

Low-Rank Structures and Tensor Approximations for Huge-Size Data Sets

Eugene Tyrtyshnikov

(Invited Paper)

Abstract

This paper suggests an elementary introduction to recent developments of low-rank matrix approximation techniques using an appropriate small sample of all data, both in the bivariate case in the classical matrix theory and in a much more involved case where the dimension is three and higher.

Index Terms

Multidimensional matrix, tensor decomposition, tensor rank, low-rank approximation, cross approximation, maximal-volume principle.

I. INTRODUCTION

Consider a multiplication problem for a *fixed* matrix A of order n and an *arbitrary* vector x . Given x on input, we are interested in an algorithm producing $y = Ax$ on output. For instance, if $A = F_n$ is the Fourier matrix of order n , then $y = Ax$ can be found in $O(n \log_2 n)$ arithmetic operations by the famous FFT, which delivers the exact y in case there are no round-offs. In practice, however, we may be satisfied with a certain ε -approximation to y . The complexity of *approximate algorithms* should obviously depend both on n and on ε .

The entries a_{ij} of A are typically defined as the values of a function $f(u, v)$ at the nodes $u = u_i, v = v_j$, where u_1, \dots, u_n and v_1, \dots, v_n are grids in a k -dimensional space:

$$a_{ij} = f(u_i, v_j), \quad 1 \leq i, j \leq n.$$

Manuscript received May 30, 2007.

The author is with the Institute of Numerical Mathematics of the Russian Academy of Sciences, Lomonosov Moscow State University and Moscow Institute of Physics and Technology.

Let all the nodes u_i, v_j belong to a domain $\mathcal{D} \subset \mathbb{R}^k$, and assume that for any $\varepsilon > 0$ the function $f(u, v)$ admits separation of variables

$$f(u, v) \approx \sum_{s=1}^r \phi_s(u) \psi_s(v), \quad r = r(\varepsilon),$$

where

$$|f(u, v) - \sum_{s=1}^r \phi_s(u) \psi_s(v)| \leq \varepsilon, \quad u, v \in \mathcal{D}.$$

Then the matrix A is approximated by a matrix A_r of the form

$$A_r = \sum_{s=1}^r \begin{bmatrix} \phi_s(u_1) \\ \dots \\ \phi_s(u_n) \end{bmatrix} \begin{bmatrix} \psi_s(v_1) & \dots & \psi_s(v_n) \end{bmatrix} \quad (1)$$

with the entry-wise accuracy

$$|a_{ij} - (A_r)_{ij}| \leq \varepsilon, \quad 1 \leq i, j \leq n.$$

Evidently, $Ax \approx A_r x$, and the multiplication of A_r by a vector x requires only $O(nr)$ operations.

As we see, the number of operations depend upon n *linearly*. It is important as well to investigate the dependence of r upon ε . This question pertains to approximation theory and may involve some delicate techniques of analysis. Nevertheless, some estimates on r can be derived by very simple means.

For example, set $k = 1$ and let $\mathcal{D} = [a, b]$ be an interval on the real axis. Assume that a function $f(u, v)$ has infinitely many derivatives in v for any fixed u . Then, if we choose and fix any u , the function $f(u, v)$ can be expanded into the Taylor series at the point $v = v_0 = (a+b)/2$:

$$f(u, v) = \sum_{s=0}^{r-1} \left. \frac{\partial^s f}{\partial v^s} \right|_{v=v_0} \frac{(v - v_0)^s}{s!} + E_r(u, v),$$

where $E_r(u, v)$ is the residue term. When neglecting the latter, we arrive at some approximation in the additive form containing r terms with separated variables u and v . Consider f as a function of v . If it belongs to the class of infinitely smooth functions for which a derivative of an order s is bounded in modulus by some quantity M^s , where M is a positive constant the same for all $u \in \mathcal{D}$, then

$$|E_r(u, v)| \leq \frac{M^r}{r!} \left(\frac{b-a}{2} \right)^r.$$

A trivial matter is to show that the right-hand side tends to zero as $r \rightarrow \infty$. Moreover, for some constants $p, q > 1$ it does not exceed p/q^r for all r . The inequality

$$p/q^r \leq \varepsilon$$

is fulfilled whenever

$$\frac{\log p + \log \varepsilon^{-1}}{\log q} \leq r.$$

In this case, we come up with an approximate matrix-vector algorithm of the complexity $O(n \log \varepsilon^{-1})$.

The low matrix-vector complexity comes together with the low number of parameters to represent the low-rank approximant of the form (1). It is defined by the values $\phi_s(u_i)$ and $\psi_s(v_j)$. Thus, the number of defining parameters is $2rn \ll n^2$.

II. MULTILEVEL MATRICES

The low-parametric approximations as above are pervasive in applications. A very helpful idea is to construct them for certain blocks of a given matrix (especially in the case when it is nonsingular). Implicitly, this very idea can be found in the first papers on approximate fast matrix-vector algorithms in the context of potential theory ([4], [8]).

Let A be a $m \times n$ matrix, $I = \{1, 2, \dots, m\}$ and $J = \{1, \dots, n\}$. Then

$$A(\hat{I}, \hat{J}), \quad \hat{I} \subset I, \quad \hat{J} \subset J,$$

denotes a submatrix (block) on the intersection of the rows with indices from \hat{I} and columns with indices from \hat{J} . Any subdivisions on disjoint subsets

$$I_1 \cup \dots \cup I_{m_1} = I, \quad J_1 \cup \dots \cup J_{n_1} = J$$

allow us to view A as a block matrix with the blocks $A(I_k, J_l)$, $1 \leq k \leq m_1$, $1 \leq l \leq n_1$. These are the *level-1 blocks*.

Further subdivisions

$$\tilde{I}_1 \cup \dots \cup \tilde{I}_{m_2} = I, \quad \tilde{J}_1 \cup \dots \cup \tilde{J}_{n_2} = J, \quad m_2 > m_1, \quad n_2 > n_1,$$

in which the subsets are disjoint and subdivide the previous ones, consist of the *level-2 blocks* $A(\tilde{I}_k, \tilde{J}_l)$. Thus, any level-1 block is considered as a block matrix composed of the level-2

blocks. Let us proceed in the same way until we get the level- p blocks. In this case we say that A has p levels.

If A has p levels, then we make use of its *multilevel splitting* of the form

$$A = A_1 + \dots + A_p,$$

where A_k is a *sparse block matrix* consisting of the level- k blocks. The sparsity means that some prescribed blocks are zero while the others may be nonzero, call them *formally nonzero*. Assume that all the blocks of A_l of a formally nonzero block A_k , where $l > k$, must be zero.

Let r be the largest rank and τ be the sum of sizes of all nonzero blocks in a multilevel splitting of A . Then the number of defining parameters for this multilevel splitting is $O(r\tau)$. If A comes from discretization of a typical integral operator, then A can be ε -approximated by a multilevel matrix with $r = r(\varepsilon)$ and $\tau \ll n^2$. The fact is easiest to understand from the following pretty general model.

Assume that $a_{st} = f(x_s, y_t)$, where x_1, \dots, x_n and y_1, \dots, y_n are the nodes in \mathbb{R}^d . Consider different grids x_s and y_t for $n = 1, 2, \dots$ and assume that all the nodes belong to the same cuboid $S = [a_1, b_1] \times \dots \times [a_d, b_d]$. Take $k = 1, 2, \dots$ and consider uniform grids along the edges of S , each is subdivided into 2^{k+1} equal sub-intervals. There are $(2^{k+1})^d$ Cartesian products of these sub-intervals. Numerate them in some way for any k . Then, for any fixed k we obtain S as a union of smaller cuboids S_{ik} of volume h_k :

$$S = \bigcup_{i=1}^{2^{(k+1)d}} S_{ik}, \quad h_k = \frac{h_0}{2^{(k+1)d}},$$

where h_0 is the volume of S .

Assumption 1. If S_{ik} contains $\mu(S_{ik})$ nodes x_s and $\nu(S_{ik})$ nodes y_t , then

$$\max\{\mu(S_{ik}), \nu(S_{ik})\} \leq ch_k n \quad \forall i, k, n, \quad (2)$$

where $c > 0$ does not depend on i, k and n .

Assumption 2. For any $\varepsilon > 0$, the function $f(x, y)$ admits approximate separation of variables

$$\left| f(x, y) - \sum_{\alpha=1}^r \Phi_\alpha(x) \Psi_\alpha(y) \right| \leq \varepsilon \quad \forall x \in S_{ik}, \quad \forall y \in S_{jl}, \quad (3)$$

whenever $S_{ik} \cap S_{jl} = \emptyset$.

For any k , assume that $\{1, \dots, n\}$ is subdivided into subsets I_{ik} and into subsets J_{ik} so that $s \in I_{ik}$ implies $x_s \in S_{ik}$ and $t \in J_{ik}$ implies $y_t \in S_{ik}$. According to (3), if $S_{ik} \cap S_{jl} = \emptyset$ then the entries of the block $\hat{A} = A(I_{ik}, J_{jl})$ are ε -approximated by the entries of a rank- r matrix \hat{A}_r :

$$|(\hat{A})_{st} - (\hat{A}_r)_{st}| \leq \varepsilon \quad s \in I_{ik}, \quad t \in J_{jl}, \quad \text{rank} \hat{A}_r \leq r = r(\varepsilon). \quad (4)$$

Theorem. *Let $A = [f(x_s, y_t)]$ be a matrix of order n . Under Assumptions 1 and 2, its entries are ε -approximated by the entries of some multilevel matrix in which any nonzero block has upper estimate on the rank $r = r(\varepsilon)$ while the sum of sizes of all nonzero blocks is $O(n \log_2 n)$.*

Proof. Let us construct a multilevel matrix B with the following properties:

- 1) B approximates A with the entry-wise accuracy ε ;
- 2) B has p levels;
- 3) B admits a multilevel splitting $B = B_1 + \dots + B_p$ in which B_k consists of the blocks $B(I_{ik}, J_{jk})$.

Let $1 \leq k < p$. Then, if $S_{ik} \cap S_{jk} \neq \emptyset$, then set $B(I_{ik}, J_{jk}) = 0$, and the corresponding nonzero entries of B are relegated to the level- l blocks with $l > k$. If $k = 1$ then all other blocks of B_1 are formally nonzero.

Let $1 < k \leq p$. Suppose that $S_{ik} \subset S_{i'k-1}$. Then, if $S_{i'k-1} \cap S_{j'k-1} = \emptyset$ and $S_{jk} \subset S_{j'k-1}$, then $B(I_{ik}, J_{jk}) = 0$, because the corresponding nonzero entries of B are already included in some level- l blocks with $l < k$. Therefore, for any fixed index subset I_{ik} , the number of formally nonzero blocks $B(I_{ik}, J_{jk})$ does not exceed 6^d . The sum of sizes of each of them is not greater than $2ch_k n$. Since the number of index subsets I_{ik} is equal to h_0/h_k , we conclude that the sum of sizes of all formally nonzero blocks in B_k is estimated from above by $2 \cdot 6^d h_0 n$. Consequently, the sum of sizes of all nonzero blocks is $O(np)$.

It remains to choose the number of levels p so that the sizes of level- p blocks be equal to or less than r :

$$cnh_k \leq r.$$

This results in the claim that $p = O(\log_2 n)$. Hence, the rank of any nonzero block does not exceed r while the sum of sizes of all nonzero blocks amounts to $O(n \log_2 n)$, which completes the proof.

In applications related to integral equations, we enjoy the estimate

$$r(\varepsilon) \leq c \log^\gamma \varepsilon^{-1}, \quad (5)$$

where $c, \gamma > 0$ are some constants depending on the context. If so, let us say that a function $f(x, y)$ is *asymptotically separable*. This notion generalizes the one of an *asymptotically smooth* function. A typical example reads

$$f(x, y) = \frac{1}{|x - y|^\theta}, \quad \theta > 0. \quad (6)$$

Thus, if a matrix is generated by an asymptotically separable function with the constant γ in (5), then it can be ε -approximated by a multilevel matrix with low-rank blocks via $O(n \log_2 n \log^\gamma \varepsilon^{-1})$ defining parameters. The approximate matrix-vector complexity is $O(n \log_2 n \log^\gamma \varepsilon^{-1})$.

III. CROSS APPROXIMATION TECHNIQUES

The above results are germane to be regarded as *existence results*. Once we are aware that a certain structured approximations exist, we may be better off in their construction if apply purely matrix techniques. A fundamental result claims that a rank- r approximation to a matrix can be obtained from a certain cross consisting of r rows and columns of this matrix ([2], [3]). Practical algorithms for low-rank approximation prove to be some variants of the Gaussian elimination with a dynamic choice of pivot [1].

Theorem. *Given a matrix A of order n , assume that it can be approximated by a rank- r matrix B so that $\|B - A\|_2 \leq \varepsilon$, and let A be written in the block form*

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix},$$

where A_{11} is a nonsingular submatrix of order r with maximal volume (determinant in modulus) among all submatrices of order r . Then

$$\left| \left(A - \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} A_{11}^{-1} \begin{bmatrix} A_{11} & A_{12} \end{bmatrix} \right)_{ij} \right| \leq (r + 1)\varepsilon, \quad 1 \leq i, j \leq n.$$

Recently, a similar result has been proposed for tensor approximations [7]. Given a three-dimensional array (tensor)

$$\mathcal{A} = [a_{ijk}], \quad i = 1, \dots, n_1, \quad j = 1, \dots, n_2, \quad k = 1, \dots, n_3,$$

we are interested to approximate it in the so-called Tucker format

$$a_{ijk} = \sum_{i'=1}^{r_1} \sum_{j'=1}^{r_2} \sum_{k'=1}^{r_3} g_{i'j'k'} u_{ii'} v_{jj'} w_{kk'} + e_{ijk} \quad (7)$$

with the error e_{ijk} to be sufficiently small for prescribed r_1, r_2, r_3 .

If $r = r_1 = r_2 = r_3$ and $n = n_1 = n_2 = n_3$, then the number of defining parameters for the Tucker format is $r^3 + 3nr \ll n^3$. The tensor dimension can be large (for example, $n = 10^4 \div 10^6$ for some tensors coming from three-dimensional integral equations). The array itself can not be even stored in the operative memory as $\mathcal{O}(n^3)$ memory cells are needed. However, if $r \ll n$ then $\mathcal{O}(rn + r^3)$ defining parameters are by all means affordable. Useful estimates on r are developed in [5], [10], [11]. We can mention also some practical algorithms using interpolation and other function approximation techniques or additional structural properties rather than the given arrays of data ([5], [6]).

The Tucker approximation (7) is defined by the *Tucker core*

$$\mathcal{G} = [g_{i'j'k'}]$$

and three matrices

$$U = [u_{ii'}], \quad V = [v_{jj'}], \quad W = [w_{kk'}].$$

Under the knowledge that the approximation (7) exists, we may try to find it or some other Tucker approximation

$$a_{ijk} = \sum_{i'=1}^{r_1} \sum_{j'=1}^{r_2} \sum_{k'=1}^{r_3} g_{i'j'k'} u'_{ii'} v'_{jj'} w'_{kk'} + e'_{ijk} \quad (8)$$

with the same level of accuracy but arising only from a small *three-dimensional cross* consisting of rows, columns, and fibers (let use this name) of the given array. The next theorem assures that this is possible [7].

Theorem. *Given an array \mathcal{A} , assume that it can be approximated in the Tucker format (7). Then there exists a Tucker approximation (8) that is determined only from some r_1 columns, r_2 rows and r_3 fibers of \mathcal{A} , and provides the accuracy*

$$|e'_{ijk}| \leq (r_1 r_2 r_3 + 2r_1 r_2 + 2r_1 + 1)\varepsilon.$$

IV. THE USE OF TENSOR APPROXIMATIONS

Suppose we need to solve the equation

$$\int_D \frac{1}{|x-y|} \phi(y) dy = f(x), \quad x, y \in D = [0 : 1]^3. \quad (9)$$

Then we can subdivide the cube D into subcubes

$$[a_{i-1}, a_i] \times [a_{j-1}, a_j] \times [a_{k-1}, a_k], \quad 0 = a_0 < a_1 < \dots < a_n = 1,$$

approximate $\phi(y)$ by a constant u_{ijk} on every subcube, and by collocation arrive at a system of linear algebraic equations

$$Au = f,$$

where the vectors are associated with discrete functions u_{ijk}, f_{ijk} on the grid with n^3 nodes. The coefficient matrix A contains n^6 nonzero entries of which none can be neglected. It possesses no special structure if the grids are nonuniform. If $n = 64$ then A would require 512Gb to be stored. If $n = 256$, then it amounts to 2Pb (1Pb = 2^{50} byte). A big problem is already here with the storage for A .

To overcome the difficulty, the only idea is to find a sufficiently close problem with a prominent low-parametric structure defined by reasonably few parameters, and then construct some gain-of-the-structure methods. For the equation (9) we can use the following *tensor format*:

$$A \approx \tilde{A}_r = U_1 \otimes V_1 \otimes W_1 + \dots + U_r \otimes V_r \otimes W_r.$$

What it may give for practice can be seen from the following table [9]:

grid size (n)	16	32	64	128	256	512
full storage for A	128Mb	8Gb	512Gb	32Tb	2Pb	128Pb
tensor format storage	74Kb	320Kb	1.1Mb	6Mb	25Mb	114Mb
approximation time	0.6sec	1.5sec	8.4 sec	54sec	5.5min	30min

TABLE 1: Tensor approximation of A with accuracy $\varepsilon = 10^{-5}$.

Thus, instead of 2Pb it is sufficient to have just 25Mb. This is available on most home computers nowadays! The tensor approximation was constructed from cleverly selected rows, columns and fibers of the 3D array representing the coefficient matrix [7].

May 30, 2007

ACKNOWLEDGMENT

This work was supported by the Russian Foundation for Basic Research (grants 05-1-00721, 06-01-08052) and the Priority Research Programme of the Department of Mathematical Sciences

of the Russian Academy of Sciences. The paper was completed during the author's visit to the School of Computational and Applied Mathematics at the University of Witwatersrand within the framework agreement between the government of South Africa and the Russian Federation.

REFERENCES

- [1] J. M. Ford, E. E. Tyrtyshnikov, "Combining Kronecker product approximation with discrete wavelet transforms to solve dense, function-related systems", *SIAM J. Sci. Comp.*, vol. 25, no. 3, pp. 961–981, 2003.
- [2] S. A. Goreinov, E. E. Tyrtyshnikov, N. L. Zamarashkin, "A theory of pseudo-skeleton approximations", *Linear Algebra Appl.*, 261, pp. 1–21, 1997.
- [3] S. A. Goreinov, E. E. Tyrtyshnikov, "The maximal-volume concept in approximation by low-rank matrices", *Contemporary Mathematics*, vol. 280, pp. 47–51, 2001.
- [4] W. Hackbusch, Z. P. Nowak, On the fast matrix multiplication in the boundary elements method by panel clustering, *Numer. Math.* **54** (4) (1989) 463–491.
- [5] W. Hackbusch, B. N. Khoromskij, E. E. Tyrtyshnikov, "Hierarchical Kronecker tensor-product approximations", *J. Numer. Math.*, vol. 13, pp. 119–156, 2005.
- [6] V. Olshevsky, I. V. Oseledets, E. E. Tyrtyshnikov, "Tensor properties of multilevel Toeplitz and related matrices", *Linear Algebra Appl.*, 412, pp. 1–21, 2006.
- [7] I. V. Oseledets, D. V. Savostianov, E. E. Tyrtyshnikov, "Tucker dimensionality reduction of three-dimensional arrays in linear time", *SIMAX*, to appear, 2007.
- [8] V. Rokhlin, Rapid solution of integral equations of classical potential theory, *J. Comput. Physics* 60: 187–207, 1985.
- [9] D. V. Savostianov, "Fast polylinear approximation of matrices and integral equations", *Ph.D. Thesis, Institute of Numerical Mathematics of the Russian Academy of Sciences*, Moscow, 2006.
- [10] E. E. Tyrtyshnikov, "Tensor approximations of matrices generated by asymptotically smooth functions", *Sbornik: Mathematics*, vol. 194, no. 5-6, pp. 941–954, 2003.
- [11] E. E. Tyrtyshnikov, "Kronecker-product approximations for some function-related matrices", *Linear Algebra Appl.*, 379, pp. 423–437, 2004.



Eugene Tyrtyshnikov was born on June 2, 1955, in Moscow; graduated in 1977 from the Lomonosov Moscow State University (Faculty of Computational Mathematics and Cybernetics); Ph. D. (1980); the Russian Doctor Degree (habilitation) in 1990; Professor at the Lomonosov Moscow State University and Moscow Institute of Physics and Technology; the member of editorial boards of *Calcolo*, *Linear Algebra and Its Applications*, *Journal of Numerical Mathematics*, *Russian Journal of Numerical Mathematics and Modelling*, *Journal of Computational Mathematics and Mathematical Physics*; elected the Correspondent Member of the Russian Academy of Sciences in 2006. The Deputy Director of Institute of Numerical Mathematics of the Russian Academy of Sciences (since 2001).