicitation
- Minimization
- 3 Learning procedures
 - Abstract Learning
 - Supervised Learning
 - Unsupervised Learning (clustering)
 - Stochastic algorithms/Approximation
 - From Robbins-Monro to Robbins-Siegmund
 - Stochastic Gradient Descent (SGD) and pseudo-SGD
- 5 Examples revisited by SFD
 - Numerical probability
 - Learning (supervised and unsupervised)
- 6 Application to Neural Networks and deep learning
 - Linear neural network
 - One hidden layer feedforward perceptron
 - Toward deep learning
 - Multilayer feedforward perceptron and Backpropagation



Optimization (deterministic, the origins)

- **Examples from Finance**
- Implicitation
- Minimization
- 3 Lear

Learning procedures

Abstract Learning

- Supervised Learning
- Unsupervised Learning (clustering)
- Stochastic algorithms/Approximation
 - From Robbins-Monro to Robbins-Siegmund
 - Stochastic Gradient Descent (SGD) and pseudo-SGD
- 5 Examples revisited by SFD
 - Numerical probability
 - Learning (supervised and unsupervised)
- 6 Application to Neural Networks and deep learning
 - Linear neural network
 - One hidden layer feedforward perceptron
 - Toward deep learning
 - Multilayer feedforward perceptron and Backpropagation

Abstract Learning

- Huge dataset $(z_k)_{k=1:N}$ with of possibly high dimension $d: N \simeq 10^6$, even 10^9 , and $d \simeq 10^3$. [Image, profile, text, ...]
- Set of parameters $\theta \in \Theta \subset \mathbb{R}^{K}$, K large (see later on).
- There exists a smooth local loss function/local predictor

 $v(\theta, z).$

• Global loss function:
$$V(\theta) = \frac{1}{N} \sum_{k=1}^{N} v(\theta, z_k)$$

with gradient $\nabla V(\theta) = \frac{1}{N} \sum_{k=1}^{N} \nabla_{\theta} v(\theta, z_k).$

• Solving the minimization problem

 $\min_{\theta\in\Theta}V(\theta).$

• Suggests a (GD) i.e. $h = \nabla V$ [or others... if $\nabla^2_{\theta} v(\theta, z)$ exists]:

$$egin{aligned} & heta_{n+1} &= heta_n - \gamma_{n+1}
abla V(heta_n) \ & = heta_n - rac{\gamma_{n+1}}{N} \sum_{k=1}^N
abla_ heta v(heta, z_k), \ n \geq 0, \end{aligned}$$

with the step sequence satisfying the (DS) assumption.



Supervised learning

- Input x_k , output y_k . Data $z_k = (x_k, y_k) \in \mathbb{R}^{d_x + d_y}$, k = 1 : n.
- Transfer function $f: \Theta \times \mathbb{R}^{d_x} \to \mathbb{R}^{d_y}$



• Prediction/loss function (local) $v(\theta, z) = \frac{1}{2} |f(\theta, x_k) - y_k|^2$, k = 1 : n so that

$$abla_{ heta} \mathbf{v}(\theta, \mathbf{z}) =
abla_{ heta} f(\theta, \mathbf{x})^{ op} \big(f(\theta, \mathbf{x}) - \mathbf{y} \big).$$

Resulting loss function gradient

$$V(\theta) = rac{1}{N} \sum_{k=1}^{N}
abla_{ heta} f(\theta, x_k)^{ op} (f(\theta, x_k) - y_k).$$

Optimization (deterministic, the origins)

- **Examples from Finance**
- Implicitation
- Minimization

Learning procedures

- Abstract Learning
- Supervised Learning

Unsupervised Learning (clustering)

- Stochastic algorithms/Approximation
 - From Robbins-Monro to Robbins-Siegmund
 - Stochastic Gradient Descent (SGD) and pseudo-SGD
- 5 Examples revisited by SFD
 - Numerical probability
 - Learning (supervised and unsupervised)
- 6 Application to Neural Networks and deep learning
 - Linear neural network
 - One hidden layer feedforward perceptron
 - Toward deep learning
 - Multilayer feedforward perceptron and Backpropagation

Unsupervised learning (clustering)

• Only input
$$z_k = x_k \in \mathbb{R}^d$$
, $k = 1 : N$.

- Prototype parameter set: $(\theta^1, \ldots, \theta^r) \in \Theta = (\mathbb{R}^d)^r$, $r \in \mathbb{N}$.
- Local loss function: $x \in \mathbb{R}^d$, $\theta \in \Theta$.

$$v(\theta, x) = \frac{1}{2} \min_{i=1:r} |\theta^i - x|^2 = \frac{1}{2} \operatorname{dist} \left(x, \{\theta^1, \dots, \theta^r\} \right)^2$$

(minimal distance to prototypes).

- $v(\theta, x)$ is not convex in θ !
- Global loss function (Distortion):

$$V(\theta) = \frac{1}{2} \sum_{k=1}^{N} \min_{i=1:r} |\theta^i - x_k|^2$$
 (mean minimal distance to prototypes).

94

Batch k-means/Forgy's algorithm

• Gradient at
$$\theta$$
 s.t. $\theta^{i} \neq \theta^{j}$: $\nabla V(\theta) = \frac{1}{2} \sum_{k=1}^{N} \nabla_{\theta} v(\theta, x_{k}).$

with, for i = 1 : r,

$$\partial_{\theta^i} v(\theta, x_k) = (\theta^i - x_k) \mathbf{1}_{\{|x_k - \theta^i| < \min_{j \neq i} |x_k - \theta^j|\}}.$$

^ /

• $\mathbf{1}_{\{|x_k - \theta^i| < \min_{j \neq i} | x_k - \theta^j|\}}$ = nearest neighbour search.

- Compute $\nabla V(\theta) = \frac{1}{N} \sum_{k=1}^{N} \nabla_{\theta} v(\theta, x_k)$
- \implies *N*·nearest neighbour searches among *r* ptotypes of dim *d*!
- Forgy's algorithm

$$\theta_{n+1} = \theta_n - \gamma_{n+1} \nabla V(\theta_n)$$



Gilles PAGÈS (LPSM)

Stochastic approximation I

Pros and cons: toward stochastic algorithm I

- Numerical Probability (for Finance): we do not know how to compute $h(\theta)$.
 - *h* always has a probabilistic presentation in our examples:

$$h(\theta) = \mathbb{E} H(\theta, Z) = \int_{\mathbb{R}^q} H(\theta, z) \mathbb{P}_z(dz) = \int_{\mathbb{R}^q} H(\theta, z) f_z(z) dz$$

where $H: \mathbb{R}^d \times \mathbb{R}^q \to \mathbb{R}^d$ is Borel, q often large...

- Cons: ... which requires the computation of (often) high dimensional integrals on \mathbb{R}^q at a reasonable cost: impossible.
- Pros: The random vector Z can be simulated.
- Pros: The function *H* is computable at a reasonable computational cost.
- Pros: Regularizing effect of \mathbb{E} : *h* smoother than the functions H(.,z). (Think to $F_X(\xi) = \mathbb{E}\mathbf{1}_{\{X \le \xi\}}$.)

Pros and Cons: toward stochastic algorithm I

- Data Science (usually V is given and h = ∇V): but we cannot compute h(θ).
 - *h* still has probabilistic representation using the empirical measure.

$$h(heta) = rac{1}{N}\sum_{k=1}^{N}
abla_{ heta} v(heta, z_k) = \int_{\mathbb{R}^q}
abla_{ heta} v(heta, z_k) \mu_{ heta}(dz) ext{ with } \mu_{ heta} = rac{1}{N}\sum_{k=1}^{N} \delta_{z_k}$$

• Cons: But N huge $\implies h(\theta)$ cannot be computed at a reasonable cost.

• Pros:

$$h(\theta) = \mathbb{E} \nabla_{\theta} v(\theta, Z)$$

where Z can be simulated by picking up a datum (uniformly) at random since

$$Z \sim z_I, \quad I \sim \mathcal{U}(\{1, \ldots, N\}).$$

- $v(\theta, z)$ and $\nabla_{\theta} v(\theta, z)$ both computable hence V and $h = \nabla V$ too.
- Cons: No regularizing effect of E: smoothness of [h = ∇V]= smoothness of H(., z).

• Cons: Transfer of convexity in θ from $v(\cdot, z)$ to V.

Toward stochastic algorithm II

- Zero search of $h(\theta) = \mathbb{E} H(\theta, Z)$ as above.
- Idea 1: Use Monte Carlo simulation

$$\theta_{n+1} = \theta_n - \gamma_{n+1} \widehat{h}_{M_{n+1}}(\theta_n)$$
$$\widehat{h}_{M_{n+1}}(\theta_n) = \frac{1}{M_n} \sum_{k=1}^{M_n} H(\theta_n, Z_k^{(n+1)}), \quad (Z_k^{(n+1)})_{k,n} \ i.i.d. \sim Z$$

• Idea 2: Robbins-Monro, 1951 (²)

Set
$$\forall n \geq 1$$
, $M_n = 1 !!$

• Idea 1.5: Mini-batch i.e. $M_n = M > 2$. Successful among practitioners.

²H. Robbins, S. Monro (1951). A stochastic approximation method, Ann. Math. Stat., 22:400-407.



Robbins-Monro framework (1951)

 \triangleright Pactitioner's corner: Replace $h(\theta_n)$ by a $H(\theta_n, Z_{n+1})$.

Let (θ_n)_{n≥0} be a sequence of ℝ^d-valued random vectors recursively defined on (Ω, A, ℙ) by

$$\theta_{n+1} = \theta_n - \gamma_{n+1} H(\theta_n, Z_{n+1}), \quad \theta_0 \in L^2(\mathbb{P}, \mathcal{A})$$

with

(i)
$$(Z_n)_{n\geq 1}$$
 is i.i.d. $\sim Z$, independent of θ_0
(ii) $||H(\theta, Z)||_2 \leq C(1 + |\theta|)$ ($\Rightarrow h$ linear growth)
(iii) $(\gamma_n)_{n\geq 1}$ is a $(0, +\infty)$ -valued deterministic step sequence
so that $(\theta_n)_{n\geq 0}$ is (\mathcal{F}_n) -adapted with $\mathcal{F}_n = \sigma(\theta_0, Z_1, \dots, Z_n)$. Then
 $\theta_{n+1} = \theta_n - \gamma_{n+1}h(\theta_n) + \gamma_{n+1}\Delta M_{n+1}, n \geq 0$.

• The key is simply: $\Delta M_{n+1} = H(\theta_n, Z_{n+1}) - h(\theta_n)$ since

$$\mathbb{E}\big(H(\theta_n, Z_{n+1}) \,|\, \mathcal{F}_n\big) \underset{Z_{n+1} \perp \perp \mathcal{F}_n}{=} \Big[\mathbb{E} \,H(\theta, Z_{n+1})\Big]_{|\theta=\theta_n} = h(\theta_n), \ n \geq 0.$$

- Idea 1 (Robbins-Monro 1951, Robbins-Siegmund 1971):
 - Perturbed zero search procedure with decreasing step for *h*.
 - The perturbation is a martingale increment.
- Idea 2 (Ljung, 1977): Perturbed Euler scheme with decreasing step of the ODE

$$\dot{\theta} = -h(\theta).$$

(not exploited here)
Idea 1: Robbins-Siegmund Lemma, 1971

Theorem (Robbins-Siegmund Lemma, 1971)

Lyapunov function: V : ℝ^d → ℝ₊, C¹, lim_∞ V = +∞, ∇V Lipschitz, |∇V|² ≤ c(1 + V) and (mean-reversion) (∇V|h) ≥ 0 and ||H(θ, Z)||₂ ≤ C√(1 + V(θ)).
Decreasing Step assumption (DS): Σ_n γ_n = +∞ and Σ_n γ_n² < +∞.
V(θ₀) ∈ L¹.

Then

(i)
$$V(\theta_n) \xrightarrow{a.s.} V_{\infty} \in L^1 \ [\Rightarrow (\theta_n)_{n\geq 0} \ pathwise bounded] \ and \ L^1-bounded.$$

(ii) $\sum_n \gamma_n (\nabla V|h)(\theta_{n-1}) \in L^1 \ a.s.$, hence $< +\infty \ a.s.$
(iii) $\sum_n |\Delta \theta_n|^2 < +\infty a.s.$ (so that $\theta_n - \theta_{n-1} \to 0 \ a.s.$).
(iv) $\sum_n \gamma_n \Delta M_n$ converges a.s. and in L^2 .

• Note that V is sub-quadratic i.e. $V(\theta) \le \kappa (1 + |\theta|^2)$ and h is sunlinear. Gilles PAGÈS (LPSM) Stochastic approximation I LPSM-Sorbonne Université

36 / 94

Proof

• Set $\mathcal{F}_n := \sigma(\theta_0, Z_1, \dots, Z_n)$, $n \ge 1$ and $\Delta \theta_n := \theta_n - \theta_{n-1}$, $n \ge 1$.

• There exists $\xi_{n+1} \in (\theta_n, \theta_{n+1})$ s.t.

$$V(\theta_{n+1}) = V(\theta_n) + (\nabla V(\xi_{n+1}) | \Delta \theta_{n+1})$$

$$\leq V(\theta_n) + (\nabla V(\theta_n) | \Delta \theta_{n+1}) + [\nabla V]_{\text{Lip}} | \Delta \theta_{n+1} |^2$$

$$= V(\theta_n) - \gamma_{n+1} (\nabla V(\theta_n) | H(\theta_n, Z_{n+1}))$$

$$+ [\nabla V]_{\text{Lip}} \gamma_{n+1}^2 | H(\theta_n, Z_{n+1}) |^2$$

$$= V(\theta_n) - \gamma_{n+1} (\nabla V(\theta_n) | h(\theta_n)) - \gamma_{n+1} (\nabla V(\theta_n) | \Delta M_{n+1})$$

$$+ [\nabla V]_{\text{Lip}} \gamma_{n+1}^2 | H(\theta_n, Z_{n+1}) |^2,$$
(*)

where

$$\Delta M_{n+1} = H(\theta_n, Z_{n+1}) - h(\theta_n).$$

Show by induction that V(θ_n)∈ L¹(ℙ), given that V(θ₀)∈ L¹(ℙ).
Key:

$$\mathbb{E}\left|\left(\nabla V(\theta_n)|H(\theta_n, Z_{n+1})\right)\right| \leq \frac{1}{2} \left(\mathbb{E}\left|\nabla V(\theta_n)\right|^2 + \mathbb{E}\left|H(\theta_n, Z_{n+1})\right|^2\right).$$

• So that $(\Delta M_n)_{n\geq 1}$ is a sequence of $L^2(\mathcal{F}_n)$ -martingale increments satisfying

 $\mathbb{E}\left(\left|\Delta M_{n+1}\right|^{2} | \mathcal{F}_{n}\right) \leq \mathbb{E}\left(\left|H(\theta_{n}, Z_{n+1})\right|^{2} | \mathcal{F}_{n}\right) \leq C(1 + V(\theta_{n})).$

• Coming back to

$$egin{aligned} & V(heta_{n+1}) \leq V(heta_n) - \gamma_{n+1}(
abla V(heta_n)) - \gamma_{n+1}(
abla V(heta_n)|\Delta M_{n+1}) \ & + [
abla V]_{ ext{Lip}} \gamma_{n+1}^2 |H(heta_n, Z_{n+1})|^2, \end{aligned}$$

• conditioning given \mathcal{F}_n yields

•
$$\mathbb{E}[(\nabla V(\theta_n)|\Delta M_{n+1})|\mathcal{F}_n]=0$$

• and, other terms in the RHS being \mathcal{F}_n -measurable,

$$\begin{split} \mathbb{E}\left(V(\theta_{n+1}) \,|\, \mathcal{F}_n\right) + \gamma_{n+1}(\nabla V |h)(\theta_n) &\leq V(\theta_n) + C_{_V} \gamma_{n+1}^2 \big(1 + V(\theta_n)\big) \\ &= V(\theta_n)(1 + C_{_V} \gamma_{n+1}^2) + C_{_V} \gamma_{n+1}^2 \end{split}$$

with $C_V = C^2 [\nabla V]_{\text{Lip}} > 0.$

Add

- the positive term $\sum_{k=1}^{n} \gamma_k (\nabla V | h)(\theta_{k-1}) + C_v \sum_{k \ge n+2} \gamma_k^2$ on the left-hand side of the above inequality,
- $(1 + C_v \gamma_{n+1}^2)$ times this term on the right-hand side .

• Divide the resulting inequality by $\prod_{k=1}^{n+1} (1 + C_V \gamma_k^2)$ shows that (the \mathcal{F}_n -adapted sequence)

$$S_{n} = \frac{V(\theta_{n}) + \sum_{k=0}^{n-1} \gamma_{k+1}(\nabla V|h)(\theta_{k}) + C_{v} \sum_{k \ge n+1} \gamma_{k}^{2}}{\prod_{k=1}^{n} (1 + C_{v} \gamma_{k}^{2})}, \ n \ge 1,$$

is a (non-negative) super-martingale with $S_0 = V(\theta_0) \in L^1(\mathbb{P})$.

- The fact that the added term is positive follows from the mean-reverting inequality (∇V|h) ≥ 0.
- Hence

$$S_n \xrightarrow{a.s.} S_\infty \in L^1_{\mathbb{R}_+}(\mathbb{P}).$$

• Consequently, using that $\sum_{k\geq n+1}\gamma_k^2\rightarrow 0,$ one gets

$$V(heta_n) + \sum_{k=0}^{n-1} \gamma_{k+1}(
abla V|h)(heta_k) \xrightarrow{a.s.} \widetilde{S}_{\infty} = S_{\infty} \prod_{n \geq 1} (1 + C_V \gamma_n^2) \in L^1(\mathbb{P}).$$

• (*i*)_a The super-martingale

 $(S_n)_{n\geq 0}$ is $L^1(\mathbb{P})$ -bounded by $\mathbb{E} S_0 = \mathbb{E} V(\theta_0) < +\infty$,

hence $(V(\theta_n))_{n\geq 0}$ is L¹-bounded since

$$V(\theta_n) \leq \left(\prod_{k=1}^n (1+C_V\gamma_k^2)\right) S_n, \quad n \geq 0,$$

and $\prod_{k\geq 1} (1 + C_V \gamma_k^2) < +\infty$ by the (*DS*) assumption on $(\gamma_n)_{n\geq 1}$.

• (ii) Now, for the same reason, the series with non-negative terms $\sum_{0 \le k \le n-1} \gamma_{k+1} (\nabla V | h)(\theta_k)$ satisfies for every $n \ge 1$,

$$\mathbb{E}\left(\sum_{k=0}^{n-1}\gamma_{k+1}(\nabla V|h)(\theta_k)\right) \leq \prod_{k=1}^n (1+C_V\gamma_k^2)\mathbb{E}S_0$$

so that, by the Beppo Levi monotone convergence Theorem for series with non-negative terms,

$$\mathbb{E}\left(\sum_{n\geq 0}\gamma_{n+1}(\nabla V|h)(\theta_n)\right)<+\infty$$

so that, in particular,

$$\sum_{n\geq 0}\gamma_{n+1}(\nabla V|h)(\theta_n)<+\infty\qquad\mathbb{P}\text{-}a.s.$$

and the series converges in L^1 to its *a.s.* limit.

- (i)_b It follows that $V(\theta_n) \longrightarrow V_{\infty}$ a.s. as $n \to +\infty$. $V_{\infty} \in L^1$ by Fatou's Lemma since $(V(\theta_n))_{n\geq 0}$ is L^1 -bounded.
- (*iii*) Again by Beppo Levi's monotone convergence Theorem for series with non-negative terms,

$$\mathbb{E}\left(\sum_{n\geq 1} |\Delta\theta_n|^2\right) = \sum_{n\geq 1} \mathbb{E} |\Delta\theta_n|^2 \leq \sum_{n\geq 1} \gamma_n^2 \mathbb{E} |H(\theta_{n-1}, Z_n)|^2$$
$$\leq C \sum_{n\geq 1} \gamma_n^2 (1 + \mathbb{E} V(\theta_{n-1})) < +\infty$$

so that

$$\sum_{n\geq 1} |\Delta heta_n|^2 \in L^1(\mathbb{P})$$
 (hence *as*. finite)

which in turns yields

$$\Delta \theta_n = \theta_n - \theta_{n-1} \to 0$$
 a.s. and in $L^2(\mathbb{P})$

• (iv) We have $M_n^{\gamma} = \sum_{k=1}^n \gamma_k \Delta M_k$ so that M^{γ} is clearly an (\mathcal{F}_n) -martingale. Moreover,

$$\begin{split} \langle \mathcal{M}^{\gamma} \rangle_{n} &= \sum_{k=1}^{n} \gamma_{k}^{2} \mathbb{E} \left(|\Delta \mathcal{M}_{k}|^{2} |\mathcal{F}_{k-1} \right) \leq \sum_{k=1}^{n} \gamma_{k}^{2} \mathbb{E} \left(|\mathcal{H}(\theta_{k-1}, Z_{k})|^{2} |\mathcal{F}_{k-1} \right) \\ &\leq C \sum_{k=1}^{n} \gamma_{k}^{2} \left(1 + \mathbb{E} V(\theta_{k-1}) \right) \end{split}$$

Consequently, owing to $(i)_a$,

$$\mathbb{E} \langle M^{\gamma} \rangle_{\infty} < +\infty,$$

Hence M_n^{γ} converges *a.s.* and in L^2 .

Table of Contents



Theorem (Robbins-Monto algorithm)

Assume the mean function h is continuous and satisfies

$$\forall \theta \in \mathbb{R}^d, \ \theta \neq heta_*, \quad (heta - heta_* | h(heta)) > 0.$$

Suppose furthermore that $\theta_0 \in L^2$ and that H satisfies

$$orall heta \in \mathbb{R}^d, \quad ig\| H(heta, Z) ig\|_2 \leq C(1+| heta|).$$

Finally, assume $(\gamma_n)_{n\geq 1}$ satisfies (DS). Then

$$\{h=0\} = \{\theta_*\}$$
 and $\theta_n \xrightarrow{a.s.} \theta_*$.

The convergence also holds in every L^p , $p \in (0,2)$ (and $(|\theta_n - \theta_*|)_{n \ge 0}$ is L^2 -bounded).

• If $H(\theta, z) = h(\theta)$, back to a deterministic zero search procedure!!!

- The function $V(\theta) = \frac{1}{2}|\theta \theta_*|^2$ is a Lyapunov function.
- The quadratic linear growth assumption on H is satisfied too.
- Robbins-Siegmund's Lemma implies

•
$$|\theta_n - \theta_*|^2 \longrightarrow V_\infty \in L^1$$
,
• $\sum_{n \ge 1} \gamma_n (h(\theta_{n-1})|\theta_{n-1} - \theta_*) < +\infty$ P-a.s

•
$$(|\theta_n - \theta_*|^2)_{n \ge 0}$$
 is L^1 -bounded.

• We keep on reasoning pathwise: let ω be generic.

On has

$$\lim_{n} \left(\theta_{n-1}(\omega) - \theta_* | h(\theta_{n-1}(\omega)) \right) = 0.$$

• If $\underline{\lim_{n}} \left(\theta_{n-1}(\omega) - \theta_* | h(\theta_{n-1}(\omega)) \right) > 0$, the above convergence induces

a contradiction with
$$\sum_{n\geq 1}\gamma_n=+\infty$$
.

• Let $\left(\phi(n,\omega)\right)_{n\geq 1}$ be a subsequence such that

$$ig(heta_{\phi(n,\omega)}(\omega)- heta_*|h(heta_{\phi(n,\omega)}(\omega))ig)\longrightarrow 0 \quad ext{ as } \quad n o+\infty.$$

Now, (θ_n(ω))_{n≥0} being bounded, one may assume, up to one further extraction,

$$heta_{\phi(n,\omega)}(\omega) o heta_{\infty} = heta_{\infty}(\omega).$$

• By continuity of h, $(\theta_{\infty} - \theta_* | h(\theta_{\infty})) = 0$ which implies $\theta_{\infty} = \theta_*$. Now, since we know that $V(\theta_n(\omega)) = \frac{1}{2} |\theta_n(\omega) - \theta_*|^2$ converges,

$$\lim_{n} |\theta_{n}(\omega) - \theta_{*}|^{2} = \lim_{n} |\theta_{\phi(n,\omega)}(\omega) - \theta_{*}|^{2} = 0.$$

• Convergence in L^p , $p \in (0,2)$ follows by uniform integrability.

Theorem (Stochastic Gradient Descent)

↓ Let V : ℝ^d → ℝ₊ be a differentiable function lim_∞ V(θ) = +∞, ∇V
Lipschitz, |∇V|² ≤ C(1 + V) and {∇V = 0} = {θ_{*}}.
↓ Let h(θ) = ℝ H(θ, Z) = ∇V with H s.t. ||H(θ, Z)||₂ ≤ C√(1 + V(θ)) and that V(θ₀) ∈ L¹(ℙ). Assume (γ_n)_{n≥1} satisfies (DS).

Then

$$V(heta_*) = \min_{\mathbb{R}^d} V \quad \textit{and} \quad heta_n \stackrel{a.s.}{\longrightarrow} heta_* \quad \textit{as} \quad n \to +\infty.$$

Moreover, $\nabla V(\theta_n)$ converges to 0 in every L^p , $p \in (0,2)$ (and $(V(\theta_n))_{n \ge 0}$ is L^1 -bounded so that $(\nabla V(\theta_n))_{n \ge 0}$ is L^2 -bounded).

- *Proof.* Use (almost) the same arguments as above but with $(\nabla V \mid h) = |\nabla V|^2$ instead of $(\theta \theta_* \mid h(\theta))$.
- If $H(\theta, z) = h(\theta) = \nabla V(\theta)$: Convergence thm for Gradient descent!!

Theorem (Multitarget Stochastic Gradient Descent)

(a) If the former assumption $\{\nabla V = 0\} = \{\theta_*\}$, one has mutatis mutandis: a.s. there exists $v_{\infty} \in \mathbb{R}_+$ and a connected component χ_{∞} of $\{\nabla V = 0\} \cap \{V = v_{\infty}\}$ such that

$$\operatorname{dist}(\theta_n, \chi_\infty) \longrightarrow 0 \quad a.s.$$

(b) In particular if $\{\nabla V = 0\} \cap \{V = v\}$ is locally finite for every $v \ge 0$ is finite, then there exists a r.v. θ_{∞} such that

$$abla V(heta_\infty) = 0 \quad \textit{ and } \quad heta_n \longrightarrow heta_\infty.$$

(c) Moreover, $\nabla V(\theta_n)$ converges to 0 in every L^p , $p \in (0,2)$ (and $(V(\theta_n))_{n\geq 0}$ is L^1 -bounded so that $(\nabla V(\theta_n))_{n\geq 0}$ is L^2 -bounded).

- By the R.-S. Lemma, one has for free that, a.s., (θ_n)_{n≥0} is pathwise bounded and θ_n − θ_{n−1} → 0 pathwise. Hence its limiting values makes up a connected compact set Θ_∞, clearly included in some {V = v_∞}.
- $\bullet\,$ But this is not enough. . . Θ_∞ is also invariant under the flow of

$$ODE \equiv \dot{ heta} = -h(heta)$$

which converges toward $\{\nabla V = 0\}$.

- Still not enough : needs to make a transfer from ODE to algorithm.
- Needs further insights based on topology and the ODE method $(^{3})$.

³G. Pagès (2018). Introduction to Numerical Probability with application to Finance, Springer-Verlag, Berlin, 576p.

Theorem (Traps (Pemantle 1984, Lazarev 1989, Brandière-Duflo 1996, Fort-P. 1997, Benaïm 1998s))

Let $\theta_* \in \{\nabla V = 0\}$. If there exists (λ, u) such that $D^2 V u = \lambda u$ such that

$$\lambda < 0$$
 and $\mathbb{E}\left(H(\theta_*, Z)|u\right)^2 > 0$

then

$$\mathbb{P}(\theta_n \to \theta_*) = 0.$$

 This allows to eliminate noisy local maxima, saddle points, monkey saddle points, etc.

Table of Contents



• Multilayer feedforward perceptron and Backpropagation

Table of Contents



• Multilayer feedforward perceptron and Backpropagation

Numerical probability: Implicit volatility II

• Let (4)
$$h(\sigma) = \operatorname{Call}_{BS}(\sigma) - \operatorname{Call}_{M2Mkt} = \mathbb{E}H(\sigma, Z)$$
 with

$$H(\theta, z) = \left(x_0 e^{-\frac{\sigma_+^2}{2}T + \sigma_+\sqrt{T}z} - e^{-rT}K\right)_+ - \operatorname{Call}_{M2Mkt}$$

(σ_+ to ensure that *h* is increasing).

Then the recursive stochastic zero search reads

$$\sigma_{n+1} = \sigma_n - \gamma_{n+1} H(\sigma_n, Z_{n+1}), \quad \sigma_0 > 0.$$

with $(Z_n)_{n\geq 1}$ i.i.d., $\sim \mathcal{N}(0,1)$ and $\sum_n \gamma_n = +\infty$, $\sum_{n\geq 1} \gamma_n^2 < +\infty$

• Try with
$$\gamma_n = rac{a}{b+n}$$
 so that $\gamma_1 imes H(\sigma_0, Z_{+1}) \simeq$ few units.

Gilles PAGÈS (LPSM)

55 / 94

⁴G. Pagès (2018). Introduction to Numerical Probability with application to Finance, Springer-Verlag, Berlin, 576p.

Numerical probability: correlation search II

• Let
$$h(\rho) = \text{Best-ofCall}_{BS}(\dots,\rho,\dots) - \text{Best-of-Call}_{M2Mkt}$$

 $= \mathbb{E} H(\rho, Z), \ Z = (Z^1, Z^2) \sim \mathcal{N}(0, I_2)$
with $H(\rho, Z) = \left(\max\left(x_0^1 e^{-\frac{\sigma_1^2 T}{2} + \sigma_1 \sqrt{T} z^1}, x_0^2 e^{-\frac{\sigma_2^2 T}{2} + \sqrt{T} \sigma_2(\rho z^1 + \sqrt{1 - \rho^2} z^2)} - e^{-rT} \kappa\right)_+ - \text{Best-of-Call}_{M2Mkt}.$

• The naive algorithm (with $(\gamma_n)_{n\geq 1}$ satisfying the (DS) assumption) $\rho_{n+1} = \rho_n - \gamma_{n+1} H(\rho_n, Z_{n+1})$

does not live inside $[-1,1] \ensuremath{\,!\!\!\!\!!} \hdots$. . .

• What to do ? Project on [-1,1] (theorems do exist) or change of variable (⁵)

$$\rho = \frac{2}{\pi} \arctan(\theta) =: \varphi(\theta)$$

so that

$$\theta_{n+1} = \theta_n - \gamma_{n+1} H(\varphi(\theta_n), Z_{n+1}), \ \theta_0 \in \mathbb{R}.$$

56 / 94

• There exists a C^1_{Lip} -Lyapunov function $V : \mathbb{R} \to \mathbb{R}_+$ such that $V'(\theta)h(\varphi(\theta))) \ge 0$ and = 0 iff $h(\varphi(\theta)) = 0$ i.e. $\varphi(\theta) = \rho_*$.

⁵G. Pagès (2018). Introduction to Numerical Probability with application to Finance, Springer-Verlag, Berlin, 576p. Gilles PAGÈS (LPSM) Stochastic approximation I LPSM-Sorbonne Université

Higher dimensions

• In higher dimension a correlation matrix *R* whose Cholesky decomposition

$$R = TT^{\top}$$
 with T lower triangular and

$$\sum_{j=1}^i t_{ij}^2 = 1.$$

• (Hyper-)spherical parametrization

$$\begin{split} t_{11} &= 1 \\ t_{21} &= \cos(\theta_2), \ t_{22} &= \sin(\theta_2) \\ t_{31} &= \cos(\theta_3)\sin(\phi_3); \ L_{32} &= \cos(\theta_3)\cos(\phi_3) \\ t_{41} &= \cos(\theta_4)\cos(\phi_4)\cos(\psi_4), \ t_{42} &= \cos(\theta_4)\cos(\phi_4)\sin(\psi_4) \\ t_{43} &= \cos(\theta_4)\sin(\phi_4), \ t_{44} &= \cos(\theta_4). \end{split}$$

etc.

- Then $\theta = (\theta_2, \theta_3, \phi_3, \theta_4, \phi_4, \psi_4).$
- More involved problem: periodicity introduces multiple solutions.

Numerical probability: VaR_{α} - $CVaR_{\alpha}$ II

• Set
$$H(\xi, x) = \partial_{\xi} v(\xi, x) = 1 - \frac{1}{1-\alpha} \mathbf{1}_{\{x \ge \xi\}} = \frac{1}{1-\alpha} (\mathbf{1}_{\{x \le \xi - \alpha\}})$$
 so that
 $V'(\xi) = \mathbb{E} H(\xi, X)$

• Set $\gamma_n = \frac{1}{n}$ and let X_n i.i.d., $\sim X$, then

$$\xi_{n+1} = \xi_n - \frac{\gamma_{n+1}}{1-\alpha} (\mathbf{1}_{\{X_{n+1} \le \xi_n\}} - \alpha) \longrightarrow \xi_\alpha = \operatorname{VaR}_\alpha(X).$$

• What about $\operatorname{CVaR}_{\alpha}(X)$? Various solutions...

$$\Xi_n = \frac{v(\xi_0, X_1) + \cdots + v(\xi_{n-1}, X_n)}{n} \longrightarrow \mathbb{E} v(\xi_\alpha, X) = \operatorname{CVaR}_{\alpha}(X).$$

• Recursive form
$$\Xi_n = \Xi_{n-1} - \frac{1}{n} \Big(\Xi_{n-1} - v(\xi_{n-1}, X_n) \Big), \ \Xi_0 = 0.$$

 Warning ! Rare events phenomenon tends to freeze the algorithm ⇒ adaptive Importance Sampling (⁶) !

• ... and try to slowly increase $\alpha = \alpha_n$ from $\alpha_0 = \frac{1}{2}$ to the target level.

Gilles PAGÈS (LPSM)

58 / 94

⁶O. Bardou, N. Frikha, G. Pagès (2009). Computing VaR and CVaR using Stochastic Approximation and Adaptive Unconstrained Importance Sampling, *Monte Carlo and Applications Journal*, 15(3):173–210.

- Low dimensional examples selected on purpose for expository.
- Not as automatic as (linear) Monte Carlo simulation: tuning of the step is mandatory.
- Many other examples : adaptive variance reduction (see Lemaire-P. 2007, AAP).
- Central-Limit Theorem, Averaging principle (Ruppert-Polyak).
- More details and results in $(^{7})$ if interested and the references therein.

⁷G. Pagès (2018). Introduction to Numerical Probability with application to Finance, Springer-Verlag, Berlin, 576p.

Table of Contents



Multilayer feedforward perceptron and Backpropagation

Learning

- Database (z_k)_{k=1:N}, parameters θ∈ Θ ⊂ ℝ^K and (local) loss function/predictor v(θ, z).
- Let $(I_k)_{k\geq 1}$ be an i.i.d. sequence $\mathcal{U}(\{1,\ldots,N\})$ -distributed.
- The stochastic gradient descent reads

$$\theta_{n+1} = \theta_n - \gamma_{n+1} \nabla_{\theta} v \big(\theta_n, z_{I_{n+1}} \big)$$

where $z_{I_{n+1}}$ means that a datum has been picked up at random in the database uniformly in $\{1, \ldots, N\}$.

Check that

$$\mathbb{E} \nabla_{\theta} v(\theta_n, z_l) = \frac{1}{N} \sum_{k=1}^{N} \nabla_{\theta} v(\theta_n, z_k)$$
$$= \int \nabla_{\theta} v(\theta_n, z) \mu_N(dz) = \nabla V(\theta_n)$$

CLVQ/k-means (unsupervised learning)

• Aim:

$$\min_{\theta^{j})_{j=1:r}} \left[V(\theta) = \frac{1}{2} \sum_{k=1}^{N} \min_{i=1:r} |\theta^{i} - x_{k}|^{2} \right]$$

(mean minimal distance to prototypes).

• Competitive Learning Vector Quantization:

$$\theta_{n+1}^{i} = \begin{cases} \theta_{n}^{i} - \gamma_{n+1} \left(\theta_{n}^{i} - x_{n+1} \right) & \text{if } |x_{n+1} - \theta_{n}^{i}| < \min_{j \neq i} |x_{n+1} - \theta_{n}^{j}| \\ = 0 & \text{otherwise} \end{cases}$$

• In other words: $\rightarrow n+1$ reads

- Nearest neighbour searchto the datum among *r* prototypes of dimension *d* .
- Moving the winner by a dilatation centered at the datum with ratio $1 \gamma_{n+1} > 0.$



Table of Contents



Multilayer feedforward perceptron and Backpropagation

Table of Contents



Multilayer feedforward perceptron and Backpropagation

Linear artificial neuron

- Mc Cullogh & Pitts in 1943 : linear "neuron": linear partitioner of two classes of date (if any...).
- Learning by Hebb's rule also known as reinforcement rule.
- The first perceptron: 1957 by Rosenblatt (⁸) as a binary classifier.
- Only able to classify linearly separable classes of data.



Figure: The Rosenblatt neural network performances.

^oRosenblatt, Frank (1957), The Perceptron: a perceiving and recognizing automaton. Report 85-460-1, Cornell Aeronautical Laboratory.

No hidden layer: ADALINE (Adaptive Linear Neuron (B. Widrow & T. Hoff, 1960))

- Input data x_1, \ldots, x_N of the form $x_k = (1, x_k^1, \ldots, x_N^d) \in \mathbb{R}^{d+1}$.
- Output (true) data are real numbers y_1, \ldots, y_N .
- Data set: $z_k = (x_k, y_k)$, k = 1 : N (supervised learning).
- Let $\Theta = \mathbb{R}^{d+1}$ the parameter set..
- The answer of the ADALINE neuron/network for $\theta \in \Theta$ "fed" with an imput datum $x = (1, x^1, \dots, x^d)$ is

$$\theta^{\top} x = \sum_{i=0}^{d} \theta^{i} x^{i} = \theta^{0} + \sum_{i=1}^{d} \theta^{i} x^{i}$$

where θ^{\top} denote the transpose of θ .

• (Convex) local prediction/loss function: $v(\theta, (x, y)) = \frac{1}{2}(y - \theta^{\top}x)^2$



Figure: The ADALINE neural network (with T. Montes)

•
$$V(\theta) = \frac{1}{2} \int (\theta^\top x - y)^2 \mu_N(d(x, y)) = \mathbb{E} (\theta^\top x_I - y_I)^2 \text{ (convex!)}$$

i.e. Global prediction/loss function = Ordinary Least squares !

•
$$\nabla V(\theta) = \int (\theta^{\top} x - y) x \mu_N(d(x, y))).$$

• The target is

$$\nabla V(\theta) = 0 \iff \int \underbrace{(\theta^{\top} x) x}_{=(xx^{\top})\theta} \mu_N(d(x, y)) = \int y x \mu_N(d(x, y))$$

$$\iff \theta_{opt}^{\top} = \left(\int x x^{\top} \mu_N(d(x, y))\right)^{-1} \left(\int y x \mu_N(d(x, y))\right)$$

• It is linear regression (as expected).

• Resulting (*SGD*):

$$\theta_{n+1} = \theta_n - \gamma_{n+1} \big(\theta_n^\top x_{l_{n+1}} - y_{l_{n+1}} \big) x_{l_{n+1}}$$

All assumptions are satisfied $\Longrightarrow \theta_n \to \theta^* a.s.$

- Conclusion : the ADALINE network performs linear regression without matrix inversion.
- Extension to q-dimensional (true) output data y_k with $\Theta = \mathbb{M}(d,q)$ and

$$v(\theta, (x, y)) = \frac{1}{2} |\theta^\top x - y|^2.$$

"Non-linear" ADALINE

- Assume the output data y_k are e.g. (0, 1)-valued.
- Let $\Psi \in C(\mathbb{R}, (0, 1))$, Ψ increasing homeomorphism from \mathbb{R} onto (0, 1).

•
$$V(\theta) = \frac{1}{2}\mathbb{E}\left(\Psi(\theta^{\top}x_I) - y_I\right)^2$$
.

•
$$\nabla V(\theta) = \int \left(\left(\Psi(\theta^{\top} x) - y \right) \Psi'(\theta^{\top} x) x \right) \mu_N(d(x, y)).$$

Loss of convexity since

$$\nabla^2 V(\theta) = \int \left[\Psi'(\theta^\top x)^2 + \left(\Psi(\theta^\top x) - y \right) \underbrace{\Psi''(\theta^\top x)}_{\geq 0} x x^\top \right] \mu_N(d(x, y))$$


Figure: The non-linear ADALINE neural network (with T. Montes).

• Historical Rosenblatt's "linear perceptron" (1957): $\Psi(u) = \mathbf{1}_{u \ge 0}$ is a linear classifier.



Figure: Historical Rosenblatt's perceptron: the hardware (without T. Montes !).

Table of Contents



Multilayer feedforward perceptron and Backpropagation

One hidden layer: universal approximation property



Figure: One hidden layer MLP (below: convention $\alpha_{ij} = w_{ij}$).

- $d = d_x$ -dimensional inputs. Switch to $x \rightsquigarrow \begin{pmatrix} 1 \\ x \end{pmatrix}$.
- L units i = 1 : L on the hidden layer, with an *activation function*

 $\Psi \in C(\mathbb{R},\mathbb{R}).$

- Unit *i* of the hidden layer receives $x = (x^j)_{j=1:d}$ and emits $\Psi((w_i \cdot | x)) = \Psi(w_{i0} + \sum_{1 \le i \le d} w_{ij} x^j).$
- The output layer receives $\left[\Psi((w_{i\cdot}|x))
 ight]_{i=1:L}$ and emits

$$\sum_{1\leq i\leq L}\lambda_{\ell}\Psi((w_{i\cdot}|x)).$$

Theorem (Cybenko, 1989)

(*) If the activation function
$$\Psi$$
 satisfies
 $(CL_{\Psi}) \equiv \Psi \in C_b(\mathbb{R}, \mathbb{R}), \quad \lim_{\xi \to -\infty} \Psi(\xi) = 0, \quad \lim_{\xi \to +\infty} \Psi(\xi) = 1,$
then $F = \left\{ x \longmapsto \sum_{\ell=1}^{L} \lambda_\ell \Psi((w_i \cdot | x)), \ L \in \mathbb{N}^*, \ \lambda \in \mathbb{R}^L, \ w \in \mathbb{R}^{L \times (d+1)} \right\}$

is $\|\cdot\|_{sup}$ -dense in $C([0,1]^d,\mathbb{R})$.

 $^a{\rm G.}$ Cybenko, Approximation by Superpositions of a Sigmoidal Function. Mathematics of Control, Signals, and Systems, 2:303-314, 1989.

• (CL_{Ψ}) can be slightly relaxed into Hornik's condition (⁹, 1989) $(CS_{\Psi}) \equiv \Psi \in C_b(\mathbb{R}, \mathbb{R}), \quad \exists x_1, x_2 \in \mathbb{R}^d \text{ such that } \Psi(x_1) \neq \Psi(x_2).$

• Extension to vector-valued outputs is straightforward.

⁹ Hornik K., Stinchcombe M. & White H., Multilayer feedforward networks are universal approximators. *Neural Networks*, **2**, 359–366

Proof I

Assume d = 1 (for simplicity).

► STEP 1. If $\overline{F}^{\|\cdot\|_{sup}} \subsetneq C([0,1],\mathbb{R})$, linear form $\Lambda \in C([0,1],\mathbb{R})^*$ s.t.

$$\Lambda \not\equiv 0$$
 and $\Lambda_{|\overline{F}} \equiv 0$ (by Hahn-Banach's Theorem).

By Riesz's representation Theorem, there exists a signed measure μ s.t.

$$\Lambda g = \int_{\mathbb{R}} g \, d\mu.$$

Hence, with for every weight $w = (\alpha, \beta)$,

$$\forall w = (\alpha, \beta) \in \mathbb{R}^2, \quad \int_{\mathbb{R}} \Psi(\alpha x + \beta) \mu(dx) = 0.$$

STEP 2. Let μ = μ⁺ − μ⁻, denote the/a decomposition of μ.
 Let β → +∞. By Lebesgue's Dominated Convergence (LDC)
 ∫ 1μ(dx) = 0 so that μ⁺(ℝ) = μ⁻(ℝ).

Proof II

• Now setting $\beta = -\alpha u$ yields

$$\forall \alpha, u \in \mathbb{R}, \quad \int \Psi(\alpha(x-u))\mu(dx) = 0.$$

Letting $\alpha \to +\infty$ implies, still by *LDC* theorem

 $\forall u \in \mathbb{R}, \quad \mu\bigl(\{x : x > u\}\bigr) + \Psi(0)\mu(\{u\}) = 0.$

• The set $D = \{u : \mu^+(\{u\}) > 0 \text{ or } \mu^-(\{u\}) > 0\}$ is at most countable and, for every $u \in D^c$,

 $\forall u \in D^c, \quad \mu(\{x : x > u\}) = 0.$

▶ STEP 3. Combined with $\mu^+(\mathbb{R}) = \mu^-(\mathbb{R})$, one has

$$\forall u \in \mathbb{R}, \quad \mu^+(\{x : x > u\}) = \mu^-(\{x : x > u\}).$$

Hence $\mu^+ = \mu^-$ since D^c is everywhere dense in \mathbb{R} i.e. $\mu \equiv 0$.

► STEP 4. Contradiction! Since $\Lambda \neq 0$.

Gilles PAGÈS (LPSM)

78 / 94

Higher dimensional outputs I

Instant result by concatenating networks



Although this is better in practice





Constructive proof/smoothness

- Both proofs are not constructive (Hornik's relies on a distribution based argument).
- Constructive approach are possible including rates for higher order derivatives like...

Theorem (Attali-Pagès, 1997)

(^a) Assume the activation function $\Psi \in C^r_b(\mathbb{R},\mathbb{R})$ and satisfies

 ψ is non-polynomial on any open interval and $\exists \xi_0 \ s.t. \ \psi^{(s)}(\xi_0) \neq 0, \ s = 0, \dots, r.$

Let $f \in C^{\infty}(\mathbb{R}^d, \mathbb{R})$. For every $L \in \mathbb{N}$, there exists weights $(w_{ij}^{(L)})_{i=1:d,j=1:L}$, $\lambda_{i=1:L}^{(L)}$ such that (with $\|g\|_{[0,1]^d} = \sup_{\xi \in [0,1]^d} |g(\xi)|$)

$$\max_{s=0:r} \left\| \partial^s \Big(\sum_{i=1}^L \lambda_i^{(L)} \psi\big((w_i^{(L)} | x) \big) \Big) - \partial^s f \right\|_{[0,1]^d} \le C_f \cdot L^{-1/d}$$

^a J.-G. Attali, G. Pagès, Approximations of Functions by a Multilayer a New Approach, *Neural Networks*, **10**(6):1069-10811,1997.

- The proof heavily relies on multi-variable Bernstein polynomials and Vandermonde determinants.
- Our 1997 paper does not avoid the curse of dimensionality and came too late: Vapnik's *SVM* were coming at the front owing to a nice mathematical framework.
- Y. Le Cun et al. stroke back in the early 2000's with outstanding performances classification results at a yearly challenge using deep learning...
- New networks with more layers, not fully connected and new type of units (convolutive units, recurrent units).
- and even more recently Generative Adversarial Networks (2014) seem to get rid of curse of dimensionality,
- but, so far, no theoretical evidence.

Table of Contents



Multilayer feedforward perceptron and Backpropagation

From learning to deep learning

- Feedforward multilayer perceptron.
- 1998 (Y. Le Cun): MNIST database of handwritten figures.
 - Input: Image = 28×28 pixels $\times 256$ grey levels $d_x \simeq 2 \times 10^5$.
 - Output : probability y_i for each figure 10 figures : $[0, 1]^{10}$ i.e. $d_y \simeq 9$.
 - Predictor : Multilayer perceptron (MLP) with 2 hidden layers with 200 units \implies $K = 200 \times 10 = 22 \times 10^3$ parameters
 - Size of the database: $N = 50\,000$ for learning, 10000 for testing.
- 2010: Same database
 - Predictor: MLP with convolutional units on GPU with 7 hidden layers

 $K = 2500 \times 2000 \times 1500 \times 1000 \times 500 \times 10 \simeq 122 \times 10^{6}$ parameters

- 2013 (Google) : ImageNet database of N = 162 × 10⁶ images of 100 × 100 pixels and 256 grey levels: d_x ~ 2.562 × 10⁶.
 - $d_y = 21$ (classifier across 21 classes).
 - Predictor: Convolutional MLP network $\implies K \simeq 1.72 \times 10^9$ parameters.
 - Error rate < 1/1000.



deep neural network

Figure: A three-hidden layer feedforward perceptron.

Table of Contents



• Multilayer feedforward perceptron and Backpropagation

Multilayer : toward Back propagation



Figure: A K - 1 hidden layer feedforward perceptron. Warning! New notation $w_{ii}^{(k)}$ instead of w_{ii} .

Multilayer and backpropagtion in one slide

- Input data (layer 0): $x^0 = x$ i.e. $x_j^0 = x_j$, $j = 1 : d_0 = d + 1$,
- K-1 hidden layers (k = 1 : K 1) and K + 1 layers (0 and K)
- Output of unit *i* of layer k 1: $x^{k-1,i}$.
- Input of unit *j* of layer *k*:

$$\sum_{1 \le i \le d_{k-1}} w_{ji}^{(k)} x^{k-1,i} = (w_{j.}^{(k)} | x^{k-1,\cdot})$$

so that, after passing through the activation function

$$x^{k,j} = \Psi_k ((w_{j}^{(k)} | x^{k-1,\cdot})).$$

• Typical activation functions

$$\Psi_k(\xi) = c \cdot rac{e^{\xi}}{e^{\xi}+1}, \quad \Psi_k(\xi) = c \cdot \arctan(\xi), \quad \Psi_k(\xi) = c \cdot \xi_+$$

In fact two slides

• Parameter to be calibrated to perform the learning of the network:

$$\mathbf{w} = (w^{(0)}, w^{(1)}, \dots, w^{(K)})$$

with $w^{(0)} = Id$, $w^{(k)} \in \mathbb{M}_{d_k, d_{k+1}}(\mathbb{R})$.

- Dimension of $\mathbf{w} = D = d_0 \cdot d_1 = \cdots + d_{K-1} \cdot d_K$.
- Local Predictor:

$$v(\mathbf{w},(x^0,y)) = \frac{1}{2}|y - x^K(\mathbf{w})|^2$$

• Global Predictor $V(w) = \mathbb{E} v(w, (x_I^0, y_I)), I \sim Unif(\{1, \dots, N\}).$

• Learning phase = Calibration

$\min_{w \in \mathbb{R}^D} V$

• (SGD) ? Needs to differentiate V in **w** i.e. $v(\mathbf{w}, (x^0, y))$!!

In fact three slides

Big Question: How to differentiate (and compute!) v(w, (x⁰, y)) ?
Easy part: from v(w, (x⁰, y)) = ¹/₂|y - x^K(w)|²

 $(^{\top}$ for transpose) \Downarrow (*J* for Jacobian matrix)

$$J_{\mathbf{w}}v(\mathbf{w},(x^0,y)) = (J_{\mathbf{w}}x^{K}(\mathbf{w}))^{\top}(x^{K}(\mathbf{w})-y),$$

• Use x^k as auxiliary variables and the induction

$$x^{k} = \varphi_{k}(w^{(k)}, x^{k-1}), \ k = 1 : K.$$

• Recursive formula: note that $J_{\mathbf{w}}x^k(\mathbf{w})^{\top} = J_{\mathbf{w}^{(1:k)}}x^k(\mathbf{w}^{(1:k)})^{\top}$ and

$$\begin{aligned} J_{\mathbf{w}} x^{k}(\mathbf{w})^{\top} &= J_{w^{(k)}} \varphi_{\kappa}(w^{(k)}, x^{k-1})^{\top} + J_{\mathbf{w}} x^{k-1}(\mathbf{w})^{\top} J_{x} \varphi_{k}(w^{(k)}, x^{k-1})^{\top} \\ &= J_{w^{(k)}} \varphi_{\kappa}(w^{(k)}, x^{k-1})^{\top} \\ &+ J_{\mathbf{w}^{(1:k-1)}} x^{k-1} (\mathbf{w}^{(1:k-1)})^{\top} J_{x} \varphi_{k}(w^{(k)}, x^{k-1})^{\top}. \end{aligned}$$

Finally ... in four slides

Lemma

In a ring $(A, +, \cdot)$ with regular multiplication, if

$$a_k = b_k + a_{k-1}c_k, \ k = 1:K$$

then: $a_{K} = b_{K} + b_{K-1}c_{K} + b_{K-2}c_{K-1}c_{K} + \dots + b_{1}c_{2}\cdots c_{k} + a_{0}c_{1}\cdots c_{k}$

which can be computed in a backward way.

The backpropagation algorithm $(^{10})$ is two-fold.

- A forward step: compute the values x_i^k at each unit j of each layer k.
- A backward step: Apply the lemma to the recursion
 - Compute

$$J_{w^{(k)}}\varphi_{k}(w^{(k)}, x^{k-1})^{\top} \quad \text{and} \quad J_{\mathbf{w}^{(1:k-1)}}x^{k-1}(\mathbf{w}^{(1:k-1)})^{\top}$$

And that's it!!

¹⁰Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. Nature, 323, 533–536. But, goes back to Paul Werbos en the 1970's for AD.

The truth...about Backpropagation of gradient

- Backpropagation of gradient is similar to the reverse mode of automatic differentiation developed independently!!
- True formulas



 $= \frac{1}{2} \left[\mu_{\nu}(h_{N}) \sum_{i}^{h_{N}} g' \left(\overline{\pi}_{h}((u_{i}^{(N)} | x^{h_{1}})) - \eta_{i} \right) \overline{\pi}_{h}((u_{i}^{(N)} | x^{h_{1}})) \right] \overline{\pi}(u_{i}^{(N)} | x^{h_{1}})$ $\frac{1}{6\pi} \frac{\mathbf{\hat{g}}(\mathbf{w}_{R_{1}}^{(\mathbf{K})} | \mathbf{x}^{\mathbf{K}_{R_{1}}})}{\sum w_{R_{1}}^{(\mathbf{K})} = \sum \frac{1}{2\pi} \frac$ $\sum_{i=1}^{\infty} \left[p_{\mathbf{x}}(l)(\mathbf{x}_{i}) \right] \sum_{i=1}^{\infty} \frac{\gamma_{\mathbf{x}}^{\mathbf{x}_{i}} \mathbf{y}^{\mathbf{x}_{i}}}{\gamma_{\mathbf{x}}^{\mathbf{x}_{i}}} \sum_{t=1}^{t} \left(\overline{\mathbf{x}}_{\mathbf{x}_{i}}(\mathbf{b}_{\mathbf{x}}^{\mathbf{x}_{i}} \mathbf{y}) - \mathbf{y}_{i} \right) \overline{\mathbf{x}}_{i}^{T} \left(\mathbf{x}_{i}^{\mathbf{x}_{i}} \mathbf{y}_{i}^{\mathbf{x}_{i}} \mathbf{y} \right) - \mathbf{y}_{i}^{T} \right) \overline{\mathbf{x}}_{i}^{T} \left(\mathbf{x}_{i}^{\mathbf{x}_{i}} \mathbf{y}_{i}^{\mathbf{x}_{i}} \mathbf{y} \right) - \mathbf{y}_{i}^{T} \right) \overline{\mathbf{x}}_{i}^{T} \left(\mathbf{x}_{i}^{\mathbf{x}_{i}} \mathbf{y}_{i}^{\mathbf{x}_{i}} \mathbf{x}_{$ =] pulling) I show I eke wight = [po(dlog)] I Jakim (ek, hall) by . $\frac{\partial x^{k_1 q_1}}{\partial u^{(k_1)}} = \frac{\partial}{\partial u^{(k_1)}} \widehat{\mathcal{L}}_{K_1} \left(\left(u^{(k_2)}_{q_1} \right| x^{k_2 q_1} \right) \right)$ $\in \overline{\mathfrak{D}}_{K_1}^{\frac{1}{2}}\left(\left|\mathbf{v}_{\mathbf{k}_1}^{\mathbf{k}_2}\right|\mathbf{z}^{\mathbf{k}_2}\right)\sum_{i=1}^{m_1}w_{\mathbf{k}_1\mathbf{k}_1}^{\mathbf{k}_2} \cdot \frac{\mathbf{j}\mathbf{x}^{\mathbf{k}_2,\mathbf{k}_1}}{\mathbf{k}_1\mathbf{k}_1}$ June coule & intransidiais on a la formels de Tell - Hill, chin = Ir ((un att)) x many but it breach dige

Lunger Von alter Vile couche le elle in, I want to Force le : Wie down $\frac{\partial \mathcal{E}}{\partial w^{M_{1}}} = \int d\mu_{W}(v_{0}) \frac{\partial \mathcal{E}}{\partial v} \frac{\partial \mathcal{R}^{K_{1}} \rho}{\partial w^{M_{2}}} \times \left(u_{1,0}^{(K_{0})} \mid e^{k u_{1,0}} \rho \right)_{R^{-1}} dv_{1}$ $\underset{\mathbf{x}}{=} \widetilde{\mathbb{E}}_{\mathbf{k}^{*}}^{\prime} \left((\mathbf{x}^{\mathbf{k}_{i}} | \mathbf{w}_{\mathbf{k}}^{(\mathbf{k})}) \right) \times$ $\frac{\Im \mathcal{E}^{(kl)}}{\Im u^{(l)}} \int \mu_{\nu}(dhy) \Phi_{k}\left(\left(z^{k,i}|u^{(l)}_{i}\right)) z^{k,i}(u^{(l)}_{i})|e^{ku_{i}}$ la etre quincel le récomme révergente (k) à juilie de la formles, a mot an course une desente 1) de gerdind (GD) will an will - You BE a particula (4) (dialor de (4) $= u \frac{\partial^2}{\partial t} u^{\mu} = u \frac{\partial^2}{\partial t} - y_{\mu\nu} \frac{\partial E}{\partial t} (u^{\mu} \partial_{t} u) \rightarrow (and u) d$ 2) Le grelient et aletique (SCB). $w_{j_1}^{(k)} = w_{j_2}^{(k),n} - \chi_{w_1} \overline{\Psi}_{k}^{i} \left(\left(w_{j_1}^{(k),n} \mid x_{u_{n_1}}^{(k),i} \right) \right) x_{u_{n_1}}^{(k),i}$ & provide la propriété d'apprennation universal : Decembra à Londrendie The sime (thank) Sol I: R - K use foretion continue, proved at man andarte (3 motors, ty za) + 200). Also $x \mapsto \overset{\mathbb{R}}{\longrightarrow} \overset{\mathbb{R}}{\longrightarrow} \overset{\mathbb{R}}{=} \left\{ (w|x)_{\pm} w_{1} \right\}, \ w_{7} \perp, \ \lambda \in \mathbb{R} \ w \in \mathbb{R}^{d} \ \} \ \text{sh down}$ dans E(las) , le venttet vide une si l'an vougles Carly you we compart guillangue the Red Revelled to a K. Hamk in 1991 deer Wand all and " Approximation availities of rultilys Fallmed Network. Tel qu'ésnesi, ce théorine demokre le propieté d'approximation surissant





Grazie per l'attenzione ! Merci de votre attention ! Thank you for your attention !