GRAN SASSO Science Institute Intensive Trimester "Particles, Fluids and Patterns: Analytical and Computational Challenges" 3-4 June 2022

High-dimensional approximation and sampling Part 1: Elements of approximation theory

Anthony Nouy

Centrale Nantes, Nantes Université, Laboratoire de Mathématiques Jean Leray

Many problems of computational science, statistics and probability require the approximation, integration or optimization of functions of many variables

 $u(x_1,\ldots,x_d)$

- High dimensional PDEs (Boltzmann, Schrödinger, Black-Scholes...)
- Parameter-dependent or stochastic equations
- Multiscale problems
- Statistical learning (density estimation, classification, regression)
- Probabilistic modelling

• ...

The goal of approximation is to replace a target function u by a simpler function (easy to evaluate and to operate with).

An approximation is searched in a set of functions X_n , where *n* is related to some complexity measure, typically the number of parameters.

Approximation

We distinguish

• linear approximation when X_n is a finite-dimensional linear space (polynomials, trigonometric polynomials, fixed knot splines...)

$$X_n = \{\sum_{i=1}^n a_i \varphi_i : a_i \in \mathbb{R}\}$$

where the φ_i form a basis of X_n .

Approximation

We distinguish

• linear approximation when X_n is a finite-dimensional linear space (polynomials, trigonometric polynomials, fixed knot splines...)

$$X_n = \{\sum_{i=1}^n a_i \varphi_i : a_i \in \mathbb{R}\}$$

where the φ_i form a basis of X_n .

• nonlinear approximation when X_n is a nonlinear set (rational functions, free knot splines, *n*-term approximation, neural networks, tensor networks...), e.g.

$$X_n = \{\sum_{i=1}^n a_i \varphi_i : a_i \in \mathbb{R}, \varphi_i \in \mathcal{D}\}$$

for *n*-term approximation from a dictionary of functions \mathcal{D} , or

$$X_n = \{g(a) : a \in \mathbb{R}^n\}$$

with some given nonlinear map g from \mathbb{R}^n to X_n .

Error of best approximation

For a given function u from a normed vector space X and a given subset X_n , the error of best approximation

$$e_n(u)_X := E(u, X_n)_X = \inf_{v \in X_n} ||u - v||_X$$

quantifies the best we can expect from X_n .



For a sequence $(X_n)_{n\geq 1}$ of sets of growing complexity, called an approximation tool, we would like to address the following questions.

• (universality) Does $e_n(u)_X$ converge to 0 for all functions u in X ?

For a sequence $(X_n)_{n\geq 1}$ of sets of growing complexity, called an approximation tool, we would like to address the following questions.

- (universality) Does $e_n(u)_X$ converge to 0 for all functions u in X ?
- (expressivity) For a certain class of functions in X, determine how fast $e_n(u)_X$ converges to 0, or determine the complexity $n = n(\epsilon, u)$ such that $e_n(u) \le \epsilon$. Typically,

$$e_n(u)_X \leq M\gamma(n)^{-1}$$

where γ is a strictly increasing function (growth function), and

$$n(\epsilon, u) \geq \gamma^{-1}(\epsilon/M)$$

For a sequence $(X_n)_{n\geq 1}$ of sets of growing complexity, called an approximation tool, we would like to address the following questions.

- (universality) Does $e_n(u)_X$ converge to 0 for all functions u in X ?
- (expressivity) For a certain class of functions in X, determine how fast $e_n(u)_X$ converges to 0, or determine the complexity $n = n(\epsilon, u)$ such that $e_n(u) \le \epsilon$. Typically,

$$e_n(u)_X \leq M\gamma(n)^{-1}$$

where γ is a strictly increasing function (growth function), and

$$n(\epsilon, u) \geq \gamma^{-1}(\epsilon/M)$$

• (approximation classes) Characterize the class of functions for which a certain convergence type is achieved, e.g.

$$\mathcal{A}^{\gamma}(X,(X_n)_{n\geq 1}) = \left\{ u: \sup_{n\geq 1} \gamma(n)e_n(u)_X < +\infty \right\}$$

for some growth function γ .

• (proximinality) Determine if for all $u \in X$, there exists an element of best approximation $u_n \in X_n$ such that

$$\|u-u_n\|_X=e_n(u)_X.$$

• (proximinality) Determine if for all $u \in X$, there exists an element of best approximation $u_n \in X_n$ such that

$$\|u-u_n\|_X=e_n(u)_X.$$

• (algorithm) Construct an approximation $u_n \in X_n$ such that

$$\|u-u_n\|_X \leq Ce_n(u)_X$$

with C independent of n or $C(n)e_n(u) \rightarrow 0$ as $n \rightarrow \infty$.

Algorithms depend on the available information, e.g. given by linear functionals such as evaluations of the function (interpolation, discrete least-squares), or equations satisfied by the function (variational/Galerkin methods).

The approximation u_n should be constructed with limited amount of information / limited complexity.

• (optimal approximation/algorithm) If we know that the function u belongs to some class of functions K, we would like to find an approximation tool X_n together with an algorithm presenting a good performance, or even the optimal performance for that class.

A fundamental problem is to quantify the best we can expect.

1. About optimal performance we can expect for the approximation of a function class *K*. That will reveal the need to exploit structures for high-dimensional approximation, beyond classical regularity.

- 1. About optimal performance we can expect for the approximation of a function class *K*. That will reveal the need to exploit structures for high-dimensional approximation, beyond classical regularity.
- 2. Classical tools X_n for high-dimensional approximation, sufficiently flexible to approximate a large class of functions. Focus on neural networks and tensor networks.

- 1. About optimal performance we can expect for the approximation of a function class *K*. That will reveal the need to exploit structures for high-dimensional approximation, beyond classical regularity.
- 2. Classical tools X_n for high-dimensional approximation, sufficiently flexible to approximate a large class of functions. Focus on neural networks and tensor networks.
- 3. How to construct dedicated approximation tools X_n for specific function classes K, using manifold approximation methods.

- 1. About optimal performance we can expect for the approximation of a function class *K*. That will reveal the need to exploit structures for high-dimensional approximation, beyond classical regularity.
- 2. Classical tools X_n for high-dimensional approximation, sufficiently flexible to approximate a large class of functions. Focus on neural networks and tensor networks.
- 3. How to construct dedicated approximation tools X_n for specific function classes K, using manifold approximation methods.
- 4. How to construct approximations from limited information, with a focus on how to generate a good information, for linear and nonlinear approximation.

Assume we want to approximate functions for some set K in X.

To measure the optimal performance we can expect from an approximation tool and associated algorithm, we rely on different measures of complexity of K depending on

• how we measure the error over K:

$$\sup_{u \in K} E(u, X_n)_X \quad \text{(worst-case error)},$$
$$\left(\int_K E(u, X_n)_X^p d\mu(u)\right)^{1/p} \quad \text{(average error)}$$

Assume we want to approximate functions for some set K in X.

To measure the optimal performance we can expect from an approximation tool and associated algorithm, we rely on different measures of complexity of K depending on

• how we measure the error over K:

$$\sup_{u \in K} E(u, X_n)_X \quad \text{(worst-case error)},$$
$$\left(\int_K E(u, X_n)_X^p d\mu(u)\right)^{1/p} \quad \text{(average error)}$$

• the type of approximation (linear or nonlinear),

Assume we want to approximate functions for some set K in X.

To measure the optimal performance we can expect from an approximation tool and associated algorithm, we rely on different measures of complexity of K depending on

• how we measure the error over K:

$$\sup_{u \in K} E(u, X_n)_X \quad \text{(worst-case error)},$$
$$\left(\int_K E(u, X_n)_X^p d\mu(u)\right)^{1/p} \quad \text{(average error)}$$

- the type of approximation (linear or nonlinear),
- and possibly the properties of the approximation process (type of information, stability...).

Outline

Optimal linear approximation

- 2 Optimal nonlinear approximation
- 3 How to beat the curse of dimensionality ?

Optimal linear approximation (worst-case setting) Kolmogorov widths

For a compact subset K of a normed vector space X and a *n*-dimensional space X_n in X, we define the worst-case error

$$E(K, X_n)_X = \sup_{u \in K} E(u, X_n)_X = \sup_{u \in K} \inf_{v \in X_n} ||u - v||_X$$



Optimal linear approximation (worst-case setting) Kolmogorov widths

The Kolmogorov n-width of K is defined as

$$d_n(K)_X = \inf_{\dim(X_n)=n} \sup_{u \in K} E(u, X_n)_X = \inf_{\dim(X_n)=n} \sup_{u \in K} \inf_{v \in X_n} ||u - v||_X$$

where the infimum is taken over all linear subspaces X_n of dimension n.



 $d_n(K)_X$ measures how well the set K can be approximated (uniformly) by a *n*-dimensional space. It measures the ideal performance that we can expect from linear approximation methods.

Near to optimal spaces can be constructed by greedy algorithms (see in part 3).

Optimal linear approximation (average setting) Average Kolmogorov widths

If K is equipped with a measure μ , an average Kolmogorov *n*-width is defined by

$$d_n^{(p)}(K,\mu)_X = \inf_{\dim(X_n)=n} \left(\int_K E(u,X_n)_X^p d\mu(u) \right)^{1/p}.$$

If the measure is finite,

$$d_n^{(p)}(K,\mu)_X \leq \mu(K)^{1/p} d_n(K)_X.$$

For X a Hilbert space, p = 2 and μ the push-forward measure of a K-valued random variable $U \in L^2(\Omega; X)$, this is equivalent to

$$\inf_{\dim(X_n)=n} \mathbb{E}(\|U - P_{X_n}U\|_X^2)^{1/2}$$

and an optimal space is given by Principal Component Analysis, that is a dominant eigenspace of the operator $v \mapsto \mathbb{E}((U, v)_X U)$ (see in part 3).

Optimal linear approximation (worst-case setting) Linear widths

Another measure of complexity taking into account the approximation process is the linear width

$$a_n(K)_X = \inf_{A_n} \sup_{u \in K} \|u - A_n u\|_X$$

where the infimum is taken over all continuous linear maps $A_n : K \to X$ with rank at most n.

Equivalently,

$$a_n(K)_X = \inf_{g,a} \sup_{u \in K} \|u - g(a(u))\|_X$$

where both $a: K \to \mathbb{R}^n$ (encoder) and $g: \mathbb{R}^n \to X$ (decoder) are linear maps.

Optimal linear approximation (worst-case setting) Linear widths

Another measure of complexity taking into account the approximation process is the linear width

$$a_n(K)_X = \inf_{A_n} \sup_{u \in K} \|u - A_n u\|_X$$

where the infimum is taken over all continuous linear maps $A_n : K \to X$ with rank at most n.

Equivalently,

$$a_n(K)_X = \inf_{g,a} \sup_{u \in K} \|u - g(a(u))\|_X$$

where both $a : K \to \mathbb{R}^n$ (encoder) and $g : \mathbb{R}^n \to X$ (decoder) are linear maps. For a Hilbert space X,

$$a_n(K)_X = d_n(K)_X = \inf_{\dim(X_n)=n} \sup_{u \in K} \|u - P_{X_n}u\|_X$$

with optimal map A_n given by the orthogonal projection P_{X_n} onto an optimal space X_n . For a general Banach space X,

$$d_n(K)_X \leq a_n(K)_X \leq \sqrt{n} d_n(K)_X$$

Optimal linear approximation from point evaluations (worst case) Linear sampling numbers

For functions defined on a set X, by restricting the information to point evaluations, the performance is characterized by sampling numbers.

For deterministic information, the worst-case optimal performance for the approximation of functions in K is measured through the linear sampling number

$$\rho_n(K)_X = \inf_{x} \inf_{R} \sup_{u \in K} ||u - R(u(x_1), \dots, u(x_n))||_X$$

where the infimum is taken over all linear maps R and points $\mathbf{x} = (x_1, \ldots, x_n) \in \mathcal{X}^n$, or equivalently

$$\rho_n(K)_X = \inf_{\mathbf{x}} \inf_{\varphi_1, \dots, \varphi_n \in X} \sup_{u \in K} \|u - \sum_{i=1}^n u(x_i)\varphi_i\|_X$$

This quantifies the best we can expect from a linear algorithm using n samples for the approximation of functions in the class K.

Clearly,

$$\rho_n(K)_X \geq a_n(K)_X \geq d_n(K)_X$$

For random information, the optimal performance can be measured in average mean squared error through the (linear) sampling number

$$\rho_n^{rand}(K)_X^2 = \inf_{\nu^n} \inf_R \sup_{u \in K} \mathbb{E}_{\mathbf{x} \sim \nu^n} (\|u - R(u(x_1), \dots, u(x_n))\|_X^2)$$

with an infimum taken over all measures ν^n on \mathcal{X}^n and linear maps R. Choosing for ν^n a Dirac measure on an optimal deterministic set of points, we deduce that

$$d_n(K)_X \leq \rho_n^{rand}(K)_X \leq \rho_n(K)_X$$

For random information, the optimal performance can be measured in average mean squared error through the (linear) sampling number

$$\rho_n^{rand}(K)_X^2 = \inf_{\nu^n} \inf_R \sup_{u \in K} \mathbb{E}_{\mathbf{x} \sim \nu^n} (\|u - R(u(x_1), \dots, u(x_n))\|_X^2)$$

with an infimum taken over all measures ν^n on \mathcal{X}^n and linear maps R. Choosing for ν^n a Dirac measure on an optimal deterministic set of points, we deduce that

$$d_n(K)_X \leq \rho_n^{rand}(K)_X \leq \rho_n(K)_X$$

The question is how far sampling numbers $\rho_n(K)_X$ or $\rho_n^{rand}(K)_X$ are from Kolmogorov widths $d_n(K)_X$, and how to generate optimal samples and algorithms in practice.

A series of results have been recently obtained for L^2 approximation, comparing sampling numbers with Kolmogorov widths, e.g. [Cohen and Dolbeault 2021, Nagel, Schafer and Ullrich 2021, Temlyakov 2021, Dolbeault, Krieg and Ullrich 2022].

These results are based on constructive approaches for the approximation of functions in a given model class.

See in the last part.

Other complexity measures can be defined with general linear information

 $\ell_1(u),\ldots,\ell_n(u)$

Information can be chosen non-adaptively or adaptively, deterministically or randomly (see [Krieg et al 2024] for comparisons).

Bounds of Kolmogorov widths $d_n(K)_X$

Upper bounds for $d_n(K)_X$ can be obtained by specific linear approximation methods. Proofs are sometimes constructive.

Lower bounds for $d_n(K)$ can be obtained using different techniques.

• Using diversity in K:

$$d_n(K)_X \ge d_n(S)_X$$

with S some subset of K whose Kolmogorov width can be bounded from below.

Bounds of Kolmogorov widths $d_n(K)_X$

Upper bounds for $d_n(K)_X$ can be obtained by specific linear approximation methods. Proofs are sometimes constructive.

Lower bounds for $d_n(K)$ can be obtained using different techniques.

• Using diversity in K:

$$d_n(K)_X \ge d_n(S)_X$$

with S some subset of K whose Kolmogorov width can be bounded from below.

Example: if X is a Hilbert space and K contains a set of orthogonal vectors $S = \{u_1, \ldots, u_m\}$ with norm $||u_i||_X = c_m$,

$$d_n(K)_X \geq d_n(S)_X = d_n(bc(S))_X = d_n(c_m B(\ell_1(\mathbb{R}^m)))_{\ell_2} \geq c_m \sqrt{1 - n/m}$$

where we used the fact that the n-width of S is equal to the n-width of its balanced convex hull

$$bc(S) = \{\sum_{i=1}^{m} \alpha_i u_i : u_1, \dots, u_m \in S, \sum_{i=1}^{m} |\alpha_i| \leq 1, m \in \mathbb{N}\}$$

which is isomorphic to $c_m B(\ell_1(\mathbb{R}^m))$, and a result of Stechkin (1954) for *n*-widths of ℓ^p balls.

• Using Bernstein width

$$b_n(K)_X = \sup_{\dim(X_{n+1})=n+1} \sup\{r : rB(X_{n+1}) \subset K\}$$

that is the largest r > 0 such that K contains the ball of radius r of some (n + 1)-dimensional space

 $d_n(K)_X \geq b_n(K)_X$



Bounds of Kolmogorov widths $d_n(K)_X$

• Using covering number $N_{\epsilon}(K)_X$ (minimal number of balls of radius ϵ for covering K) or entropy numbers

$$\epsilon_n(\mathcal{K})_X = \inf\{\epsilon : \mathcal{K} \subset \bigcup_{i=1}^{2^n} B(u_i, \epsilon), u_i \in \mathcal{K}\} = \inf\{\epsilon : \log_2(N_\epsilon(\mathcal{K})_X) \le n\}$$

that is the smallest ϵ such that K can be covered by 2^n balls of radius ϵ . Any $u \in K$ can be encoded with n bits up to precision $\epsilon_n(K)$.



Bounds of Kolmogorov widths $d_n(K)_X$

• Using covering number $N_{\epsilon}(K)_X$ (minimal number of balls of radius ϵ for covering K) or entropy numbers

$$\epsilon_n(K)_X = \inf\{\epsilon : K \subset \bigcup_{i=1}^{2^n} B(u_i, \epsilon), u_i \in K\} = \inf\{\epsilon : \log_2(N_\epsilon(K)_X) \le n\}$$

that is the smallest ϵ such that K can be covered by 2^n balls of radius ϵ . Any $u \in K$ can be encoded with n bits up to precision $\epsilon_n(K)$.



Carl's inequality: for all s > 0,

$$(n+1)^{s}\epsilon_{n}(\mathcal{K})_{X} \leq C_{s} \sup_{0\leq m\leq n} (m+1)^{s}d_{m}(\mathcal{K})_{X}$$

Therefore, if $\epsilon_n(K)_X \gtrsim n^{-s}$, then $d_n(K)_X \lesssim n^{-r}$ can not hold with r > s.

Kolmogorov width of a set of discontinuous functions

Letting $\chi_{(a,b)}$ denote the indicator function of the interval (a, b), consider for K the set of indicator functions in $X = L^2(-1, 1)$:

$$K = \{\chi_{[-1,s]} : s \in [-1,1]\}.$$

The balanced convex hull bc(K) contains the set of functions

$$S = \{\psi_i = \frac{1}{2}\chi_{(\mathsf{x}_i,\mathsf{x}_{i+1}]} : 1 \le i \le m\}$$

with $x_i = -1 + 2i/m$ and $\|\psi_i\|_{L^2} = (2m)^{-1/2} := c_m$.

Therefore it holds

$$d_n(K)_X = d_n(bc(K))_X \ge d_n(S)_X = d_n(bc(S))_X = d_n(c_m B(\ell_1(\mathbb{R}^m)))_{\ell^2} \ge c_m \sqrt{1 - n/m}$$

Taking m = 2n, we obtain

$$d_n(K)_X \geq 2^{-3/2} n^{-1/2}$$
For $X = L^{p}(\mathcal{X})$, $\mathcal{X} = [0, 1]^{d}$, $1 \leq p \leq \infty$, and K the unit ball of $W^{k,p}(\mathcal{X})$, that are functions u such that

$$\max_{|\alpha|_{\mathbf{1}}\leq k}\|D^{\alpha}u\|_{L^{p}}\leq 1,$$

it holds

$$d_n(K)_X \sim n^{-k/d}$$

and optimal performance is obtained e.g. by fixed knot splines (with degree adapted to the regularity).

For $X = L^{p}(\mathcal{X})$, $\mathcal{X} = [0, 1]^{d}$, $1 \leq p \leq \infty$, and K the unit ball of $W^{k,p}(\mathcal{X})$, that are functions u such that

$$\max_{|\alpha|_{\mathbf{1}}\leq k}\|D^{\alpha}u\|_{L^{p}}\leq 1,$$

it holds

$$d_n(K)_X \sim n^{-k/d}$$

and optimal performance is obtained e.g. by fixed knot splines (with degree adapted to the regularity).

We observe

- the curse of dimensionality: deterioration of the rate of approximation when *d* increases. Exponential growth with *d* of the complexity for reaching a given accuracy.
- the blessing of smoothness: improvement of the rate of approximation when k increases.

For $X = L^{p}(\mathcal{X})$, $\mathcal{X} = [0, 1]^{d}$, $1 \le p \le \infty$, and K the unit ball of $MW^{k,p}(\mathcal{X})$ (Sobolev space with dominating mixed smoothness), that are functions u such that

 $\max_{|\alpha|_{\infty}\leq k}\|D^{\alpha}u\|_{L^{p}}\leq 1.$

we have

$$d_n(K)_X \sim n^{-k} \log(n)^{k(d-1)}.$$

with optimal performance achieved by hyperbolic cross approximation (*n*-term approximation with tensor products of dilated splines) [Dung et al 2016].

For $X = L^{p}(\mathcal{X})$, $\mathcal{X} = [0, 1]^{d}$, $1 \le p \le \infty$, and K the unit ball of $MW^{k,p}(\mathcal{X})$ (Sobolev space with dominating mixed smoothness), that are functions u such that

 $\max_{|\alpha|_{\infty}\leq k}\|D^{\alpha}u\|_{L^{p}}\leq 1.$

we have

$$d_n(K)_X \sim n^{-k} \log(n)^{k(d-1)}.$$

with optimal performance achieved by hyperbolic cross approximation (*n*-term approximation with tensor products of dilated splines) [Dung et al 2016].

Curse of dimensionality is milder but still present.

Outline

- 1 Optimal linear approximation
- Optimal nonlinear approximation
- 3 How to beat the curse of dimensionality ?

For evaluating the ideal performance of nonlinear methods for the approximation of functions from a class K, different notions of widths have been introduced.

Optimal nonlinear approximation (worst case setting) Nonlinear Kolmogorov widths

A measure of complexity closely related to *n*-term approximation and relevant for nonlinear model reduction is the nonlinear Kolmogorov width [Temlyakov 1998] or library width

$$d_n(K,N)_X = \inf_{\#\mathcal{L}_{N,n}=N} \sup_{u \in K} \inf_{V_n \in \mathcal{L}_n} E(u,V_n)_X$$

where the infimum is taken over all libraries $\mathcal{L}_{N,n}$ of N linear spaces of dimension n.



For a given library $\mathcal{L}_{N,n} = \{V_1, \dots, V_N\}$, this corresponds to an approximation in $M = \bigcup_{k=1}^N V_k$.

Optimal nonlinear approximation (worst case setting) Nonlinear Kolmogorov widths

Choosing N = N(n), this yields a width only depending on n. Interesting regimes are $N(n) = b^n$ or $N(n) = n^{\alpha n}$.

It clearly holds

$$d_1(K,2^n)_X \leq \epsilon_n(K)_X$$

Also, we have a Carl's type inequality: for all r > 0,

$$n^r \epsilon_n(K)_X \leq C(r,b) \max_{1 \leq k \leq n} k^r d_{k-1}(K,b^k)_X.$$

Therefore if for some b > 0, $d_{n-1}(K, b^n)_X \leq n^{-r}$, then $\epsilon_n(K)_X \leq n^{-r}$.

For unit balls K of Besov spaces $B_q^{\alpha}(L^{\tau})$ compactly embedding in $L^p((0,1)^d)$, since $\epsilon_n(K) \gtrsim n^{-\alpha/d}$, we deduce that $d_n(K, b^n)_X \lesssim n^{-\beta}$ can not hold with $\beta > \alpha/d$.

Consider the approximation from a *n*-dimensional "manifold"

 $X_n = \{g(a) : a \in \mathbb{R}^n\}$

parametrized by a nonlinear map $g : \mathbb{R}^n \to X$. We could consider the problem of finding the best manifold of dimension *n* for approximating functions from *K*:

 $\inf_{g} \sup_{u \in K} \inf_{a \in \mathbb{R}^n} \|u - g(a)\|_X := \eta_n$

where the infimum is taken among all maps g from \mathbb{R}^n to X.

Consider the approximation from a *n*-dimensional "manifold"

 $X_n = \{g(a) : a \in \mathbb{R}^n\}$

parametrized by a nonlinear map $g : \mathbb{R}^n \to X$. We could consider the problem of finding the best manifold of dimension *n* for approximating functions from *K*:

 $\inf_{g} \sup_{u \in K} \inf_{a \in \mathbb{R}^n} \|u - g(a)\|_X := \eta_n$

where the infimum is taken among all maps g from \mathbb{R}^n to X.

For any compact set K, $\eta_n = 0$ for all $n \ge 1$. Indeed, K admits a countable dense subset $\{u_i\}_{i\in\mathbb{N}}$ (space-filling manifold). For n = 1, letting $g(a) = u_k$ for $a \in [k, k + 1)$, we obtain $\eta_1 = 0$.

Consider the approximation from a *n*-dimensional "manifold"

 $X_n = \{g(a) : a \in \mathbb{R}^n\}$

parametrized by a nonlinear map $g : \mathbb{R}^n \to X$. We could consider the problem of finding the best manifold of dimension *n* for approximating functions from *K*:

 $\inf_{g} \sup_{u \in K} \inf_{a \in \mathbb{R}^n} \|u - g(a)\|_X := \eta_n$

where the infimum is taken among all maps g from \mathbb{R}^n to X.

For any compact set K, $\eta_n = 0$ for all $n \ge 1$. Indeed, K admits a countable dense subset $\{u_i\}_{i\in\mathbb{N}}$ (space-filling manifold). For n = 1, letting $g(a) = u_k$ for $a \in [k, k + 1)$, we obtain $\eta_1 = 0$.

We can even provide a continuous parametrization, by considering a dense subset $\{u_i\}_{i \in \mathbb{Z}}$ and $g(a) = (a - k)u_{k+1} + (k + 1 - a)u_k$ for $a \in [k, k + 1]$.

Consider the approximation from a *n*-dimensional "manifold"

 $X_n = \{g(a) : a \in \mathbb{R}^n\}$

parametrized by a nonlinear map $g : \mathbb{R}^n \to X$. We could consider the problem of finding the best manifold of dimension *n* for approximating functions from *K*:

 $\inf_{g} \sup_{u \in K} \inf_{a \in \mathbb{R}^n} \|u - g(a)\|_X := \eta_n$

where the infimum is taken among all maps g from \mathbb{R}^n to X.

For any compact set K, $\eta_n = 0$ for all $n \ge 1$. Indeed, K admits a countable dense subset $\{u_i\}_{i\in\mathbb{N}}$ (space-filling manifold). For n = 1, letting $g(a) = u_k$ for $a \in [k, k + 1)$, we obtain $\eta_1 = 0$.

We can even provide a continuous parametrization, by considering a dense subset $\{u_i\}_{i \in \mathbb{Z}}$ and $g(a) = (a - k)u_{k+1} + (k + 1 - a)u_k$ for $a \in [k, k + 1]$.

In general, the map which associates to $u \in K$ the coefficients a(u) of its best approximation (if it exists) is not continuous, which makes the approximation process not reasonable.

Optimal nonlinear approximation (worst case setting) Manifold widths

The following definition of manifold width [DeVore, Howard and Micchelli 1989] quantifies how well the set K can be approximated by *n*-dimensional nonlinear manifolds having continuous parametrization and a continuous parameter selection

$$\delta_n(K)_X = \inf_{\substack{g,a\\ u \in K}} \sup \|u - g(a(u))\|_X$$

where the infimum is taken over all continuous maps a (encoder) from K to \mathbb{R}^n and all continuous maps g (decoder) from \mathbb{R}^n to K.



Optimal nonlinear approximation (worst case setting) Manifold widths

Further assuming that encoders a and decoders g are Lipschitz continuous yields the notion of stable width introduced in [Cohen et al 2022].

As for linear widths, manifold widths are lower bounded by Bernstein widths [DeVore, Howard and Micchelli 1989, Theorem 3.1]

 $\delta_n(K)_X \geq b_n(K)_X.$

Optimal nonlinear approximation (worst-case setting) Sensing numbers and Gelfand widths

Sensing numbers measure the optimal performance of nonlinear approximation using linear information

$$s_n(K)_X = \inf_{\substack{g,a\\ u \in K}} \sup \|u - g(a(u))\|_X$$

where a is a linear continuous map from K to \mathbb{R}^n extracting n linear information

$$a(u) = (\ell_1(u), \ldots, \ell_n(u)) \in \mathbb{R}^n$$

and g is an arbitrary nonlinear map from \mathbb{R}^n to X.

It is closely related with Gelfand widths

$$d^n(K)_X = \inf_{a \in L(X;\mathbb{R}^n)} \sup_{v \in K \cap Ker(a)} \|v\|_X$$

which are such that

$$s_n(K)_X \leq d^n(K-K)_X \leq 2s_n(K)_X$$

and $s_n(K)_X = d^n(K)_X$ when K convex and centrally symmetric.

For $X = L^p(\mathcal{X})$, $\mathcal{X} = [0, 1]^d$, and K the unit ball of Sobolev spaces $W^{s,q}$ or Besov spaces $B_a^s(L^{\tau})$ which compactly embed in L^p

$$\delta_n(K)_X \sim n^{-s/d}$$

Rate $O(n^{-s/d})$ is achieved for a larger class of functions than for linear methods (functions with regularity measured in norms weaker than L^p).

Optimal performance is achieved by free knot splines or best n-term approximation with a dictionary of tensor products of dilated splines.

For $X = L^p(\mathcal{X})$, $\mathcal{X} = [0, 1]^d$, and K the unit ball of Sobolev spaces $W^{s,q}$ or Besov spaces $B_a^s(L^{\tau})$ which compactly embed in L^p

$$\delta_n(K)_X \sim n^{-s/d}$$

Rate $O(n^{-s/d})$ is achieved for a larger class of functions than for linear methods (functions with regularity measured in norms weaker than L^p).

Optimal performance is achieved by free knot splines or best n-term approximation with a dictionary of tensor products of dilated splines.

Again, we observe the curse of dimensionality, which can not be avoided by such nonlinear methods.

- 1 Optimal linear approximation
- 2 Optimal nonlinear approximation
- 3 How to beat the curse of dimensionality ?

Could extra regularity help ?

Consider
$$X = L^{\infty}(\mathcal{X})$$
 with $\mathcal{X} = [0, 1]^d$ and
 $\mathcal{K} = \{ v \in C^{\infty}([0, 1]^d) : \sup_{\alpha} \|D^{\alpha}u\|_{L^{\infty}} < \infty \},$

It holds

$$K \subset B(W^{sd,\infty}) \quad \forall s > 0,$$

so that for all s > 0, the sensing numbers

$$s_n(K)_{L^{\infty}} \lesssim n^{-s}$$
.

Could extra regularity help ?

Consider
$$X = L^{\infty}(\mathcal{X})$$
 with $\mathcal{X} = [0, 1]^d$ and
 $\mathcal{K} = \{ v \in C^{\infty}([0, 1]^d) : \sup_{\alpha} \|D^{\alpha}u\|_{L^{\infty}} < \infty \},$

It holds

$$K \subset B(W^{sd,\infty}) \quad \forall s > 0,$$

so that for all s > 0, the sensing numbers

 $s_n(K)_{L^{\infty}} \lesssim n^{-s}.$

However, it holds [Novak and Wozniakowski 2009]

$$s_n(K)_{L^{\infty}} = 1$$
 for all $n = 0, 1, \dots, 2^{\lfloor d/2 \rfloor} - 1$

or

$$\min\{n: s_n(K)_{L^{\infty}} < 1\} \geq 2^{\lfloor d/2 \rfloor}$$

This even holds when information is chosen adaptively.

Could extra regularity help ?

Consider
$$X = L^{\infty}(\mathcal{X})$$
 with $\mathcal{X} = [0, 1]^d$ and
 $\mathcal{K} = \{ v \in C^{\infty}([0, 1]^d) : \sup_{\alpha} \|D^{\alpha}u\|_{L^{\infty}} < \infty \},$

It holds

$$K \subset B(W^{sd,\infty}) \quad \forall s > 0,$$

so that for all s > 0, the sensing numbers

 $s_n(K)_{L^{\infty}} \lesssim n^{-s}.$

However, it holds [Novak and Wozniakowski 2009]

$$s_n(K)_{L^{\infty}} = 1$$
 for all $n = 0, 1, \dots, 2^{\lfloor d/2 \rfloor} - 1$

or

$$\min\{n: s_n(K)_{L^{\infty}} < 1\} \ge 2^{\lfloor d/2 \rfloor}$$

This even holds when information is chosen adaptively.

Extra regularity can not help...

More assumptions on model classes K are needed...

Consider a parameter-dependent equation

$$\mathcal{P}(u(y); y) = 0, \quad u(y) \in X$$

with $y \in Y$ some parameter.

The objective is to approximate the solution manifold (model reduction methods)

$$K = \{u(y) : y \in Y\}$$

or to approximate explicitly the solution map $y \mapsto u(y)$.

Consider a parameter-dependent equation

$$\mathcal{P}(u(y); y) = 0, \quad u(y) \in X$$

with $y \in Y$ some parameter.

The objective is to approximate the solution manifold (model reduction methods)

$$K = \{u(y) : y \in Y\}$$

or to approximate explicitly the solution map $y \mapsto u(y)$.

As an example, consider the elliptic diffusion equation on a convex domain $D \subset \mathbb{R}^d$

$$-div(a(y)\nabla u(y)) = f$$

with $f \in H^{-1}$, $0 < \underline{a} \le a(y) \le \overline{a} < \infty$, and homogeneous Dirichlet boundary conditions. The solutions

$$u(y)\in H_0^1:=X.$$

• Assuming $f \in L^2$ and a(y) sufficiently smooth, we know that K is in some ball of $H^2(D)$, so that

$$d_n(K)_{H^1} \lesssim n^{-1/d}$$

with optimal performance achieved by splines (finite elements with uniform mesh).

• Assuming $f \in L^2$ and a(y) sufficiently smooth, we know that K is in some ball of $H^2(D)$, so that

$$d_n(K)_{H^1} \lesssim n^{-1/d}$$

with optimal performance achieved by splines (finite elements with uniform mesh).

• If $a(y) = a_0 + \sum_{i=1}^m a_i y_i$ with $(||a_i||_{L^{\infty}})_{i \ge 1} \in \ell_p$ for some p > 1, then

$$d_n(K)_{H^1} \leq Cn^{-s}, \quad s = p^{-1} - 1$$

with constant C independent of d (no curse of dimensionality).

These rates are achieved by sparse polynomial expansions of $y \mapsto u(y)$, exploiting anisotropic analyticity of the solution map.

• Assuming $f \in L^2$ and a(y) sufficiently smooth, we know that K is in some ball of $H^2(D)$, so that

$$d_n(K)_{H^1} \lesssim n^{-1/d}$$

with optimal performance achieved by splines (finite elements with uniform mesh).

• If $a(y) = a_0 + \sum_{i=1}^m a_i y_i$ with $(||a_i||_{L^{\infty}})_{i \ge 1} \in \ell_p$ for some p > 1, then

$$d_n(K)_{H^1} \leq Cn^{-s}, \quad s = p^{-1} - 1$$

with constant C independent of d (no curse of dimensionality).

These rates are achieved by sparse polynomial expansions of $y \mapsto u(y)$, exploiting anisotropic analyticity of the solution map.

• More generally, letting $\mathcal{A} = \{a(y) : y \in Y\}$, we have [Cohen and DeVore 2015]

$$\sup_{n \geq 1} n^s d_n(K)_{H^1} \lesssim \sup_{n \geq 1} n^r d_n(\mathcal{A})_{L^{\infty}}, \quad \forall s < r-1.$$

• Assuming $f \in L^2$ and a(y) sufficiently smooth, we know that K is in some ball of $H^2(D)$, so that

$$d_n(K)_{H^1} \lesssim n^{-1/d}$$

with optimal performance achieved by splines (finite elements with uniform mesh).

• If $a(y) = a_0 + \sum_{i=1}^m a_i y_i$ with $(||a_i||_{L^{\infty}})_{i \ge 1} \in \ell_p$ for some p > 1, then

$$d_n(K)_{H^1} \leq Cn^{-s}, \quad s = p^{-1} - 1$$

with constant C independent of d (no curse of dimensionality).

These rates are achieved by sparse polynomial expansions of $y \mapsto u(y)$, exploiting anisotropic analyticity of the solution map.

• More generally, letting $\mathcal{A} = \{a(y) : y \in Y\}$, we have [Cohen and DeVore 2015]

$$\sup_{n \geq 1} n^{s} d_{n}(\mathcal{K})_{H^{1}} \lesssim \sup_{n \geq 1} n^{r} d_{n}(\mathcal{A})_{L^{\infty}}, \quad \forall s < r-1.$$

• Optimal spaces X_n are data-dependent. Almost optimal spaces can be constructed using greedy algorithms (reduced basis methods) or sparse polynomial expansions.

• Assuming $f \in L^2$ and a(y) sufficiently smooth, we know that K is in some ball of $H^2(D)$, so that

$$d_n(K)_{H^1} \lesssim n^{-1/d}$$

with optimal performance achieved by splines (finite elements with uniform mesh).

• If $a(y) = a_0 + \sum_{i=1}^m a_i y_i$ with $(||a_i||_{L^{\infty}})_{i \ge 1} \in \ell_p$ for some p > 1, then

$$d_n(K)_{H^1} \leq Cn^{-s}, \quad s = p^{-1} - 1$$

with constant C independent of d (no curse of dimensionality).

These rates are achieved by sparse polynomial expansions of $y \mapsto u(y)$, exploiting anisotropic analyticity of the solution map.

• More generally, letting $\mathcal{A} = \{a(y) : y \in Y\}$, we have [Cohen and DeVore 2015]

$$\sup_{n \geq 1} n^{s} d_{n}(K)_{H^{1}} \lesssim \sup_{n \geq 1} n^{r} d_{n}(\mathcal{A})_{L^{\infty}}, \quad \forall s < r-1.$$

- Optimal spaces X_n are data-dependent. Almost optimal spaces can be constructed using greedy algorithms (reduced basis methods) or sparse polynomial expansions.
- Similar results between nonlinear widths $\delta_n(K)_{H^1}$ and $\delta_n(\mathcal{A})_{L^q}$.

• No (reasonable) approximation tool is able to overcome the curse of dimensionality for standard regularity classes.

- No (reasonable) approximation tool is able to overcome the curse of dimensionality for standard regularity classes.
- The key is to make more assumptions on model classes of functions and to provide ad-hoc approximation tools.

- No (reasonable) approximation tool is able to overcome the curse of dimensionality for standard regularity classes.
- The key is to make more assumptions on model classes of functions and to provide ad-hoc approximation tools.
- We would like flexible approximation tools X_n that perform well for a wide range of applications, that usually require a high number of parameters (Part 2),

- No (reasonable) approximation tool is able to overcome the curse of dimensionality for standard regularity classes.
- The key is to make more assumptions on model classes of functions and to provide ad-hoc approximation tools.
- We would like flexible approximation tools X_n that perform well for a wide range of applications, that usually require a high number of parameters (Part 2),
- Or construct "low-dimensional" sets X_n that approximate a specific class K (Part 3).

References I

Approximation theory



A. Pinkus.

N-widths in Approximation Theory, volume 7. Springer Science & Business Media, 2012.

R. A. DeVore, R. Howard, and C. Micchelli. Optimal nonlinear approximation. *Manuscripta mathematica*, 63(4):469–478, 1989.



R. A. DeVore and G. G. Lorentz.

Constructive approximation, volume 303. Springer Science & Business Media, 1993.



R. A. DeVore.

Nonlinear approximation.

Acta Numerica, 7:51-150, 1998.

J. Creutzig.

Relations between Classical, Average, and Probabilistic Kolmogorov Widths.

J. Complexity, 18:287-303, Mar. 2002.

References II



V. N. Temlyakov.

Nonlinear Kolmogorov widths. Math. Notes, 63(6):785–795, June 1998.



V. Temlyakov.

On optimal recovery in L2. Journal of Complexity, 65:101545, 2021.



N. Nagel, M. Schäfer, and T. Ullrich.

A new upper bound for sampling numbers. Foundations of Computational Mathematics, pages 1–24, 2021.



Optimal pointwise sampling for l^2 approximation, 2021.



M. Dolbeault, D. Krieg, and M. Ullrich.

A sharp upper bound for sampling numbers in L_2 , 2022.



D. Krieg, E. Novak, and M. Ullrich.

On the power of adaption and randomization, 2024.

References III



A. Cohen, R. DeVore, G. Petrova, and P. Wojtaszczyk.
 Optimal Stable Nonlinear Approximation.
 Found. Comput. Math., 22(3):607–648, June 2022.

High-dimensional approximation and model reduction



D. Dũng, V. N. Temlyakov, and T. Ullrich.

Hyperbolic Cross Approximation. arXiv e-prints, page arXiv:1601.03978, Jan. 2016.

A. Cohen and R. DeVore.

Approximation of high-dimensional parametric pdes. *Acta Numerica*, 24:1–159, 2015.



P. Benner, A. Cohen, M. Ohlberger, and K. Willcox, editors. *Model Reduction and Approximation: Theory and Algorithms.* SIAM, Philadelphia, PA, 2017.

E. Novak and H. Woźniakowski.

Approximation of infinitely differentiable multivariate functions is intractable. *Journal of Complexity*, 25(4):398–404, 2009.

Proofs of some results

Theorem

Let K be a subset of X and bc(K) its balanced convex hull defined by

$$bc(\mathcal{K}) = \{\sum_{i=1}^{m} \alpha_i u_i : u_1, \dots, u_m \in \mathcal{K}, \sum_{i=1}^{m} |\alpha_i| \leq 1, m \in \mathbb{N}\}$$

It holds $d_n(K)_X = d_n(bc(K))_X$.

Proof.

The inclusion $K \subset bc(K)$ implies $d_n(K)_X \leq d_n(bc(K))_X$. Then let X_n be an optimal space such that $d_n(K)_X = \sup_{u \in K} \inf_{v \in X_n} ||u - v||_X$. Then

$$d_n(bc(\mathcal{K}))_X \leq \sup_{u \in bc(\mathcal{K})} E(u, X_n)_X = \sup_m \sup_{u_1, \dots, u_m} \sup_{|\alpha_1| + \dots + |\alpha_m| \leq 1} \inf_{v \in X_n} \|\sum_{i=1}^n \alpha_i u_i - v\|_X$$
$$\leq \sup_m \sup_{u_1, \dots, u_m} \sup_{|\alpha_1| + \dots + |\alpha_m| \leq 1} \sum_{i=1}^m |\alpha_i| \|u_i - v_i\|_X$$

with $v_i \in X_n$ such that $\|u_i - v_i\|_X = \inf_{v \in X_n} \|u_i - v\|_X \le \sup_{u \in K} \inf_{v \in X_n} \|u - v\|_X = d_n(K)_X$. We deduce $d_n(bc(K))_X \le d_n(K)_X$.

m
Theorem

For any $n \ge 1$, it holds

 $d_n(K)_X \ge b_n(K)_X$

Proof.

Let $\lambda > 0$ and X_{n+1} such that $\lambda B(X_{n+1}) \subset K$. It holds

$$d_n(K)_X \geq d_n(\lambda B(X_{n+1}))_X = \lambda d_n(B(X_{n+1}))_X,$$

and from [Pinkus 2012, Theorem 1.5], $d_n(B(X_{n+1}))_X = 1$. Theorefore, it holds $d_n(K)_X \ge \lambda$, and the result follows by taking the supremum over λ and X_{n+1} .

Proofs of some results

Theorem

$$s_n(K)_X \leq d^n(K-K)_X \leq 2s_n(K)_X$$

Proof.

$$s_n(K)_X = \inf_{a \in L(X;\mathbb{R}^n)} \inf_{g} \sup_{u \in K} ||u - g(a(u))||_X$$

$$\leq \inf_{a \in L(X;\mathbb{R}^n)} \sup_{u,v \in K, a(u) = a(v)} ||u - v||_X = d^n(K - K)_X$$

that is the first inequality. Now consider a^* and g^* realizing the infimum $s_n(K)_X$. It holds

$$d^{n}(K - K)_{X} \leq \sup_{\substack{u, v \in K, a^{*}(u) = a^{*}(v) \\ u, v \in K}} \|u - g^{*}(a^{*}(u))\|_{X} + \|g^{*}(a^{*}(v)) - v\|_{X}$$
$$\leq \sup_{\substack{u, v \in K \\ u, v \in K}} \|u - g^{*}(a^{*}(u))\|_{X} + \|g^{*}(a^{*}(v)) - v\|_{X}$$
$$= 2s_{n}(K)_{X}$$