

NEW (non-classical) METHODS
in
APPROXIMATION THEORY

BORIS KASHIN

(Steklov Math. Inst., MSU, MOSCOW)

ULM

September, 2013

Approximation theory —

more than 100 years of very active research.

"The main problem of approximation consists in finding, for a complicated function f from a large space X , a close-by, simple function φ from a small subset Φ of X "

R. DeVore, G Lorentz

"Constructive Approximation", Springer, 1993

Two sources of development:

- I) Needs of Applications;
- II) Aspiration to the Ideal which leads to the attempts to solve the Extremal Problems

Sometimes $I) \Rightarrow II)$ and $II) \rightarrow I)$



Classical period of A.T.:

large space X - space of functions

($C(\Delta)$, $L^p(\Delta)$, ...)

small subset - the set \mathcal{T}_n of all algebraic polynomials of degree $\leq n$, or trigonometric polynomials..., or rational functions...

Definition 1 For $f \in C[-1, 1]$, $n=1, 2, \dots$

$$E_n(f) \equiv \inf_{P \in \mathcal{T}_n} \|f - P\|_{C[-1, 1]}$$

The solution $P = P_n$ of the corresponding extremal problem is called the polynomial of best uniform approximation;

$$R_{n,n}(f) \equiv \inf_{Q \in R_{n,n}} \|f - Q\|_{C[-1, 1]}$$

P. L. Chebyshev (1855) trying to solve some applied problem gave a characterization of the polynomial of best uniform approximation to f from $C[-1, 1]$ (Chebyshev's alternance)

In the case when $f(x) = x^n$
he founded P_{n-1} explicitly:

$$x^n - P_{n-1} = 2^{-n+1} \cdot \cos(n \arccos x)$$

$$= \frac{1}{2} \left\{ (x + \sqrt{x^2 - 1})^n + (x - \sqrt{x^2 - 1})^n \right\}$$

- Chebyshev polynomials - the solution
of the extremal problem

$$\inf_{\{a_{n-1}, \dots, a_0\}} \|x^n + a_{n-1}x^{n-1} + \dots + a_0\|_{C[-1,1]} \rightarrow \min$$

In the classical A.T. important
results concerning approximation
order of concrete functions and
functional classes were obtained.

Some examples:

a) direct and inverse theorems for
polynomial approximation

Th. (Jackson 1912, Bernstein 1912) For $0 < \alpha < 1$

$$f \in \text{Lip } \alpha \iff E_n^T(f) = \underline{\underline{O}}(n^{-\alpha})$$

where f is 2π -periodic and

$$\text{Lip } \alpha = \left\{ f \in C(-\pi, \pi), \omega(f, \delta) \leq C\delta^\alpha \right\}$$

b) Bernstein, 1913 : there exists $\beta > 0$ such that

$$n \cdot E_n(|x|, C[-1,1]) \rightarrow \beta \text{ if } n \rightarrow \infty$$

(Bernstein conjecture : $\beta = 1/2\sqrt{\pi}$ - false, $\beta = 0.2801$). The same result for $E_n^T(|x|, (-\pi, \pi))$.

c) Herbert Stahl (1942-2013), in Sbornik Mathematics v.183, N8, 1992 :

$$e^{\pi\sqrt{n}} \cdot R_{n,n}(|x|, C^{\infty}[-1,1]) \rightarrow 8 \text{ if } n \rightarrow \infty.$$

Many more important extremal problems were studied and solved in the Classical Approximation Theory.

In the last 30-40 years the main stream of studies in A.T. moved away from Classical Analysis.

The main reason is the needs of Applied mathematics.

Now we start to consider the problems of Approximation in much more general setting.

Let B be a normed space, $\phi \subset B$ be subset of B - we will call it a dictionary, and f be any element of B

Definition 2 n -term approximation of $f \in B$ with respect to the dictionary ϕ is

$$e_n(f, \phi, B) \equiv \inf_{P \in \Sigma_n} \|f - P\|_B,$$

where for $n=1, 2, \dots$

$$\Sigma_n \equiv \left\{ \sum_{j=1}^n a_j x_j, a_j \in \mathbb{R}, x_j \in \phi \right\}$$

Further, if K is a subset of B , then

$$e_n(K, \phi, B) = \sup_{f \in K} e_n(f, \phi, B)$$

Nowdays the notion of n -term approximation is one of the basic notion of Approximation Theory.

Various dictionaries have been considered. Here are some examples important for application:

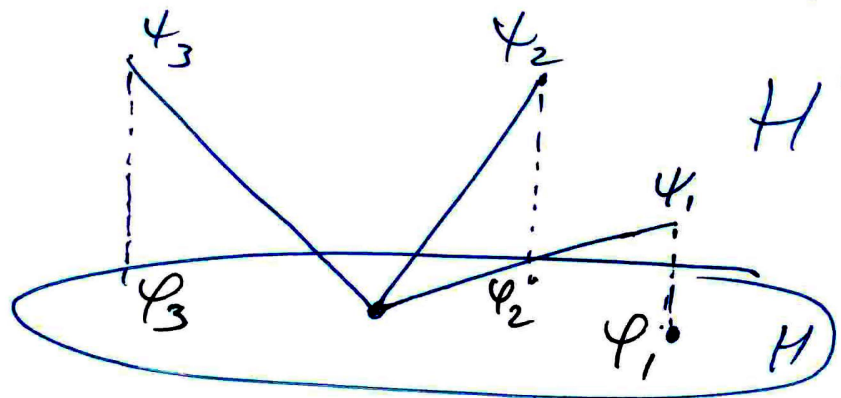
a) Φ is a complete orthonormal set of elements (O.N.S) in some Hilbert Space H (H could be the space of functions $L^2(\Omega)$); in particular, Φ can be a trigonometric system

a') Φ is a frame in H

($\Phi = \{\psi_i\}$ is a frame if for some constants $A > 0, B$ and any $f \in H$

$$A \|f\|_H^2 \leq \sum_{i=1}^{\infty} \langle f, \psi_i \rangle^2 \leq B \|f\|_H^2;$$

tight frame if $A = B$)



$\overline{\mathcal{I}}_{H' \rightarrow H}$
operator of
orthogonal projection

$\{\psi_i\}$ - O.N.S. in H' ; $H \subset H'$, $\psi_i = \overline{\mathcal{I}}_{H' \rightarrow H}(\psi_i)$

b) for the approximation of functions of d variables, defined on the cube $I^d = (0, 1)^d$ we can consider the dictionary \mathcal{D}_k that consists of all functions of the type

$$u(x_1, \dots, x_k) \cdot v(x_{k+1}, \dots, x_d); \quad 1 \leq k \leq d;$$

c) for the approximation of functions of d variables by "free knots splines" we can consider dictionaries that consist of functions of the type

$$P \cdot \chi_{\Delta}$$

where P is a polynomial of the degree $\leq r$ of d variables and χ_{Δ} is the characteristic function of the segment $\Delta \subset \mathbb{R}^d$

d) \mathcal{D} is the set of ridge functions, i.e. functions in $L^p(\Omega)$, $\Omega \subset \mathbb{R}^d$ of the type

$$u(x) = f(\langle x, \theta \rangle)$$

where f is a function of one variable, $x \in \Omega$, $\theta \in \mathbb{R}^d$, $|\theta| = 1$

e) \mathcal{D} is the set of Gabor functions

$$\{g_{a,b}(x-c); g_{a,b}(x) = e^{iax} \cdot e^{-bx^2}; a, c \in \mathbb{R} \quad b > 0\}$$

For each of those families of dictionaries there are results in A.T. that are valuable from the theoretical as well as practical point of view.

Following examples support the concept of n-term approximation.

1) Earlier we mentioned that

$$E_n^T(|x|, C[-\pi, \pi]) \asymp \frac{1}{n} \text{ if } n \rightarrow \infty$$

Theorem (Ismagilov, Maiorov... 1970th)

$$e_n(|x|, T, C[-\pi, \pi]) \asymp \frac{1}{n^{3/2}} \text{ if } n \rightarrow \infty$$

In this case the construction of the approximating polynomial is explicit.

So classical approximation is not the best in this case.

- 11 -

2) If ϕ - complete O.N.S. in H
and $f \in H$ then

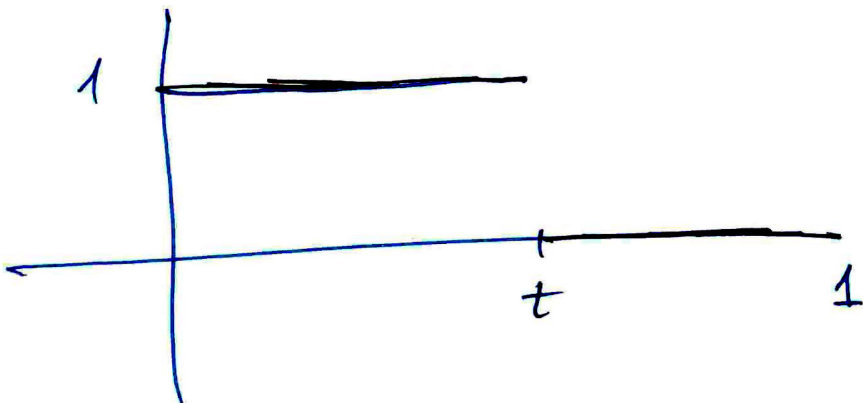
$$e_n(f, \phi, H) = \left(\sum_{k \geq n+1} |c_k^*(f)|^2 \right)^{1/2},$$

where $\{c_k^*(f)\}_{k=1}^{\infty}$ is a non-increasing rearrangement of the sequence of absolute values of the Fourier coefficients of the function f with respect to complete O.N.S. ϕ .

In many cases n -term approximation works much better than classical one.

Example. Let $\chi = \{\chi_t\}_{0 < t < 1} \subset L^2(0,1)$
where

$$\chi_t(x) = \begin{cases} 1, & \text{if } 0 \leq x \leq t \\ 0, & \text{if } x > t \end{cases}$$



Following result is well known

Theorem (Rudin, 1950th) For any O.N.S. $\Psi = \{\Psi_k\} \subset L^2(0,1)$

$$E_n(X, \Psi, L^2(0,1)) \geq c n^{-1/2},$$

$c > 0, n = 1, 2, \dots$

But if $\Psi = \mathcal{H} = \{h_k\}$ is Haar system, then

$$e_n(X, \mathcal{H}, L^2(0,1)) \leq C 2^{-n/2},$$

$n = 1, 2, \dots$

Haar system is historically first and simplest wavelet system.

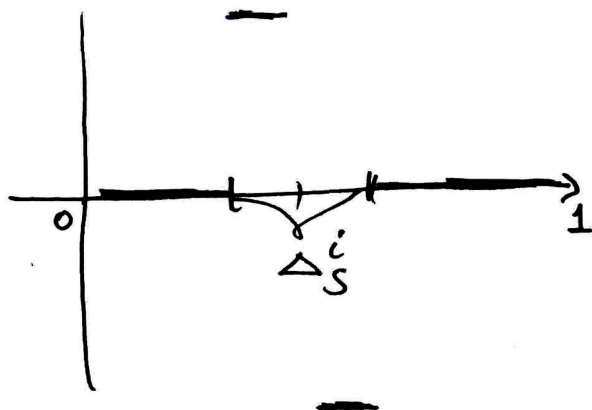
Haar functions: $h_i \equiv 1$ and if $k = 2^s + i$, $i = 1, \dots, 2^s$, $s = 0, 1, \dots$ then

$$h_k(x) = \begin{cases} 0, & \text{if } x \notin \Delta_s^i \\ 2^{s/2}, & \text{if } x \in (\Delta_s^i)^+ \\ -2^{s/2}, & \text{if } x \in (\Delta_s^i)^- \end{cases}$$

where

$$\Delta_s^i = \left(\frac{i-1}{2^s}, \frac{i}{2^s}\right) \text{ and}$$

$(\Delta_s^i)^+$, $(\Delta_s^i)^-$ - left and right halves of Δ_s^i correspondingly



3) Let us remark that in the previous example the decomposition of any $f \in X$ into Haar series is very sparse.

In general Sparsity is a very important concept in A.T. and applications.

If any f from functional class K has sparse representation with respect to O.N.S. ϕ then n -term approximation works very well and we can get very good estimates of

$$e_n(K, \phi, L^2(\Omega)).$$

Current International Standard of Images Compression JPEG 2000 is based on n -term approximation with respect to some 2-d wavelet system Ψ because engineers agree that usual picture has sparse representation in this basis

(usual model of 2-d digital picture is a discrete function of two variables)

Let us consider the problem of approximation of the set

$$K_2 = \{ \chi_A, A \subset [0,1]^2 \text{-convex} \} \subset L^2(0,1)^2$$

This problem is very natural from practical point of view.

For classical O. N. S. ϕ

$$e_n(K_2, \phi, L^2(0,1)^2) \geq C n^{-1/2}, \quad n=1,2,\dots$$

Theorem (Candes, Donoho, 2002)

There exists tight frame ϕ such that

$$e_n(K_2, \phi, L^2(0,1)^2) \leq C \frac{\ln n}{n^{3/2}}, \quad n=1,2,\dots$$

(Here for any $f \in K_2$ we consider expansion

$$f = \sum_{k=1}^{\infty} \langle f, \psi_k \rangle \psi_k$$

and take as n -term approximation the polynomial of the type

$$\sum_{k \in \Lambda_n} \langle f, \psi_k \rangle \psi_k, \quad \#\Lambda_n = n$$

with n biggest (in absolute value) coefficients.

Last result is almost optimal.

To prove it we need a method to get lower estimates for n-term approximation. In the case of orthonormal dictionary (or frame) the following theorem is widely used

Theorem (noncompressibility of hypercubes
B. Kashin, 1985)

Let $\phi \subset H$ - O.N.S. Suppose that for some $n=1,2,\dots$ the set K contains the set of vertices of hypercube:

$$K \supset Q = \left\{ \sum_{i=1}^{2n} \varepsilon_i \psi_i, \varepsilon_i = \pm 1, \{\psi_i\}_{i=1}^{2n} \text{ - O.N.S.} \right\}$$

then

$$e_n(K, \phi, H) \geq c \cdot n^{1/2},$$

$c > 0$ - absolute constant.

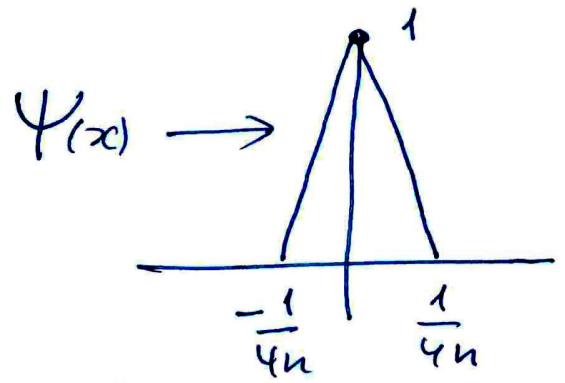
How it works? Let us show using Theorem above that for any O.N.S. ϕ

$$e_n(\text{Lip } \alpha, \phi, L^2(0,1)) \geq c \cdot n^{-\alpha}, \quad 0 < \alpha \leq 1.$$

This result shows that in this case n -term approximation is not better than usual approximation by trigonometric polynomials of degree n .

For a given n consider the set of $2n$ functions

$$\Psi_i(x) = \Psi(x - \frac{i}{2n})$$
$$i = 1, 2, \dots, 2n$$



It is easy to check that for $i = 1, 2, \dots, 2n$

$$\|\Psi_i\|_{L^2(0,1)} = (6n)^{-1/2}; \quad \Psi_i (4n)^{-d} \in \text{Lip } d$$

and \Rightarrow for any signs $\varepsilon_i = \pm 1$

$$(8n)^{-d} \cdot \sum_{i=1}^{2n} \varepsilon_i \Psi_i \in \text{Lip } d$$

So

$$(8n)^d \cdot \text{Lip } d \supset (6n)^{-1/2} \underbrace{\sum_{i=1}^{2n} \varepsilon_i \Psi_i (6n)^{1/2}}_{\text{hypercube}}$$

Th. $\Rightarrow \sup_{f \in \text{Lip } d} e_n(f, \phi, L^2(0,1)) \geq C n^{-d}$ □

GREEDY APPROXIMATION

In Applications of A.T. the results on existence of good approximation are not enough. We need a method of constructing such approximation.

That is why following topic became very popular.

We consider only the basic results concerning Greedy Approximation. (See also V. Temlyakov, "Greedy Approximation", Cambridge Univ. Press, 2011)

Below H is a Hilbert space, $D \subset H$ is a dictionary. Suppose that

$$\|g\|_H = 1 \text{ for any } g \in D$$

For $f \in H$ try to find best 1-term approximation with respect to dictionary D

This best 1-term approximation is

$$\langle f, g_0 \rangle g_0$$

where g_0 is a solution of extremal problem

$$|\langle f, g \rangle| \rightarrow \max, g \in D \quad (*)$$

(we maximize the angle between f and elements of dictionary and suppose that this solution exists)

The idea of Greedy Approximation is to construct inductively approximants to f such way that at each step we take the best possible "1-term correction" to the already constructed. So on n -th step of Greedy Algorithm we have special n -term approximation and we expect this sequence of approximants to converge to f .

To be precise define:

$$G(f) := G(f, D) = \langle f, g_0 \rangle g_0$$

where g_0 is a solution of extremal problem $(*)$;

$$R(f) := R(f, D) = f - G(f);$$

$$G_0(f) = 0, \quad f_0 \equiv R_0(f) := f$$

Then for $m=1, 2, \dots$ we inductively define

$$G_m(f) = G_{m-1}(f) + G(R_{m-1}(f));$$

$$f_m = R_m(f) = f - G_m(f) = R(R_{m-1}(f)).$$

Obviously

$$G_m(f) \in \sum_m = \left\{ \sum_{j=1}^m a_j g_j, a_j \in \mathbb{R}, g_j \in D \right\}.$$

Theorem (Jones, 1987) For any D and $f \in H$

$$G_m(f) \xrightarrow{H} f \text{ if } m \rightarrow \infty.$$

Another version of G. A.

(weak greedy algorithm):

fix $t: 0 < t < 1$ and instead of $G(f)$ defined by (*) take any element $\tilde{g} \in D$ such that

$$|\langle f, \tilde{g} \rangle| \geq t \sup_{g \in D} |\langle f, g \rangle|$$

If $\tilde{G}(f) = \langle f, \tilde{g} \rangle \tilde{g}$, $\tilde{R}(f) = f - \tilde{G}(f)$

then inductive process completely similar to G. A. converges to f for any $f \in H$ and any dictionary D .

If $D = \phi = \{\psi_k\}$ - O.N.S.
and $f = \sum a_k \psi_k$ (**)

then

$$G_m(f, \phi) = \sum_{k \in \Lambda_m} a_k \psi_k$$

where

Λ_m - the support of m "biggest" coefficients of the series (**):

$$\forall k \in \Lambda_m, \forall k' \notin \Lambda_m \quad |a_k| \geq |a_{k'}|.$$

Rate of convergence of G. A.

To precise the setting of the problem we need the following notations:

$$A_1^{\circ}(D, M) =$$

$$= \left\{ f \in H : f = \sum_{k \in \Lambda} c_k g_k, g_k \in D, \sum_{k \in \Lambda} |c_k| \leq M, \#\Lambda < \infty \right\}$$

$$A_1(D, M) = \text{closure of } A_1^{\circ}(D, M)$$

If $D = \{\psi_k\}$ is O.N.S, then

$$\bigcup_M A_1(D, M) = \left\{ f \in H : \sum a_k \psi_k = f, \sum |a_k| < \infty \right\}$$

In this case for any $f \in A_1(D, M)$

$$\|f - G_m(f)\|_H \leq C \cdot M \cdot m^{-1/2}, m=1,2,\dots$$

What happens in the case of general dictionary? It turns out that the rate we can guarantee is about $m^{-1/2}$ only.

The first result was obtained by De Vore, Temlyakov:

for any D and $f \in A_1(D, M)$

$$\|f - G_m(f)\|_H \leq C \cdot M \cdot m^{-1/6}, \quad m=1,2,\dots$$

Best known result (Sil'nichenko)

$$1/6 \rightarrow 0.182$$

Lower estimate (Livshitz, Temlyakov):

$\exists D, f$ such that

$$\lim_{m \rightarrow \infty} \|f - G_m(f)\|_H \cdot m^{0.1898} > 0.$$

A.T. in Compressed Sensing

Compressed Sensing is new and very dynamic branch of signal processing with wide spectrum of Applications. Crucial role in this area is played by Kolmogorov widths estimates.

Definition 3 If E is a normed space,

$F \subset E$, then for $n = 0, 1, \dots$ the Kolmogorov width is

$$d_n(F, E) = \inf_{\substack{L \subset E \\ \dim L \leq n}} \sup_{f \in F} \text{dist}_E(f, L),$$

where $\text{dist}_E(f, L) = \inf_{y \in L} \|f - y\|_E$.

This definition was proposed by A.N. Kolmogorov in 1936.

Starting from 1970s a lot of widths estimates were obtained in A.T.

Recently it became clear that this results could be successfully applied in signal processing.

Compressed Sensing refers to the problem of economical recovery of an unknown vector $u \in \mathbb{R}^N$ from the information provided by linear measurements $\langle u, \varphi_j \rangle$, $\varphi_j \in \mathbb{R}^N$, $j=1, 2, \dots, n$ in the case when

N is much bigger than n

The natural variant of such a setting uses the concept of sparsity.

Notations:

$$\ell_p^N = (\mathbb{R}^N, \|\cdot\|_{\ell_p^N})$$

where for $x = \{x_j\}_{j=1}^N \in \mathbb{R}^N$

$$\|x\|_{\ell_p^N} = \begin{cases} (\sum_{j=1}^N |x_j|^p)^{1/p}, & \text{if } 1 \leq p < \infty \\ \max_{1 \leq j \leq N} |x_j|, & \text{if } p = \infty \end{cases}$$

$$B_p^N = \{x \in \mathbb{R}^N : \|x\|_{\ell_p^N} = 1\}.$$

We call a vector $u \in \mathbb{R}^N$ k -sparse if it has at most k nonzero coordinates.

We start with a problem: to find for a given pair (N, n) the biggest sparsity $k(N, n)$ such that there exists a set of vectors

$$\varphi_j \in \mathbb{R}^N, j=1, \dots, n$$

and an economical algorithm $A = A_{\alpha}$ mapping vector

$$y = (\langle u, \varphi_1 \rangle, \dots, \langle u, \varphi_n \rangle) \in \mathbb{R}^n$$

into \mathbb{R}^N in such a way that for any u of sparsity $k(N, n)$ one would have an exact recovery

$$A(y(u)) = u.$$

In other words we want to describe matrices A with rows $\varphi_j \in \mathbb{R}^N, j=1, \dots, n$ such that there exists an economical algorithm of solving the above formulated sparse recovery problem.

$$\mathbb{R}^n \ni \begin{array}{|c} y \end{array} = \begin{array}{|c|} \hline \phi \\ \hline \end{array} \times \begin{array}{|c} u \in \mathbb{R}^n \end{array}$$

$$\begin{cases} u \text{ is } k\text{-sparse} \\ A_{\phi}(y) = u \end{cases} \quad (***)$$

If $2k \leq n$ it is easy to construct ϕ and noneconomical algorithm

A such that $(***)$ holds true

But such a solution is useless for Applications

(Such Algorithm is connected with the solution of extremal problem

where $\|v\|_0 := |\text{supp}(v)|$, which is NP hard) P_0

In 2005-2006 Donoho and Candes, Romberg, Tao proved that if

$$K \leq C \frac{n}{\ln \frac{2N}{n}} \quad \text{then} \quad (****)$$

there exists matrix Φ such that P_0 is equivalent to the following extremal problem

$$\|v\|_{\ell_1} \rightarrow \min \quad \text{subject to } \Phi v = y \quad P_1$$

P_1 can be solved by linear programming techniques in a polynomial time

The estimate $(****)$ is best possible.

These results can be obtained using some widths estimates.

Notation For a positive a and $v \in \mathbb{R}^N$.

$$\sigma_a(v)_1 := \min_{\substack{w \in \mathbb{R}^N \\ |\text{supp}(w)| \leq a}} \|v - w\|_{\ell_1^N}$$

We say that a measurement matrix Φ has **Strong Compressed Sensing Property (SCSP)**

if for any $u \in \mathbb{R}^N$ we have

$$\|u - A_\Phi(\Phi u)\|_{\ell_2^N} \leq (k^{-1/2} \cdot \sigma_k(u))_1$$

for $k \leq n / \log \frac{2N}{n}$

$$\text{SCSP} \implies (***)$$

Theorem (Kashin, Temlyakov 2007)

Φ has SCSP if and only if

$$\text{dist}_{\ell_\infty^N} (B_2^N, L_\Phi) \leq \frac{C'}{n^{1/2}} \log \frac{2N}{n}$$

where L_Φ is the subspace in \mathbb{R}^N generated by rows of Φ .

The existence of ϕ such as in the theorem above follows from width estimates

Theorem (Kashin, 1977 improved by Garnuev, Gluskin, 1989,

$$d_n(B_2^N, \ell_\infty^N) \leq \frac{C}{n^{1/2}} \log^{1/2} \frac{2N}{n}$$

(and random n -dimensional subspace in \mathbb{R}^N provides this estimate)