

Convergence Acceleration of Alternating Least Squares with a Matrix Polynomial Predictive Model for PARAFAC Decomposition of a Tensor

Ming Shi, JianQiu Zhang*, Bo Hu, Bin Wang, and Qiyong Lu
 The Research Center of Smart Networks and Systems
 School of Information Science and Technology, Fudan University
 No. 220, Handan Rd., Shanghai, 200433, P.R. China
 Email: jqzhang01@fudan.edu.cn

Abstract—In this paper, a matrix polynomial whose coefficients are matrices is first defined. Its predictive model, called as the Matrix Polynomial Predictive Model (MPPM), is then derived. When the loading matrices of a decomposed tensor in the Alternating Least Squares (ALS) are replaced by the predicted ones of the MPPM, a new ALS algorithm with the MPPM (ALS-MPPM) is proposed. Analyses show that the convergent rate of the proposed ALS-MPPM is closely related to the degree of the matrix polynomial. Namely, when an accelerative convergence rate is expected, the polynomial with a high degree is preferred. Although a high degree means a high possibility of prediction failure, a simple solution can be used to handle such failure. Moreover, the relationship between our ALS-MPPM and the existing ALS-based algorithms is also analyzed. The results of numerical simulations show that the proposed ALS-MPPM outperforms the reported ALS-based algorithms in the literature while the analytical results are verified.

I. INTRODUCTION

It is well known that the PARAllel FACtor decomposition (PARAFAC) [1], also called as the CANonical DECOMPosition (CANDECOMP) [2] or Canonical Polyadic Decomposition (CPD) [3] is one of the most commonly used tools for processing multidimensional data. The PARAFAC decomposition factorizes an N-way array (i.e., a tensor) into the sum of vector outer products. It has theoretically been shown that it is deterministic and essentially unique (that is the uniqueness with scaling and permutation ambiguity) [4] [5] under mild conditions, which makes it powerful for numerous applications.

As the most commonly used algorithm for decomposing a tensor, the *Alternating Least Squares* (ALS) [1] [2] can be intuitively understood and straightforwardly implemented. However, ALS usually needs to take multiple iterations to make a PARAFAC decomposition converge and sometimes it could encounter some numerical issues keeping it from further convergence [6] [7]. To improve the convergence of the ALS algorithm, the *Line Search* (LS) [8] [9] and the *Enhanced Line Search* (ELS) [10] [11] techniques have been reported. The LS and ELS accelerate the convergence of the

ALS by predicting the loading matrices of a tensor via a direction from the estimates of previous iterations. Although they can speed up the convergence of the ALS, they have their own drawbacks. For example, the relaxation factor in the LS is usually chosen empirically, and therefore lacks solid mathematical foundations. Although the relaxation factor in the ELS can optimally be determined by solving a high degree polynomial equation, there is no general solution available for the equation. It means that it is very difficult and time-consuming for one to finding such an optimal factor.

The *Polynomial Predictive Model* (PPM) is originally reported in [12] for tracking maneuvering targets. It assumes that, when a signal is described as a polynomial, one can predict the signal by the previous sample values of the signal with an FIR filter, whose optimal coefficients can be calculated offline. In this paper, a matrix polynomial whose coefficients are matrices is first defined. Its predictive model, called as the *Matrix Polynomial Predictive Model* (MPPM), is then derived. When the loading matrices of a decomposed tensor in the Alternating Least Squares (ALS) are replaced by the predicted ones of the MPPM, a new ALS algorithm with the MPPM (ALS-MPPM) for the PARAFAC decomposition of a tensor is obtained. The relationship between the degree of the MPPM and the rate of convergence is analyzed. It is shown that, when an accelerative convergence rate is expected, the polynomial with a high degree is preferred. Although a high degree means a high possibility of prediction failure, a simple solution to handle such failure is given. The analyses also show that there is a close relationship between our ALS-MPPM and the existing ALS-based algorithms. For example, the original ALS can be taken as a special case of our ALS-MPPM and the LS and ELS is the version of our ALS-MPPM with a polynomial of degree one. The numerical simulations show that the proposed ALS-MPPM outperforms the reported ALS-based algorithms in the literature while the analytical results are verified.

The rest of the paper is organized as follows. The PARAFAC decomposition of a tensor and ALS will be briefly reviewed in Section II. In Section III, the MPPM for accelerating the convergence of a ALS algorithm and the optimal coefficients

This work was supported by the National Natural Science Foundation of China by Grant 61571131.

for our MPPMs with a low degree polynomial are given. In Section IV, our MPPM is incorporated with the ALS, and the relationships between our ALS-MPPM and the existing algorithms are studied. Numerical simulation results are presented in Section V. Finally, conclusions are drawn in Section VI.

Notations: The scalars, vectors, matrices, and tensors are denoted by lowercase letters x , lowercase boldface letters \mathbf{x} , boldface capitals \mathbf{A} , and calligraphic letters \mathcal{X} , respectively. The superscript T , H , and \dagger respectively stand for the transpose, conjugate transpose and pseudo-inverse. The symbols \odot and \circ denote the Khatri-Rao and vector outer product. The norm of a matrix is represented by $\|\cdot\|$.

II. PROBLEM FORMULATION

For convenience of description and without loss of generalization, a third-order tensor will be taken as an example to give our approach since it is plain to extend our results into a higher order one.

Given a third order tensor $\mathcal{X} \in \mathbb{C}^{I \times J \times K}$, the objective of its PARAFAC decomposition is to find the best approximation $\hat{\mathcal{X}}$ as [4] [5]

$$\min_{\hat{\mathcal{X}}} \|\mathcal{X} - \hat{\mathcal{X}}\| \quad \text{with} \quad \hat{\mathcal{X}} = \sum_{r=1}^R \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r, \quad (1)$$

where the minimal positive integer R that holds the equation denotes the rank of \mathcal{X} , \mathbf{a}_r , \mathbf{b}_r , and \mathbf{c}_r represent the r -th column of the loading matrices $\mathbf{A} \in \mathbb{C}^{I \times R}$, $\mathbf{B} \in \mathbb{C}^{J \times R}$, and $\mathbf{C} \in \mathbb{C}^{K \times R}$, respectively. In the matricized form, the PARAFAC decomposition can also be written as [4] [5]

$$\begin{cases} \mathbf{X}_{(1)} = \mathbf{A}(\mathbf{C} \odot \mathbf{B})^T \\ \mathbf{X}_{(2)} = \mathbf{B}(\mathbf{C} \odot \mathbf{A})^T \\ \mathbf{X}_{(3)} = \mathbf{C}(\mathbf{B} \odot \mathbf{A})^T \end{cases}, \quad (2)$$

where $\mathbf{X}_{(i)}$, $i = 1, 2, 3$ are the mode- n unfolding of the tensor \mathcal{X} .

The most widely used PARAFAC algorithm is the Alternating Least Squares, which alternately update only one loading matrix while keeping the others fixed. Therefore the PARAFAC decomposition of a tensor is taken as a linear least-squares problem $\min_{\hat{\mathbf{A}}} \|\mathbf{X}_{(1)} - \hat{\mathbf{A}}(\mathbf{C} \odot \mathbf{B})^T\|_F$, whose optimal estimate can be given as [4]

$$\hat{\mathbf{A}} = \mathbf{X}_{(1)} \left[(\mathbf{C} \odot \mathbf{B})^T \right]^\dagger. \quad (3)$$

The optimal estimates of \mathbf{B} and \mathbf{C} can be similarly derived. However, the ALS will usually take numerous iterations to converge when the tensor is mixed with an additive noise. Whats even worse, it could encounter some numerical issues keeping it from further convergence [6] [7]. In the next section, a predictive model will be introduced to accelerate the convergence of the ALS.

III. MATRIX POLYNOMIAL PREDICTIVE MODEL

A. PPM in Matrix Form

Given an unknown signal $x_k, k = 1, \dots, K$, the PPM in [12] assumes that the signal is a polynomial of degree L , then it can be written as

$$x_k = \sum_{l=0}^L p(l)k^l. \quad (4)$$

It is shown in [13] that the value of x_k can be predicted by the M previous values x_{k-M}, \dots, x_{k-1} via an FIR filter as

$$x_k^- = \sum_{m=1}^M h_L(m)x_{k-m}. \quad (5)$$

Similarly, let us define a matrix polynomial signal as follows

$$\mathbf{A}_k = \sum_{l=0}^L \mathbf{P}(l)k^l, \quad (6)$$

where $\mathbf{A}_k, k = 1, \dots, K$ are a series of matrix-valued signals, and $\mathbf{P}(l), l = 0, \dots, L$ denote matrix-valued coefficients. Following the derivation of (5) in [13], one can obtain an FIR filter for predicting the matrix polynomial signal easily as

$$\mathbf{A}_k^- = \sum_{m=1}^M h_L(m)\mathbf{A}_{k-m}, \quad (7)$$

which is called as the *Matrix Polynomial Predictive Model* (MPPM).

B. The coefficients of the MPPM

Despite the simplicity of the MPPM, the degree of the polynomial L , the length of the FIR filter M , and the coefficients $h_L(m), m = 1, \dots, M$ should firstly be determined. The optimal coefficients in (7) can be derived as follows.

If the matrix polynomial signal modeled as (6) can be predicted by (7), the right side of (6) can be rewritten as

$$\sum_{l=0}^L \mathbf{P}(l)k^l = \sum_{m=1}^M h_L(m) \sum_{l=0}^L \mathbf{P}(l)(k-m)^l. \quad (8)$$

By exchanging the summation on the right side of (8), one can obtain $L+1$ separate equations as

$$\mathbf{P}(l)k^l = \sum_{m=1}^M h_L(m)\mathbf{P}(l)(k-m)^l. \quad (9)$$

When the two sides of (9) are left multiplied with the pseudo-inverse of $\mathbf{P}(l)$ simultaneously, (9) can be rewritten as

$$k^l = \sum_{m=1}^M h_L(m)(k-m)^l, \quad (10)$$

which is no longer a matrix-valued equation. This equation also shares the same form as the one in the coefficient derivation of the original PPM or FIR filter in [12] or [13]. Therefore, the rest of the derivation should be the same as the

one in those two references. If the details are needed, please refer to the two references.

The optimal coefficients for the polynomials whose degrees are no greater than 2 are provided in this paper. If one need the optimal coefficients for the polynomials whose degrees are larger than 2, please refer to the derivation in [13]. The simplest case is $L = 0$, that is, the matrix signal is assumed constant. In this situation, the optimal coefficients for the filter are given as [13]

$$h_0(m) = \frac{1}{M}. \quad (11)$$

In other word, \mathbf{A}_k can be predicted by (7) via the mean of the previous M values. When $L = 1$ and 2, the optimal coefficients are respectively given as [13]

$$h_1(m) = \frac{4M - 6m + 2}{M(M - 1)}, \quad (12)$$

and

$$h_2(m) = \frac{9M^2 - (9 - 36m)M + (30m^2 - 18m + 6)}{M(M - 1)(M - 2)}. \quad (13)$$

From (11) to (13), one can see that the coefficients of the predictive FIR filters have nothing to do with the matrix coefficients of the polynomials, which will make their following applications very convenient.

IV. MATRIX POLYNOMIAL PREDICTIVE MODEL IN THE PARAFAC DECOMPOSITION

A. ALS-MPPM

Let \mathbf{A}_k , \mathbf{B}_k , and \mathbf{C}_k respectively denote the estimates of the corresponding loading matrices in the k -th iteration of the ALS. If the convergent progresses of the \mathbf{A}_k , \mathbf{B}_k , and \mathbf{C}_k are described by the polynomials as given (6), they can be predicted by (7) via their previous M estimates \mathbf{A}_{k-m} , \mathbf{B}_{k-m} , and \mathbf{C}_{k-m} ($m = 1, \dots, M$). When the predictions, \mathbf{A}_k^- , \mathbf{B}_k^- , and \mathbf{C}_k^- are taken as the known values of the next iteration for the ALS, a new version of the ALS denoted as ALS-MPPM is found.

It should be emphasized that the rate of the convergence can be controlled by the degree of the polynomial model (6). It means that the higher the degree of the polynomial is, the faster the algorithm convergence will be. Such an acceleration can be illustrated taking the PPM for a scalar polynomial signal as an example. Assume that x_{k-2} and x_{k-1} are the estimates of x in the previous iterations, whereas the x_k^- is the prediction of x_k using (5) with $L = 1$ and $M = 2$. Let ε_k^- denote the absolute error of the estimation x_k^- , i.e., $\varepsilon_k^- = |x_k^- - x|$. Then, ε_k^- can be written as $\varepsilon_k^- = \varepsilon_{k-1} + \Delta x_k^-$, where $\Delta x_k^- = x_k^- - x_{k-1} < 0$ is the error increment of x_k^- , and ε_{k-1} is the absolute error of the estimation x_{k-1} . It is obvious that $\varepsilon_k^- < \varepsilon_{k-1}$. In other words the predicted x_k^- should be a better estimate of x than x_{k-1} . Therefore, the convergence is accelerated via such a prediction. Nevertheless, the prediction could fail. Such failure usually happens when a rapid convergence occurs near the values of the convergence. When $L = 1$, the progresses of the signal x to its convergence value is assumed to be linear.

In a rapid convergence progress, once $|\Delta x_{k-1}| > 2\varepsilon_{k-1}$, the prediction will fail since $\varepsilon_k^- > \varepsilon_{k-1}$, as illustrated in Fig.1. As a result, a faster convergence means a higher possibility of prediction failure. Therefore, a higher degree of polynomials does not guarantee better performance.

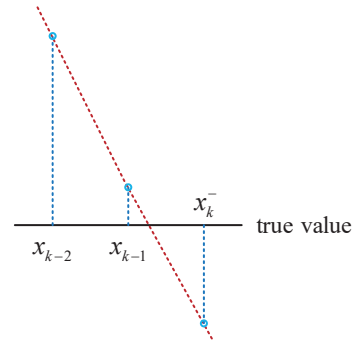


Fig. 1. An illustration of prediction failure using PPM in a rapid convergence near its value of convergence.

The simplest solution to handle such a failure is to abandon the prediction when $\varepsilon_k^- > \varepsilon_{k-1}$. Namely, a standard ALS iteration is applied after the prediction fails have been detected. The ALS-MPPM is summarized in TABLE I.

TABLE I
SUMMARIZATION OF THE ALS-MPPM ALGORITHM

Given: tensor data \mathcal{X} , initial estimates \mathbf{A}_0 , \mathbf{B}_0 , and \mathbf{C}_0 , degree of the polynomial L , and the length of the filter M .
Initialize: Calculate $h_L(m)$, $m = 1, \dots, M$ offline; Define the collection $\mathcal{A} = [\mathbf{A}_0, \dots, \mathbf{A}_0]$, $\mathcal{B} = [\mathbf{B}_0, \dots, \mathbf{B}_0]$, and $\mathcal{C} = [\mathbf{C}_0, \dots, \mathbf{C}_0]$; Set $k = 1$.
Do Predict \mathbf{A}_k^- , \mathbf{B}_k^- , and \mathbf{C}_k^- with \mathcal{A} , \mathcal{B} , and \mathcal{C} using (7); Calculate ε_k^- ; If $\varepsilon_k^- > \varepsilon_{k-1}$ Adopt the re-prediction strategy; End If Do ALS Update: $\mathbf{A}_k = \mathbf{X}_{(1)}[(\mathbf{C}_k^- \odot \mathbf{B}_k^-)^T]^\dagger$; $\mathbf{B}_k = \mathbf{X}_{(2)}[(\mathbf{C}_k^- \odot \mathbf{A}_k)^T]^\dagger$; $\mathbf{C}_k = \mathbf{X}_{(3)}[(\mathbf{A}_k \odot \mathbf{B}_k)^T]^\dagger$; Calculate ε_k ; Update the collection of previous estimates $\mathcal{A} = [\mathbf{A}_k, \mathcal{A}_{2:M}]$, $\mathcal{B} = [\mathbf{B}_k, \mathcal{B}_{2:M}]$, and $\mathcal{C} = [\mathbf{C}_k, \mathcal{C}_{2:M}]$; $k = k + 1$;
Until some termination conditions are met.
Output: final estimation $\hat{\mathbf{A}}$, $\hat{\mathbf{B}}$, and $\hat{\mathbf{C}}$.

B. Connection to Existing Algorithms

The relationship between the PPM and the *Constant-Velocity* (CV) or *Constant-Acceleration* (CA) models has been covered in [12]. Moreover, when our ALS-MPPM is used to

do the PARAFAC decomposition, it is also closely related to the existing ALS-based algorithms.

1) *ALS*: When $L = 0$ and $M = 1$, the coefficient of the MPPM is $h_0(1) = 1$, according to (11). In this case, it can be seen that the prediction is

$$\mathbf{A}_k^- = \mathbf{A}_{k-1}, \tag{14}$$

which makes the ALS-MPPM equivalent to the original ALS. In this sense, ALS can be seen as a special case of our ALS-MPPM.

2) *ALS-(E)LS*: When $L = 1$ and $M = 2$, the coefficients of the MPPM can be calculated by (12) as

$$\mathbf{A}_k^- = 2\mathbf{A}_{k-1} - \mathbf{A}_{k-2}, \tag{15}$$

which can also be written as

$$\mathbf{A}_k^- = \mathbf{A}_{k-2} + 2(\mathbf{A}_{k-1} - \mathbf{A}_{k-2}). \tag{16}$$

It can be seen that (16) shares the same form as the (*Enhanced*) *Line Search* ((E)LS) techniques given in [8]- [11] with a relaxation factor of 2. Since such a factor is the maximum number allowed for a LS algorithm, it implies that our ALS-MPPM will converge much faster than the existing LS-based algorithms although a higher possibility of prediction failure could occur. However, unlike the original LS, in which the relaxation factor is chosen empirically, the MPPM is deterministic. Moreover, unlike the ELS, the MPPM does not need to solve a high degree polynomial equation, which means MPPM takes less computation than ELS does in each iteration.

V. NUMERICAL SIMULATIONS

In this section, Monte-Carlo simulations will be conducted to test the proposed ALS-MPPM algorithm against the existing ALS-based algorithms under different SNR conditions. The tensor \mathcal{X} in (1) is randomly generated as follows. The size of the tensor is set to $4 \times 4 \times 4$ with rank 3. Each element in the loading matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} is drawn from a uniform distribution on $[0, 1]$, and then \mathcal{X} is obtained by the matricization (2). \mathcal{X} is mixed with a Gaussian noise under different SNR level. The SNR (in dB) is defined as

$$\text{SNR} = 10 \log_{10} \frac{\|\mathcal{X}\|_F^2}{\sigma^2 IJK}, \tag{17}$$

where σ^2 denotes the variance of the Gaussian noise and $\|\cdot\|_F$ denotes the Frobenius norm. For each SNR, 1000 tensor samples are generated, and a 5% truncated mean of results is taken when processing the data to rule out the impact of the numerical issues. For each algorithm tested here, the convergence criteria include:

- Normalized Mean Square Error (NMSE) threshold: $\epsilon_{\text{NMSE}} = 10^{-8}$;
- Relative NMSE threshold: $\epsilon_k = |\epsilon_k - \epsilon_{k-1}| / \epsilon_{k-1} = 10^{-6}$;
- Maximum number of iterations $K = 5000$,

where the NMSE is calculated by

$$\epsilon_k = \frac{\|\mathbf{X}_{(1)} - \mathbf{A}_k (\mathbf{C}_k \odot \mathbf{B}_k)^T\|_F^2}{\|\mathbf{X}_{(1)}\|_F^2}. \tag{18}$$

In the first simulation, different configurations of the ALS-MPPM algorithm are tested, as shown in Fig. 2. The configuration with $L = 0$ and $M = 1$ is seen as the baseline since in this case the ALS-MPPM is the same as the original ALS. It can be seen that our ALS-MPPM is almost 10 times faster (in iterations) than the baseline configuration. It is also worth noticing that, when the degree remains unchanged, the number of iterations grows with the length of the MPPM. Therefore, it is recommended to use the shortest filter, that is, the length of the filter is recommended to set as $M = L + 1$. On the other hand, the matrix polynomial with a higher degree in this simulation does not guarantee a faster convergence, since a higher order MPPM usually results in a higher rate of prediction failure, as shown in Fig. 3.

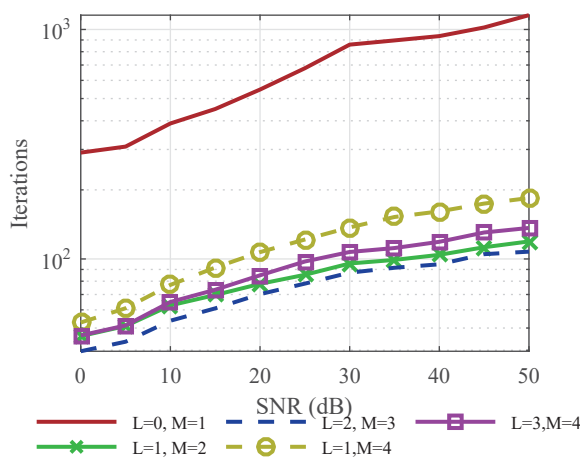


Fig. 2. Average number of iterations of several configurations of ALS-MPPM under different SNR condition. (Lower is better.)

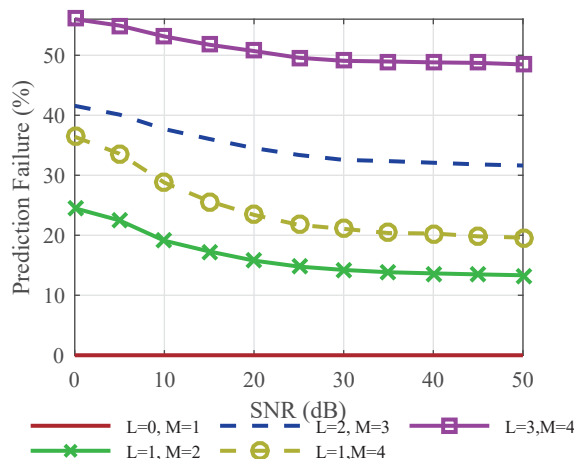


Fig. 3. Average prediction failure rate of several configurations of ALS-MPPM under different SNR condition.

In the second simulation, the ALS-MPPM is tested against several existing acceleration techniques for the ALS, including

two versions of Line Search (denoted as the ALS-LSB [9] and ALS-LSH [1]) and the Enhanced Line Search (denoted as the ALS-ELS [10]). The test results are shown in Fig. 4. It can be observed that our ALS-MPPM outperforms all reported acceleration techniques of the ALS.

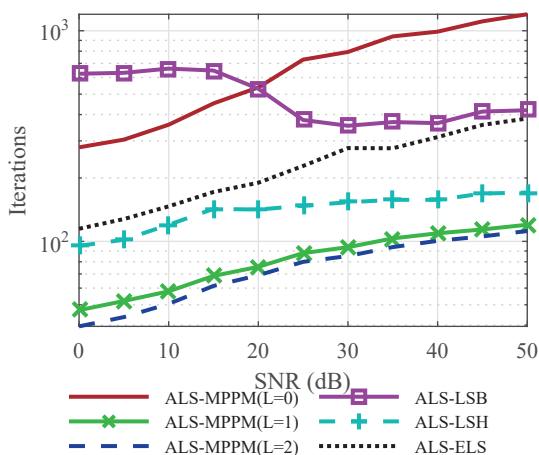


Fig. 4. Average number of iterations of the ALS-based algorithms. (Lower is better.)

VI. CONCLUSION AND FUTURE WORK

To accelerate the convergence of the PARAFAC decomposition of a tensor, a new ALS algorithm, called as the ALS-MPPM, has been proposed. Analyses show that the rate of the convergence of the proposed ALS-MPPM is closely related to the degree of the matrix polynomial. The polynomial with a high degree is preferred when an accelerative rate of convergence is expected. Although a high degree means a high possibility of prediction failure, a simple solution has also been provided to handle such failure. Moreover, the relationship between our ALS-MPPM and the existing ALS-based algorithms has also been analyzed. The results of numerical simulations show that the proposed ALS-MPPM outperforms the reported ALS-based algorithms in the literature while the analytical results are verified. Our future work includes the better strategies for handling the prediction failure and the derivation of the detailed relationship between the degree of the matrix polynomial and the rate of convergence.

REFERENCES

- [1] R. A. Harshman, "Foundations of the parafac procedure: Models and conditions for an "explanatory" multi-modal factor analysis," 1970.
- [2] J. D. Carroll and J.-J. Chang, "Analysis of individual differences in multidimensional scaling via an n-way generalization of eckart-young decomposition," *Psychometrika*, vol. 35, no. 3, pp. 283–319, 1970.
- [3] M. Sørensen, L. D. Lathauwer, P. Comon, S. Icart, and L. Deneire, "Canonical polyadic decomposition with a columnwise orthonormal factor matrix," *SIAM Journal on Matrix Analysis and Applications*, vol. 33, no. 4, pp. 1190–1213, 2012.
- [4] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.
- [5] A. Cichocki, D. Mandic, L. De Lathauwer, G. Zhou, Q. Zhao, C. Caiafa, and H. A. Phan, "Tensor decompositions for signal processing applications: From two-way to multiway component analysis," *IEEE Signal Processing Magazine*, vol. 32, no. 2, pp. 145–163, 2015.
- [6] B. C. Mitchell and D. S. Burdick, "Slowly converging PARAFAC sequences: Swamps and two-factor degeneracies," *Journal of Chemometrics*, vol. 8, no. 2, pp. 155–168, 1994.
- [7] W. S. Rayens and B. C. Mitchell, "Two-factor degeneracies and a stabilization of parafac," *Chemometrics and Intelligent Laboratory Systems*, vol. 38, no. 2, pp. 173–181, 1997.
- [8] R. T. Ross and S. Leurgans, "Component resolution using multilinear models," *Methods in enzymology*, vol. 246, pp. 679–700, 1995.
- [9] R. Bro, "Multi-way analysis in the food industry: models, algorithms, and applications," 1998.
- [10] M. Rajih, P. Comon, and R. A. Harshman, "Enhanced line search: A novel method to accelerate parafac," *SIAM journal on matrix analysis and applications*, vol. 30, no. 3, pp. 1128–1147, 2008.
- [11] D. Nion and L. De Lathauwer, "An enhanced line search scheme for complex-valued tensor decompositions. application in ds-cdma," *Signal Processing*, vol. 88, no. 3, pp. 749–755, 2008.
- [12] Y. Gao, J. Zhang, and J. Yin, "Polynomial prediction model and tracking algorithm of maneuver target," *Acta Aeronautica et Astronautica Sinica*, vol. 30, no. 8, pp. 1479–1489, 2009.
- [13] P. Heinonen and Y. Neuvo, "Fir-median hybrid filters with predictive fir substructures," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 36, no. 6, pp. 892–899, 1988.