# KOVARIK'S FUNCTION ORTHOGONALIZATION ALGORITHM WITH APPROXIMATE INVERSION*

CONSTANTIN POPA[†]

**Abstract.** Z. Kovarik proposed in 1970 a method for approximate orthogonalization of a finite set of linearly independent vectors from a Hilbert space. This method uses at each iteration a symmetric and positive definite matrix inversion. In this paper we describe an algorithm in which the above matrix inversion step is replaced by an arbitrary odd degree polynomial matrix expression. We prove that this new algorithm converges to the same orthonormal set of vectors as the original Kovarik's method. Some numerical experiments presented in the last section of the paper show us that, even for small degree polynomial expressions the convergence properties of the new algorithm are comparable with those of the original one.

**MSC 2000.** 65F10, 65F20.

**Keywords.** Approximate orthogonalization of functions, Kovarik's algorithm, Gram matrix, approximate inverse.

## 1. A MATRIX APPROACH OF KOVARIK'S ALGORITHM

Let $H$ be a (real) Hilbert space with the scalar product and the associated norm denoted by $\langle \cdot, \cdot \rangle$ and $\| \cdot \|$, respectively. In some space $\mathbb{R}^q$ the vectors will be considered as column vectors and the euclidean scalar product and norm will be denoted by $\langle \cdot, \cdot \rangle_2$ and $\| \cdot \|_2$, respectively. For a $q \times q$ matrix $A$ we shall denote by $A^t, (A)_i, (A)_{ij}, \sigma(A), \rho(A)$ the transpose, $i$-th row, $(i, j)$-th element, spectrum and spectral radius, respectively, and by $\| A \|_2, \| A \|_\infty$ the matrix norms defined by (see e.g. [1])

$$(1) \qquad \| A \|_2 = \sqrt{\rho(A^t A)}, \ \| A \|_\infty = \max_{1 \le i \le q} \sum_{j=1}^{q} |(A)_{ij}|.$$

For a linearly independent system of vectors $\Phi = \{\phi_1, \dots, \phi_n\} \subset H$ the Gram matrix $G(\Phi)$, defined by

$$(2) \qquad (G(\Phi))_{ij} = \langle \phi_j, \phi_i \rangle, \ i, j = 1, \dots, n$$

is symmetric and positive definite (SPD, for short), thus its square root $G(\Phi)^{\frac{1}{2}}$ exists and is also SPD. We define the number $\text{trace}(G(\Phi))$ by

$$(3) \qquad \text{trace}(G(\Phi)) = \sum_{i=1}^{n}(G(\Phi))_{ii} = \sum_{i=1}^{n}\langle\phi_i,\phi_i\rangle = \sum_{i=1}^{n}\parallel\phi_i\parallel^2 .$$

If $\sigma(G(\Phi)) = \{\sigma_1,\ldots,\sigma_n\}$, then $\text{trace}(G(\Phi)) = \sigma_1 + \cdots + \sigma_n$. For the above system $\Phi$ the author considered in [1] the norm $|||\cdot|||$ defined by

$$(4) \qquad |||\Phi|||^2 = \frac{1}{n}\text{trace}(G(\Phi)) = \frac{1}{n}\sum_{i=1}^{n}\parallel\phi_i\parallel^2$$

and proved the inequality

$$(5) \qquad |||\Phi|||^2 \leq \rho(G(\Phi)) = \parallel G(\Phi)\parallel_2,$$

where $\parallel G(\Phi)\parallel_2$ is the spectral norm of the matrix $G(\Phi)$. Moreover, for an $n \times n$ matrix $C = (c_{ij})_{i,j}$, he defined the product $\Psi = \Phi \cdot C$ by

$$(6) \qquad \Psi = \{\psi_1,\ldots,\psi_n\}, \ \psi_i = \sum_{j=1}^{n}c_{ij}\phi_j, \ i = 1,\ldots,n.$$

REMARK 1. If we consider the systems $\Psi = (\phi_1,\ldots,\phi_n), \Phi = (\psi_1,\ldots,\psi_n)$ as (ordered) vectors in $H^n = H \times H \times \cdots \times H$, we can write (6) in the following "matrix by vector multiplication form", very useful in computations (see Section 2)

$$(7) \qquad \Psi = \begin{bmatrix} \psi_1 \\ \psi_2 \\ \ldots \\ \psi_n \end{bmatrix} = C^{\text{t}}\begin{bmatrix} \phi_1 \\ \phi_2 \\ \ldots \\ \phi_n \end{bmatrix}.$$

Let $\Phi^\infty = \{\phi_1^\infty,\ldots,\phi_n^\infty\}$ be given by (see (6)–(7))

$$(8) \qquad \Phi^\infty = \Phi \cdot (G(\Phi))^{-\frac{1}{2}}.$$

According to [2], $\Phi^\infty$ is an orthonormal system, i.e.

$$(9) \qquad \langle\phi_j^\infty,\phi_i^\infty\rangle = \begin{cases} 1, & j = i \\ 0, & j \neq i \end{cases}$$

and the following result holds.

LEMMA 2. *Let $C = (c_{ij})_{i,j}$ be an arbitrary $n \times n$ matrix. Then*

$$(10) \qquad |||\Phi^\infty \cdot C||| \ \leq \ \parallel C\parallel_2 .$$

In [2] Z. Kovarik proposed his approximate orthogonalization algorithm. It will be briefly presented in what follows, together with the corresponding convergence results.

ALGORITHM 3. Let $\Phi^0 = \Phi$; for $k = 0,1,\ldots$ do

$$(11) \quad G_k = G(\Phi^k); K_k = (I - G_k)(I + G_k)^{-1}; S_k = I + K_k; \Phi^{k+1} = \Phi^k \cdot S_k.$$

THEOREM 4. *For any linearly independent system $\Phi$, the sequence $(\Phi^k)_{k \geq 0}$ generated by* (11) *converges to $\Phi^\infty$ (in $H$). Moreover, the following estimate holds*

$$(12) \qquad |||\Phi^\infty - \Phi^k||| \ \leq \ \| K_0 \|_2^{2^k}, \ \forall k \geq 1.$$

REMARK 5. The above convergence of $\Phi^k = \{\phi_1^{(k)}, \ldots, \phi_n^{(k)}\}$ to $\Phi^\infty$ is equivalent with (see (4))

$$(13) \qquad \lim_{k \to \infty} \| \phi_i^{(k)} - \phi_i^\infty \| = 0, \ i = 1, \ldots, n.$$

Now, if we denote by $\lambda_i^{(0)}$ the eigenvalues of $G_0$, i.e.

$$(14) \qquad \sigma(G_0) = \{\lambda_1^{(0)}, \ldots, \lambda_n^{(0)}\} \subset (0, \infty),$$

then, from the second equality in (11), we obtain

$$\sigma(K_0) = \left\{ \frac{1-\lambda_1^{(0)}}{1+\lambda_1^{(0)}}, \ldots, \frac{1-\lambda_n^{(0)}}{1+\lambda_n^{(0)}} \right\} \subset (-1, 1),$$

thus

$$(15) \qquad \| K_0 \|_2 < \ 1,$$

which tells us that the convergence in (12) is at least quadratic. Moreover, from the equalities (see [2])

$$(16) \qquad \Phi^\infty = \Phi^k \cdot G_k^{-\frac{1}{2}}, \ \forall k \geq 0,$$

we obtain

$$(17) \qquad \Phi^\infty - \Phi^k = \Phi^\infty \cdot (I - G_k^{\frac{1}{2}}), \ \forall k \geq 0,$$

which together with (9) tell us that the limit in (13) is equivalent with

$$(18) \qquad \lim_{k \to \infty} \| I - G_k^{\frac{1}{2}} \|_2 = 0.$$

LEMMA 6. *The Gram matrices $G_k = G(\Phi^k)$, $k \geq 0$, can be recursively generated by the following formulas*

$$(19) \qquad G_{k+1} = S_k G_k S_k, \ k \geq 0,$$

*with $K_k$ and $S_k$ computed as in* (11).

*Proof.* From (11) and (6) we get (also using the symmetry of $G_k$)

$$(G_{k+1})_{ij} = \langle \phi_j^{k+1}, \phi_i^{k+1} \rangle = \sum_{q=1}^{n} (S_k)_{iq} (\sum_{p=1}^{n} (S_k)_{jp} \langle \phi_p^{(k)}, \phi_q^{(k)} \rangle) =$$

$$= \sum_{q=1}^{n} (S_k)_{iq} (\sum_{p=1}^{n} (S_k)_{jp} (G_k)_{qp}) = \langle (S_k)_i, G_k(S_k)_j \rangle_2 = (S_k G_k S_k)_{ij},$$

which completes the proof.                                                    $\square$

## 2. THE MODIFIED KOVARIK ALGORITHM

As already observed by Z. Kovarik in [2], for applying the algorithm (3) we need to compute at each iteration the inverse $(I + G_k)^{-1}$ of the (SPD) matrix $I + G_k$. We shall avoid this difficulty by replacing the matrix inverse with a polynomial expression with respect to $G_k$. For this we shall first observe that after the scaling of the system $\Phi = \Phi^0 = \{\phi_1^{(0)}, \ldots, \phi_n^{(0)}\}$

$$(20) \qquad \psi_i^{(0)} = \frac{1}{\sqrt{\| G_0 \|_\infty + 1}} \phi_i^{(0)}, \; i = 1, \ldots, n.$$

If $\Psi^0 = \{\psi_1^{(0)}, \ldots, \psi_n^{(0)}\}$ we get

$$(21) \qquad \| G(\Psi^0) \|_2 = \| \frac{1}{\| G_0 \|_\infty + 1} G(\Phi^0) \|_2 = \frac{\| G_0 \|_2}{\| G_0 \|_\infty + 1} < 1.$$

Then, for a given sequence of integers $q_k \geq 1$, $k \geq 0$, we shall approximate the inverse $(I + G_k)^{-1}$ in (11) by the truncated Neumann series $S(q_k; G_k)$ defined by

$$(22) \qquad S(q_k; G_k) = \sum_{i=0}^{q_k} (-G_k)^i.$$

Then, the modified Kovarik algorithm is the following.

ALGORITHM 7. Starting with $\Psi^0 = \{\psi_1^{(0)}, \ldots, \psi_n^{(0)}\}$ from (20), for $k \geq 0$, do

$$(23) \qquad \Gamma_k = G(\Psi^k); K_k = (I - \Gamma_k) S(q_k; \Gamma_k); S_k = I + K_k; \Psi^{k+1} = \Psi^k \cdot S_k.$$

The next (main) result of the paper shows the convergence of the algorithm (7).

THEOREM 8. *If $\Psi^0 = \{\psi_1^{(0)}, \ldots, \psi_n^{(0)}\}$ is defined as in (20) and all the integers $q_k, k \geq 0$, are odd, then the sequence $(\Psi^k)_{k \geq 0}$ generated in (23) converges to $\Phi^\infty$ from (8) in the sense (13), i.e.*

$$(24) \qquad \lim_{k \to \infty} \| \psi_i^{(k)} - \phi_i^\infty \| = 0, \; \forall i = 1, \ldots, n.$$

In order to prove it we need some auxiliary results which will be presented below.

LEMMA 9. *Let $\Gamma_0 = G(\Psi^0)$ and $Q$ be an $n \times n$ orthonormal matrix such that*

$$(25) \qquad \Gamma_0 = Q^t \operatorname{diag}(\gamma_1^{(0)}, \ldots, \gamma_n^{(0)}) \; Q$$

*with (see (21))*

$$(26) \qquad \sigma(\Gamma_0) = \{\gamma_1^{(0)}, \ldots, \gamma_n^{(0)}\} \subset (0, 1).$$

*Then, if $q_k$ are all odd numbers we have*

$$(27) \qquad \Gamma_k = Q^t \operatorname{diag}(\gamma_1^{(k)}, \ldots, \gamma_n^{(k)}) \; Q,$$

with $\gamma_i^{(k)}, i = 1, \ldots, n$, recursively defined by

$$(28) \qquad \gamma_i^{(k+1)} = \left( \frac{2 + (\gamma_i^{(k)})^{q_k+1}(\gamma_i^{(k)} - 1)}{1 + \gamma_i^{(k)}} \right)^2 \cdot \gamma_i^{(k)} \in (0,1), \forall k \geq 1.$$

*Proof.* According to Lemma 6 and (23) we obtain

$$(29) \qquad\qquad\qquad \Gamma_{k+1} = S_k \Gamma_k S_k,$$

with

$$(30) \qquad\qquad\qquad S_k = I + (I - \Gamma_k)S(q_k; \Gamma_k).$$

Now we shall use an induction argument. For this, let $\Gamma_k$ be as in (27). Using the orthonormality of $Q$, (22) and (29)–(30) we obtain

$$(31) \qquad\qquad\qquad S_k = Q^{\mathrm{t}} \operatorname{diag}(d_1^{(k)}, \ldots, d_n^{(k)}) \, Q$$

with

$$(32) \qquad d_i^{(k)} = 1 + (1 - \gamma_i^{(k)}) \sum_{j=0}^{q_k} (-1)^j (\gamma_i^{(k)})^j > 0, \; i = 1, \ldots, n.$$

Then, because $\gamma_i^{(k)} \in (0,1)$ we obtain (by also taking into account that $q_k$ is odd)

$$(33) \qquad\qquad d_i^{(k)} = 1 + \frac{(1 - \gamma_i^{(k)})(1 - (\gamma_i^{(k)})^{q_k+1})}{1 + \gamma_i^{(k)}},$$

which gives us, together with (27) and (29),

$$\Gamma_{k+1} = Q^{\mathrm{t}} \operatorname{diag}(\gamma_1^{k+1}, \ldots, \gamma_n^{(k+1)}) \, Q,$$

with $\gamma_i^{(k+1)}$ as in (28). So, it rests us to prove that

$$(34) \qquad\qquad \gamma_i^{(k+1)} \in (0,1), \; \forall i = 1, \ldots, n.$$

For this, we consider the function $f : [0, \infty) \longrightarrow \mathbb{R}$

$$(35) \qquad\qquad f(x) = \left( \frac{2 + x^{q_k+1}(x - 1)}{1 + x} \right)^2 \cdot x$$

and observe that, $\forall x \in (0,1)$,

$$(36) \qquad f(x) - 1 = \frac{-(1-x)^2 - x(1-x)x^{q_k+1}[4 - (1-x)x^{q_k+1}]}{(1+x)^2} < 0.$$

Then, from (28) and (36), it results (34) and the proof is complete. $\qquad\square$

LEMMA 10. *For any $i \in \{1, \ldots, n\}$, if $q_k$ are all odd numbers then the sequence $(\gamma_i^{(k)})_{k \geq 0}$, recursively defined in (28), is strictly increasing and*

$$(37) \qquad\qquad\qquad \lim_{k \to \infty} \gamma_i^{(k)} = 1.$$

*Proof.* Using again $f(x)$ from (35) we get, for $x \in (0,1)$,

$$(38) \qquad f(x) - x = x \cdot \frac{3 + x - (1-x)x^{q_k+1}}{1+x} \cdot \frac{(1-x)(1-x^{q_k+1})}{1+x} > 0.$$

Then, by a recursive argument we obtain that the sequence $(\gamma_i^{(k)})_{k \geq 0}$ is strictly increasing and bounded,

$$(39) \qquad \gamma_i^{(k)} \in (\gamma_i^{(0)}, 1), \ \forall k \geq 0.$$

If $\gamma_i \in [0,1]$ will denote its limit, i.e.

$$(40) \qquad \gamma_i = \lim_{k \to \infty} \gamma_i^{(k)}, \ i = 1, \ldots, n,$$

and we shall prove that $\gamma_i = 1$. If this is not true true, i.e.

$$(41) \qquad 0 < \gamma_i < 1$$

($\gamma_i \neq 0$ because $\gamma_i^{(0)} > 0$ and $(\gamma_i^{(k)})_{k \geq 0}$ is strictly increasing), then from (28) and (38) we obtain

$$(42) \qquad \gamma_i^{(k+1)} - \gamma_i^{(k)} > \frac{\gamma_i^{(k)}(1 - \gamma_i^{(k)})^2}{1 + \gamma_i^{(k)}}.$$

The variation of the real function $g(x) = \frac{x(1-x)^2}{1+x}$, $x \in (0,1)$ gives us the relation

$$(43) \qquad g(x) \geq \min\{g(\gamma_i^{(0)}), g(\gamma_i)\} > 0, \ \forall x \in [\gamma_i^{(0)}, \gamma_i].$$

Thus, if we define $\epsilon_0 > 0$ by

$$(44) \qquad \epsilon_0 = \min\{g(\gamma_i^{(0)}), g(\gamma_i)\},$$

from (42) we obtain

$$\gamma_i^{(k+1)} > \gamma_i^{(k)} + \epsilon_0, \ \forall k \geq 0,$$

i.e.

$$(45) \qquad \gamma_i^{(k+1)} > \gamma_i^{(0)} + k \cdot \epsilon_0, \ \forall k \geq 0.$$

The above inequality tells us that it exists an integer $k_0 \geq 1$ such that

$$\gamma_i^{(k+1)} > \gamma_i^{(k_0+1)} > \gamma_i,$$

which contradicts (39)–(41). It results that (37) is true and the proof is complete. $\qquad \square$

*Proof.* (of Theorem 8.) From (27) and (37) we obtain that $\lim_{k \to \infty} \Gamma_k = I$, thus

$$(46) \qquad \lim_{k \to \infty} \| I - \Gamma_k^{\frac{1}{2}} \|_2 = 0.$$

Then, according to (16)–(18) we shall obtain (24) if we can prove for the systems $\Psi^k$ equalities similar with (16), i.e.

$$(47) \qquad \Phi^\infty = \Psi^k \cdot \Gamma_k^{-\frac{1}{2}}, \ \forall k \geq 0.$$

For this we shall use the mathematical induction. First, for $k = 0$ from (20), (21), (8) we get

$$(48) \quad \Psi^0 \cdot \Gamma_0^{-\frac{1}{2}} = \left( \frac{1}{\sqrt{\| G_0 \|_\infty + 1}} \, \Phi^0 \right) \cdot (\sqrt{\| G_0 \|_\infty + 1} \, G_0^{-\frac{1}{2}}) = \Psi^0 \cdot G_0^{-\frac{1}{2}} = \Phi^\infty.$$

Let now $k \geq 0$ be an integer such that (47) holds for it. Then, using (23), (29) and (7) we successively obtain

$$\Psi^{k+1} \cdot \Gamma_{k+1}^{-\frac{1}{2}} = (\Psi^k \cdot S_k)(S_k \Gamma_k S_k)^{-\frac{1}{2}} = S_k^t \Psi^k (S_k \Gamma_k S_k)^{-\frac{1}{2}} =$$

$$(49) \quad = [S_k^t (S_k \Gamma_k S_k)^{-\frac{1}{2}}]\Psi^k = ((S_k \Gamma_k S_k)^{-\frac{1}{2}} S_k)^t \Psi^k = \Psi^k \cdot [(S_k \Gamma_k S_k)^{-\frac{1}{2}} S_k],$$

which together with (27), (31) and (32) gives us the equality

$$(50) \qquad (S_k \Gamma_k S_k)^{-\frac{1}{2}} S_k = (Q^t D_k T_k D_k Q)^{-\frac{1}{2}} Q^t D_k Q,$$

where we used the notations

$$(51) \qquad D_k = \operatorname{diag}(d_1^{(k)}, \ldots, d_n^{(k)}), \ T_k = \operatorname{diag}(\gamma_1^{(k)}, \ldots, \gamma_n^{(k)}).$$

But, because $D_k$ and $T_k$ are diagonal matrices they commute, thus ($Q$ being orthonormal)

$$Q^t D_k T_k D_k Q = (Q^t (D_k T_k D_k)^{\frac{1}{2}} Q)^2$$

and from (50), (27) and (51) we obtain

$$(52) \qquad (S_k \Gamma_k S_k)^{-\frac{1}{2}} S_k = Q^t D_k^{-\frac{1}{2}} T_k^{-\frac{1}{2}} D_k^{-\frac{1}{2}} Q Q^t D_k Q = Q^t T_k^{-\frac{1}{2}} Q = \Gamma_k^{-\frac{1}{2}}.$$

Then, the relations (51), (49) and the induction hypothesis give us $\Phi^\infty = \Psi^{k+1} \cdot \Gamma_{k+1}^{-\frac{1}{2}}$ and the proof is complete. $\qquad \square$

REMARK 11. In the other cases for the integers $q_k$, i.e. $q_k = $ constant (even) or $q_k = $ arbitrary, the modified algorithm (22)–(23) does not always converge.

### 3. NUMERICAL EXPERIMENTS

We considered in our experiments a regular discretization of $(0, 1)$: $N \geq 2$, $h = 1/N$, $n = N - 1$, $x_i = ih$, $i = 1, \ldots, n$ and the piecewise linear (one dimensional) finite element basis $\Phi = \Phi^0 = \{\phi_1^{(0)}, \ldots, \phi_n^{(0)}\}$, $\phi_i^{(0)} : [0, 1] \longrightarrow \mathbb{R}$ defined by (see e.g. [3])

$$\phi_i^{(0)}(x) = \begin{cases} \frac{x - x_{i-1}}{h}, & x \in (x_{i-1}, x_i] \\ \frac{x_{i+1} - x}{h}, & x \in (x_i, x_{i+1}) \\ 0, & \text{else.} \end{cases}$$

As the Hilbert space $H$ we considered $H_0^1((0,1))$, with the scalar product

$$\langle \phi, \psi \rangle = \int_0^1 \frac{\mathrm{d}\phi}{\mathrm{d}x} \frac{\mathrm{d}\psi}{\mathrm{d}x}.$$

In this context, the Gram matrix $G_0 = (\langle \phi_j^{(0)}, \phi_i^{(0)} \rangle)_{i,j=1,\ldots,n}$ is given by

$$G_0 = \begin{bmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ \ldots & \ldots & \ldots & \ldots & \ldots & \ldots \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 \end{bmatrix}.$$

We first applied the original Kovarik algorithm (11) (using (19)), for different values of $N \geq 2$ and the stopping test

(53) $$\| G_{k+1} - G_k \|_\infty \leq 10^{-3}.$$

The numbers of iterations for obtaining (53) are described in Table 1.

Then, for $N = 128$ we applied the modified Kovarik algorithm (22)–(23) (with (29)) for $q_k = $ constant (odd), $\forall k \geq 0$ and the same stopping rule (53). The numbers of iterations for these tests are presented in Table 2.

The last tests were made as those from above, but only for the first three (odd) values of $q_k$ and different values of $N \geq 2$. The numbers of iterations are described in Table 3.

FINAL REMARKS AND COMMENTS

(1) We observe that for very small values of $q_k$ (which means less computational effort per iteration in (22)–(23)), we got good enough results (Tables 2 and 3) by comparing them with those for the algorithm (11) (Table 1).

(2) All the tests indicate a "mesh-independent" behaviour for both Kovarik, original and modified algorithms.

(3) All the numerical experiments were performed with the numerical linear algebra software "OCTAVE", freely available on the Internet.

Table 1. Algorithm 3

| $N = 16$ | $N = 32$ | $N = 64$ | $N = 128$ | $N = 256$ |
|---|---|---|---|---|
| 7 | 8 | 9 | 10 | 11 |

Table 2. Algorithm 7 with $N = 128$

| $q_k = 1$ | $q_k = 3$ | $q_k = 5$ | $q_k = 7$ | $q_k = 9$ | $q_k = 11$ |
|---|---|---|---|---|---|
| 28 | 22 | 19 | 17 | 16 | 15 |

Table 3.  Algorithm 7

|           | $N = 16$ | $N = 32$ | $N = 64$ | $N = 128$ | $N = 256$ |
|-----------|----------|----------|----------|-----------|-----------|
| $q_k = 1$ | 16       | 26       | 27       | 28        | 29        |
| $q_k = 3$ | 19       | 20       | 21       | 22        | 23        |
| $q_k = 5$ | 16       | 17       | 18       | 19        | 20        |

## REFERENCES

[1] BJÖRCK, A., *Numerical methods for least squares problems*, SIAM, Philadelphia, 1996.

[2] KOVARIK, Z., *Some iterative methods for improving orthogonality*, SIAM J. Numer. Anal **7(3)**, pp. 386–389, 1970.

[3] TROTTENBERG, U., OOSTERLEE, C. and SCHÜLLER, A., *Multigrid*, Academic Press, New York, 2001.