

An extension of PARAFAC to analyze multi-group three-way data

Rotari, Marta; Diaz, Valeria Fonseca; De Ketelaere, Bart; Kulahci, Murat

Published in: Chemometrics and Intelligent Laboratory Systems

Link to article, DOI: 10.1016/j.chemolab.2024.105089

Publication date: 2024

Document Version Publisher's PDF, also known as Version of record

Link back to DTU Orbit

Citation (APA): Rotari, M., Diaz, V. F., De Ketelaere, B., & Kulahci, M. (2024). An extension of PARAFAC to analyze multi-group three-way data. *Chemometrics and Intelligent Laboratory Systems, 246, Article 105089.* https://doi.org/10.1016/j.chemolab.2024.105089

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

• Users may download and print one copy of any publication from the public portal for the purpose of private study or research.

- · You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Contents lists available at ScienceDirect



Chemometrics and Intelligent Laboratory Systems

journal homepage: www.elsevier.com/locate/chemometrics



An extension of PARAFAC to analyze multi-group three-way data

Marta Rotari^{a,*}, Valeria Fonseca Diaz^c, Bart De Ketelaere^c, Murat Kulahci^{a,b}

^a Department of Applied Mathematics and Computer Science, Technical University of Denmark, Kgs. Lyngby, Denmark

^b Department of Business Administration, Technology and Social Sciences, Luleå University of Technology, Luleå, Sweden

^c Department of Biosystems, MeBioS division, KU Leuven, Leuven, Belgium

ARTICLE INFO

Keywords:

PARAFAC

Factor analysis

Additive manufacturing

Multi-group data set

ABSTRACT

This paper introduces a novel methodology for analyzing three-way array data with a multi-group structure. Three-way arrays are commonly observed in various domains, including image analysis, chemometrics, and real-world applications. In this paper, we use a practical case study of process modeling in additive manufacturing, where batches are structured according to multiple groups. Vast volumes of data for multiple variables and process stages are recorded by sensors installed on the production line for each batch. For these three-way arrays, the link between the final product and the observations creates a grouping structure in the observations. This grouping may hamper gaining insight into the process if only some of the groups dominate the controlled variability of the products. In this study, we develop an extension of the PARAFAC model that takes into account the grouping structure of three-way data sets. With this extension, it is possible to estimate a model that is representative of all the groups simultaneously by finding their common structure. The proposed model has been applied to three simulation data sets and a real manufacturing case study. The capability to find the common structure of the groups is compared to PARAFAC and the insights into the importance of variables delivered by the models are discussed.

Introduction

Manufacturing organizations strive to adopt innovative production methods to enhance their efficiency and design flexibility. However, modern production techniques can be complex and not yet fully understood. Achieving a complete comprehension and optimization of these processes requires identifying the variables and components of the process, as well as the conditions under which production levels can be optimized. To achieve such an understanding, it is essential to develop descriptive and interpretable data-driven models that provide clear insights into the data.

Modern manufacturing industrial systems are often equipped with technology that can collect data from a large number of sources, resulting in a voluminous amount of information. This can produce datasets that are structured in three-way arrays, organized in observations collected for multiple variables and several conditions. For example, in a batch fermentation process, data collected in time for multiple variables and multiple batches will be in the form of a three-dimensional array [1,2]. A second example is processes based on sensory data, which may come from different channels and multiple devices simultaneously [3–5].

Several approaches have been proposed to model three-way data using supervised or unsupervised methods, ranging from linear models [6–10] to black-box models for tensors [11,12]. However, to find

the optimal settings of the system, it is required to fully understand the conditions and the interactions among the variables involved in the production processes. Linear models remain a predominant tool for three-way data as they can identify key factors that impact the process and guide decision-making. Although black-box models may offer high prediction accuracy, they lack the interpretability and transparency necessary to understand the system.

The most predominant methodologies within the class of linear models for unsupervised three-way data have been factor analysis methods such as PARAFAC [8,13] and Tucker3 model [9,14]. These models are used by analyzing the covariance structure among the variables and among different experimental conditions providing the so-called loadings matrices whose values are indicative of the importance of each variable and each condition in the total variability of the three-way data [2,15]. The scores that are delivered for the observations of the data allow for the identification of the dispersion among the samples and the differences in the variance explained by each component of the model.

In high-dimensional data analysis, it is not uncommon to encounter datasets that exhibit a clear and distinguishable organizational structure, which is categorized into blocks or groups. In the context of multi-block data, each block consists of a collection of variables that

https://doi.org/10.1016/j.chemolab.2024.105089

Received 2 August 2023; Received in revised form 20 January 2024; Accepted 18 February 2024 Available online 21 February 2024

0169-7439/© 2024 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY license (http://creativecommons.org/licenses/by/4.0/).

^{*} Corresponding author. *E-mail address:* mrot@dtu.dk (M. Rotari).

exhibit certain characteristics or thematic significance, assessed on the same individuals. This situation may also be represented as a composite of datasets, often acquired from the same group of individuals but including distinct variables across the many datasets. This may be the case with batch production where different kinds of measurements are acquired along the steps of the manufacturing process. The literature extensively covers various methods of data analysis for multi-block data. Examples of unsupervised methods include multi-block principal components analysis (MB-PCA) [16] and multiple co-inertia analysis (MCOA) [17], on the other hand, multi-block partial least squares (MB-PLS) is an example of supervised method [18].

In many industrial processes, the existence of sub-groups within the observations is a common occurrence that can greatly impact the analysis and interpretation of the results. The multi-group structure in the data refers to the presence of distinct groups that may exhibit unique characteristics and behaviors. Groups of observations are likely to happen when the observations or samples share a specific characteristic that identifies the group [19-21]. Various factors, such as demographic variables, geographic location, or other relevant criteria, may define these groups. In the manufacturing processes, the observations collected during production are linked to the final products, which individuate a multi-group structure of the observations. Understanding the multi-group structure in the data is crucial for accurate analysis, as failing to consider these groups can lead to biased results and erroneous conclusions. Thus, it is important to carefully identify and account for the multi-group structure in the model in order to ensure that research findings are valid and applicable to all relevant groups [21,22].

When grouping structures exist, the sources of variability explained in the resulting model can be dominated by specific groups without representing other groups' variability. Traditional models for three-way arrays that do not take into account the group structure can be affected by the greater variance of a dominant group and fail to represent the entire data set. Therefore, the importance of specific variables or conditions that are extracted from the model parameters may not hold for all the groups simultaneously, leading to potentially erroneous conclusions. One practical solution to take into account the presence of multiple groups is to analyze the groups separately. However, this strategy would result in an abundance of parameters and the lack of a single unifying model of the complete system.

Extensive research has been conducted on linear models for highdimensional data that take into account the grouping structure in the case of two-way data. Some examples of unsupervised models are the Common principal components analysis (CPCA) [23], that involves examining the variance-covariance matrices of the groups and identifying shared orthogonal vectors of loadings; Multi-group principal component analysis (MGPCA) [24], which involves deriving the common loadings by performing an eigenanalysis on the variance-covariance matrix within each group. Additional techniques for analyzing two-way multi-group datasets are outlined in the works of Kallus et al. [25], Tenenhaus et al. [21] and Eslami et al. propose a multi-group PLS model [26]. It has been shown how the resulting model, when considering the group structure of the observations, can be informative in understanding the common cause of variability across all groups [19,27]. Through multi-group models, it is possible to interpret the functionality of the production systems for several groups simultaneously.

The current study presents an extension of the PARAFAC method that adjusts the model parameters according to the multi-group structure of the data. The model aims to identify a common structure by calculating common loadings in the presence of multi-groups in the three-way data. The article is organized as follows. First, the extended PARAFAC method is presented with the proposed algorithm to fit the model. Then, the case study with simulated and real data is presented, highlighting the multi-group structure for three-way data. The results of the PARAFAC and the extended PARAFAC models are presented and compared. The models are used to analyze the most important variables and conditions in the datasets delivered by each model. Finally, the conclusions are presented.

1. Methods

This section provides a brief description of the PARAFAC model and an algorithm for its solution, Alternating Least Squares (ALS). Following that, we present the proposed extension aimed to handle a multi-group structure.

Three-way arrays are represented by the bold italic uppercase letter $X \in \mathbb{R}^{I \times J \times K}$ whose elements are x_{ijk} where i = 1, ..., I, j = 1, ..., Jand k = 1, ..., K. The three-way entries are denoted by modes a, band c as shown in Fig. 1(a). In this structure, mode a represents the observations, mode b represents the variables, and mode c represents different conditions, such as time, different levels of temperatures, etc. The elements that belong to modes a and b with mode c fixed will be denoted by $X_{..k} = (x_{ijk})$ where k is fixed, i = 1, ..., I and j = 1, ..., J. Scalars are denoted in nonbold italic characters, such as x_{ijk} or I. Column vectors are denoted by bold lowercase characters, such as a, whereas two-way arrays or matrices are denoted by a capital bold letter, such as A. The letters I, J, K are reserved for indicating the dimensions and N for indicating the number of components. Superindex T represents the transpose operator.

Let us consider a three-way array X, its unfolded (matricized) matrix in one of the modes is **X**. Specifically, the unfolded matrix in mode *a* is represented as X_a with dimensions $I \times JK$, the unfolded matrix with respect to mode *b* is X_b with dimensions $J \times IK$. Lastly, the unfolding in mode *c* is referred to as X_c with dimensions $K \times IJ$.

1.1. The PARAFAC method

The PARAFAC method is a decomposition method of the three-way array into three matrices: the score matrix **A** and two loadings matrices **B** and **C**. Matrix **A** is referred to as scores as its entries represent the numerical value of the observations [8,15]. Matrices **B** and **C** are commonly referred to as loadings as they represent the numerical value of variables and conditions, respectively. The PARAFAC model can be seen as a particular case of the Tucker model, where the core array has the same dimensions $N \times N \times N$, has ones on the superdiagonal and zeros otherwise [14,28]. Considering a three-way array *X*, a PARAFAC model with *N* components can be written as follows:

$$x_{ijk} = \sum_{n=1}^{N} a_{in} b_{jn} c_{kn} + e_{ijk}$$
(1)

where $\mathbf{A} = (a_{in})$ is the score matrix $\mathbb{R}^{I \times N}$, $\mathbf{B} = (b_{jn}) \in \mathbb{R}^{J \times N}$ is the loadings matrix of mode *b*, and $\mathbf{C} = (c_{kn}) \in \mathbb{R}^{K \times N}$ is the loadings matrix of mode *c*. A visual representation of the PARAFAC model is shown in Fig. 1(b). The PARAFAC model results from finding the model matrices that minimize the sum of squares of the residuals:

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C}} \sum_{ijk} (x_{ijk} - \sum_{n=1}^{N} a_{in} b_{jn} c_{kn})^2 = \min_{ijk} \sum_{ijk} (e_{ijk})^2$$
(2)

Several algorithms have been proposed to solve Eq. (2) for A, B, and C such as Alternating Least Squares (ALS) [8,29,30], Derivative Computations [31] and direct (non-iterative) procedures [32,33]. We present here the ALS algorithm as it is the original and most widely used method to estimate the PARAFAC model [8].

ALS algorithm

The ALS algorithm aims at calculating the model matrices **A**, **B** and **C** that minimize the sum of squared residuals in Eq. (2) using an iterative procedure. At each step *t*, it calculates and updates the matrices until convergence or until the maximum number of iterations is reached rendering a series of matrices $\{\mathbf{A}^{(t)}\}$, $\{\mathbf{B}^{(t)}\}$ and $\{\mathbf{C}^{(t)}\}$ [8,15,28]. Let us consider the unfolded matrices $\mathbf{X}_{\mathbf{a}}$ ($I \times JK$), $\mathbf{X}_{\mathbf{b}}$ ($J \times IK$) and $\mathbf{X}_{\mathbf{c}}$ ($K \times IJ$) of the three-way array X. The first step consists in initializing the



Fig. 1. (a): Three-way array X represented by the modes a, b and c. (b): PARAFAC model representation, X three-way array, A, B and C model loadings matrices and E three-way residual array.



Fig. 2. (*a*): Three-way array $X = x_{ijk} \in \mathbb{R}^{I \times J \times K}$ which present a multi-group structure. Along the mode *a* the observations are divided into *G* groups, such that group $g = \{1, ..., G\}$ is represented by $1, ..., I_g$ rows and $\sum_{g=1}^G I_g = I$. (*b*): Representation of the new proposed model.

S

matrices $B^{(1)}$ and $C^{(1)}$ as random matrices. The matrix $A^{(t)}$ is calculated by solving the minimization problem

$$\min_{\mathbf{A}} \| \mathbf{X}_{\mathbf{a}} - \mathbf{A}^{(t)} \mathbf{Z}^{T} \| \text{ where } \mathbf{Z} = \mathbf{C}^{(t-1)} \otimes \mathbf{B}^{(t-1)}, \ \mathbf{X}_{\mathbf{a}} \in \mathbb{R}^{I \times JK}$$
(3)

 \otimes represent the Khatri–Rao product. Matrices $B^{(t)}$ and $C^{(t)}$ are calculated in an analogous way.

$$\begin{split} \min_{\mathbf{B}} \| \mathbf{X}_{\mathbf{b}} - \mathbf{B}^{(t)} \mathbf{Z}^{T} \| & \text{where } \mathbf{Z} = \mathbf{C}^{(t-1)} \otimes \mathbf{A}^{(t)}, \ \mathbf{X}_{\mathbf{b}} \in \mathbb{R}^{J \times IK} \\ & \text{and} \\ \\ \min_{\mathbf{C}} \| \mathbf{X}_{\mathbf{c}} - \mathbf{C}^{(t)} \mathbf{Z}^{T} \| & \text{where } \mathbf{Z} = \mathbf{B}^{(t)} \otimes \mathbf{A}^{(t)}, \ \mathbf{X}_{\mathbf{c}} \in \mathbb{R}^{K \times IJ} \end{split}$$

1.2. Proposed method for multi-group data

The multi-group three-way array X is composed by the same three modes a, b, and c, when the observations in mode a are divided into G groups. A visual representation of this structure is shown in Fig. 2(a). Each group $g \in \{1, ..., G\}$ is represented by $1, ..., I_g$ rows, such that $\sum_{g=1}^{G} I_g = I$. The sub-three way array that contains the observations belonging to group g is represented by $X_g \in \mathbb{R}^{I_g \times J \times K}$ where I_g is the number of observations in such group. The two main characteristics of the proposed method are:

- A, B, and C loadings matrices are estimated from the decomposition of the three-way array *X*, constrained by its group structure.
- The resulting loadings matrices **B** and **C** focus on the common variability among the groups.

Thus, the resulting model is not dominated by a specific group but is representative of all the groups uniformly. The optimization function of the proposed method is defined as in Eq. (2) with an adjustment for the grouping structure

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C}} \sum_{ijk} (x_{ijk} - \sum_{n=1}^{N} a_{in} b_{jn} c_{kn})^2$$

uch that $\mathbf{B} = \sum_{g=1}^{G} w_g \mathbf{B}_{\mathbf{g}}$ and $\mathbf{C} = \sum_{g=1}^{G} w_g \mathbf{C}_{\mathbf{g}}, g \in 1, \dots, G$ (4)

and the score matrix **A** is a multi-group matrix and $\mathbf{A}_{\mathbf{g}} = (a_{in})$ where $i = 1, \ldots, I_g$ and $n = 1, \ldots, N$ is the sub-matrix containing the observations that belong to the group g. The matrices \mathbf{B}_g and \mathbf{C}_g are the loadings matrices that correlate to the observations in the group g and w_g are the weights associated with each group. The matrices **B** and **C** are common loading matrices for all groups in X. They are determined using the weighted scheme approach described in the next section. We present here the extension of the ALS algorithm to solve the problem presented in Eq. (4).

Extended ALS algorithm

Similar to the ALS algorithm for the PARAFAC model, the matrices **A**, **B** and **C** are updated at each step of the algorithm by minimizing the sum of squares of residuals between the input matrix *X* and the product of the three model matrices. Thus, we shall generate a series of matrices $\{\mathbf{A}^{(t)}\}, \{\mathbf{B}^{(t)}\}$ and $\{\mathbf{C}^{(t)}\}$ that reach convergence at each successive step.



Fig. 3. Representation of the calculation of a matrix common to all groups in a generic step t of the algorithm. This scheme applies also to the calculation of C.

The algorithm is initialized by calculating the PARAFAC model with solution denoted by $\mathbf{A}^{(1)}, \mathbf{B}^{(1)}, \mathbf{C}^{(1)}$. Therefore, in the absence of groups, this extension returns to the conventional PARAFAC model. In the presence of multiple groups, the algorithm first solves the PARAFAC model for each group and then finds the common solution across all groups. We describe here the process of calculating **B** at step *t*.

We consider the unfolded two-way array $\mathbf{X}_{\mathbf{b}}$ along the mode *b* of dimension $J \times IK$ and reorganize the columns by grouping the observations that belong to the same group. The matrix $\mathbf{X}_{\mathbf{b}}$ is a multigroup structure matrix where each group is of dimensions $J \times I_g K$. For each group thus created, the group loadings matrix $\mathbf{B}_{\mathbf{g}}^{(i)}$ is calculated by solving the objective function:

$$\min_{\mathbf{B}_{\mathbf{g}}} \| \mathbf{X}_{\mathbf{b}\mathbf{g}} - \mathbf{B}_{\mathbf{g}}^{(t)} \mathbf{Z}_{\mathbf{g}}^{T} \| \text{ where } \mathbf{Z}_{\mathbf{g}} = \mathbf{C}^{(t-1)} \otimes \mathbf{A}_{\mathbf{g}}^{(t-1)}$$
(5)

where $\mathbf{A}_{\mathbf{g}}^{(t-1)}$ denotes the rows of the matrix $\mathbf{A}^{(t-1)}$ that correspond to the group g. The procedure is repeated for $g = 1, \ldots, G$ obtaining the group loadings matrices $\{\mathbf{B}_{1}^{(t)}, \ldots, \mathbf{B}_{G}^{(t)}\}$. The matrix $\mathbf{B}^{(t)}$ corresponds to the common loadings matrix as represented in Fig. 3. These common loadings $\mathbf{B}_{\mathbf{g}}$.

The calculation of C_g and C at step t follows the same procedure. The three-way array X is unfolded along the c mode as the matrix X_c of dimensions ($K \times IJ$), and the observations are reorganized so that X_c shows a multi-group structure matrix. For each group g = 1, ..., G we calculate the corresponding group loadings matrix $C_g^{(t)}$ by solving the problem given by:

$$\min_{\mathbf{C}_{\mathbf{g}}} \|\mathbf{X}_{\mathbf{c}\mathbf{g}} - \mathbf{C}_{\mathbf{g}}^{(t)} \mathbf{Z}_{\mathbf{g}}^{T}\| \text{ where } \mathbf{Z}_{\mathbf{g}} = \mathbf{B}^{(t)} \otimes \mathbf{A}_{\mathbf{g}}^{(t-1)}$$
(6)

Here, the algorithm uses the matrix $\mathbf{A}^{(t-1)}$ of the previous iteration and the latest updated loadings $\mathbf{B}^{(t)}$. For each group $g = 1, \ldots, G$ we calculate the group loadings matrices $\{\mathbf{C}_1, \ldots, \mathbf{C}_G\}$ and the matrix $\mathbf{C}^{(t)}$ becomes the common loadings matrix to all the groups obtained by weighted mean.

Finally, to update the score matrix **A**, the unfolded matrix **X**_a is considered and solve the same PARAFAC step for each group *g* to obtain $\{\mathbf{A}_{1}^{(t)}, \dots, \mathbf{A}_{G}^{(t)}\}$. Finally we repeat the above steps until convergence or the maximum number of iterations is reached. Algorithm 1 describes the entire procedure of the proposed method.

The weighting scheme is determined based on the optimization problem of Eq. (4). There are several options to obtain a solution for the weights in order to determine a common **B** and **C**. One option is that the solution is given by a simple mean over the groups. In this scenario, the weights w_g will be uniform across all the groups. A second option would consist of a weighted average, where the weights correspond to the inverse of the variance of each group explained within the model. Because the model needs to represent all groups uniformly, these weights are determined based on the variance explained for each group. If a group has a high variance explained, its weight will be lower. Correspondingly, the groups with lower variance explained will have a higher weight. This renders a common loadings matrix represented by variance homogeneity. An alternative option is to determine the exact weight values for each group. The usefulness of this solution is case-specific dependent. Indeed, it may be appropriate to assign more significance to the group that we know has more relevance or a higher concentration of the chemical compound being studied, among other possible factors.

Algorithm 1 Extended ALS algorithm for PARAFAC model with multi-group data

Input: *X* data, list *G* group identification, *N* number of components **A**, **B**, **C** \leftarrow parafac(*X*, *N*) group-weights = ($w_1, ..., w_G$)

while convergence or max_iterations do

for
$$g = 1, ..., G$$
 do
 $Z_g = C \otimes A_g$
 $\min_{B_g} ||X_{bg} - B_g Z_g^T||$
end for
For each model component $n = 1, ..., N$
 $B = \sum_{g=1}^G w_i B_i$
for $g = 1, ..., G$ do
 $Z_g = B \otimes A_g$
 $\min_{C_g} ||X_{cg} - C_g Z_g^T||$
end for
For each model component $f = 1, ..., N$
 $C = \sum_{g=1}^G w_i C_i$
for $g = 1, ..., G$ do
 $Z = C \otimes B$
 $\min_{A_g} ||X_{ag} - A_g Z^T||$
end for
 $A = [A_1, ..., A_G]^T$
hd while

Output: A, B and C

et

2. Data overview

To evaluate the performance of the proposed model, three simulated data sets and one real case study from additive manufacturing were used. Each simulation reflects a different scenario for the grouping structure of three-way arrays. The first simulation involves simulating data in batches for k = 1, 2, 3, where each batch is composed of three groups. The second and third simulations use model loadings **A**, **B**, and **C** to generate the three-way array. Various shifts in group variance, observations and variables were introduced to test the algorithm's performance. In the real case study, we examined a three-way array **X** that represents four batches from an additive manufacturing process.

The PARAFAC model and the proposed extension were fitted to each dataset. The two methods were compared based on their ability to homogenize the variance across groups, create common loadings for all the groups and provide informative insights through variable loadings. Our goal was to demonstrate the efficacy of the proposed approach in accurately analyzing complex multi-group array data and identifying a common model. All the analyses were performed in R 4.0 with in-house codes.

2.1. Data: Simulated data 1

The first simulation defines a three-way array $X \in \mathbb{R}^{(90\times40\times3)}$ consisting of three groups of observations. In this simulation we define the matrices \mathbf{X}_{k} (k = 1, 2, 3), each composed of three groups. Each group of each matrix is simulated using a multivariate normal distribution $N_{40}(\cdot, \cdot)$ with mean $\mu = 0$ and standard deviations $\sigma = 2, 5, 1$ as shown below:





where **I** represents the identity matrix of the corresponding order. After concatenating the groups and the matrices $X_{..k}$, an error with a normal distribution of $\mu = 0$ and standard deviation $\sigma = 0.05$ was added to obtain the final array X.

It is worth noticing that the second group, $X_{..k_{g=2}}$, has been assigned a larger standard deviation to challenge the ability of the proposed model to homogenize the variances across the groups. This will test the algorithm's capability to overcome any potential biases towards a particular group with larger variation and to generate informative loadings that reflect the entire three-way array.

2.2. Data: Simulated data 2

The second simulation consists of a three-way array $X \in \mathbb{R}^{(90\times20\times3)}$ with a multi-group structure consisting of three groups as well. In this case, we simulated each sub-array $X_g \in \mathbb{R}^{I_g \times 20\times3}$ for $I_g = 30, 30, 30$ and concatenate them vertically to obtain the final three-way array. Each sub-array X_g is obtained as the product of matrices **A**, **B** and **C** generated using the Normal distribution for **A** and Uniform distribution for **B** and **C**.



 $\begin{aligned} \mathbf{A} &\sim N(0, \mathbf{I}) \\ \mathbf{B} &\sim U(0, \mathbf{I}) \\ \mathbf{C} &\sim U(0, \mathbf{I}) \end{aligned}$

Within the established framework, a shift in the model matrix was introduced to create separation between groups. The aim is to examine the effect of the shifts on both models. In this case, a shift was introduced in the score matrix **A** corresponding to the first and third groups. Specifically, $\mathbf{A}_{3_{g=1}}$ was multiplied by factor of 4 and $\mathbf{A}_{3_{g=3}}$ was multiplied by -1. After the concatenation of the groups, an error with a normal distribution of $\mu = 0$ and standard deviation $\sigma = 0.05$ was added to obtain the final array.

2.3. Data: Simulated data 3

A third simulation is obtained using a similar framework as in simulated data 2 of Section 2.2. We generate a three-way array $X \in \mathbb{R}^{(90\times20\times3)}$, vertical concatenation of three sub-arrays $X_g \in \mathbb{R}^{I_g \times 20\times3}$ for $I_g = 30, 30, 30$. Each sub-way array X_g is obtained as the product of matrices **A**, **B** and **C**, as in the previous example.

Here, a shift in the loadings matrix **B** was introduced. Specifically, the first column of the first group was multiplied by 5, (i.e. $5\mathbf{B}_{1_{g=1}}$). After the concatenation of the groups, an error with a normal distribution of $\mu = 0$ and standard deviation $\sigma = 0.05$ was added to obtain the final array.

The objective of this simulation study is to test the capability of the proposed model in generating common loadings that mitigate the variations unique to each group, thereby enabling the recovery of common information shared across all groups. We aim to showcase that the model captures the underlying common patterns while minimizing the impact of group-specific variations.

2.4. Case study: Additive manufacturing

The real case study corresponds to an additive manufacturing process for high-volume 3D printing. The so-called Selective Thermoplastic Electrophotographic Process (STEP) [34] takes place by rendering a three-dimensional object from a digital model. This process generates a 3D bulk structure by fusing and pressing super-thin layers. Multiple sensors are positioned throughout the production chain on the new manufacturing line measuring several variables for each super-thin layer.

Our study focuses on the analysis of a data set consisting of four batches that form a three-way array denoted as X. In the X three-way array, mode b contains 53 continuous variables collected through the sensors located on the printing machine. The variables are labeled from V1 to V53 due to confidentiality. Mode a represents the observations of the variables collected for each super-thin layer. Finally, mode c represents the different batches, as shown in Fig. 4(b). Each batch consists of three final groups of products. The association between the final products and the observations creates a multi-group structure of three groups represented in the X three-way array, Fig. 4(a). The three-way array was previously scaled by group [35]. Group scaling involves the scaling of each variable within the sub-array X_{g} based on the mean and standard deviation calculated within the respective group denoted as $X_{\cdot,j_{e_g}}$ for $g \in 1, ..., G$. This approach ensures that variables within each group are standardized in relation to the group-specific distribution characteristics.

This case study aims at gaining insights into the covariance structure among the variables and the different batches. Specifically, we aim to analyze the loadings matrix **B**, which reveals the relative importance of each variable in explaining the total variability of the three-way array. To accomplish this, we employ both the PARAFAC model and the proposed extension.



(a) Representations of the multi-group structure.

(b) Three-way array representation in the Additive manufacturing studying case.

Fig. 4. (a): The creation of a multi-group structure in the process data as a result of the relationship between the batch products and the observations in the process data. (b) Observations, variables and batches organization in the three-way array in the Additive manufacturing case study.

 Table 1

 Cumulative % of Variance Explained by Groups in Simulated dataset 1.

 PAPATAC Algorithm

PARAFAC Algorithin										
	Comp 1	Comp 2	Comp 3	Comp 4	Comp 5	Comp 6	Comp 7	Comp 8	Comp 9	Comp 10
Group1	1.26	2.76	3.71	5.10	6.25	7.89	8.60	10.74	12.40	14.15
Group2	3.61	7.43	11.33	14.72	18.44	21.37	25.28	27.71	31.26	33.71
Group3	1.15	2.33	3.44	4.19	4.91	5.24	6.61	7.76	8.68	10.18
Proposed Algorithm										
	Comp 1	Comp 2	Comp 3	Comp 4	Comp 5	Comp 6	Comp 7	Comp 8	Comp 9	Comp 10
Group 1	2.27	4.45	6.39	9.05	9.42	12.93	15.71	16.65	17.98	20.94
Group2	2.40	4.79	7.00	8.50	10.41	13.17	15.67	18.00	18.67	21.60
Group3	2.63	4.82	6.71	9.50	11.30	13.82	14.90	16.83	17.78	22.04

3. Results and discussion

As a rule of thumb, we present all the results across all data sets using models with only 10 components. We adopted this choice to effectively showcase two key objectives. Firstly, through Simulated data set 1, our aim was to demonstrate the consistency of the explained variance across all groups. Additionally, Simulated data sets 2 and 3 were conducted to provide further evidence of the proposed model's ability to generate common loadings that are shared among all groups. It is important to note, however, that the determination of the number of components is context-dependent and should be tailored to the specific characteristics of each individual case, and may vary accordingly. The determination of the components can be done using a similar approach as in the PARAFAC model, which involves residual analysis, loading analysis, split-half analysis, cross-validation etc. [28].

Simulated data sets

The percentage of variance explained by both models for each group in Simulation 1 is summarized in Table 1. The PARAFAC model reveals that the second group has a higher representation than the other two groups in all of the components. This result is in line with the ranking of the standard deviations set for each group as presented in Section 2.1.

To compare the representation of the groups across all of the components, the ratio of the explained variances for each group with respect to group 3 was calculated (See Fig. 5). Throughout the different components, the second group dominated the representativity with a percentage of explained variance 3 times the variability group 3. This result showed the dominance of the second group for the resulting PARAFAC model with a lower representation of the other groups. This is indicative of a model that represents only one group rather than being representative of the entirety of the data. In contrast, the proposed model demonstrates a uniform explained variance across all 3 groups, providing evidence of a model that uniformly represents the entire dataset. This is supported by the ratios in Fig. 5, which depict

a more uniform explained variance across all groups throughout the components.

Fig. 6 displays the scores of the 1st and 2nd components for Simulation 2, described in Section 2.2. In this case, a shift in the score matrix **A** was applied in the data simulation. The analysis of the results of the two models, therefore, focuses on the visualization of the first two scores.

In the case of the PARAFAC model, while the scores of the second and third groups nearly overlap and vary along the same direction, the scores of the third group are almost perpendicular to those of the first two groups, with a much larger dispersion. In contrast, the proposed model demonstrates nearly overlapping scores across all groups. Additionally, it can be noticed that the direction of variation for all groups is closer to one another. This provides evidence that the proposed model effectively mitigates the inherent shifts between groups and captures shared directions.

In Fig. 7, the **B** loadings matrices are presented for both the PARAFAC and proposed models referred to Simulation 3, described in Section 2.3. The black profile in the figures represents the original loadings matrix and a significant shift is highlighted represented by the first group. This extreme shift scenario was deliberately chosen to assess the algorithm's performance under challenging conditions. However, in practical scenarios, the variations between groups are expected to be comparatively smaller.

Upon comparing the results of the two models, noticeable differences can be observed. The PARAFAC model exhibits a greater bias towards the shift, resulting in **B** loadings profiles with higher levels of noise. These loadings exhibit extremely high peaks, influenced by the larger variability present in the first group. Conversely, the proposed model displays less pronounced peaks and generally smoother trend profiles. The loadings of the proposed model are closer to those of groups 2 and 3, indicating reduced bias from group 1.

These findings suggest that the proposed model can identify the common patterns present in all groups while mitigating any bias towards a specific group. Overall, these results highlight the potential of



Fig. 5. The ratios of the explained variances for each group with respect to group 3 in simulated data set 1.



Fig. 6. First two scores plot for the PARAFAC model and the Proposed model with a 99% confidence interval ellipse related to Simulation 2, described in Section 2.2.



Fig. 7. The PARAFAC and proposed model loadings matrices B for the dataset 3, Section 2.3.

the proposed model to serve as a reliable tool for analyzing multi-group datasets.

Case study: Additive manufacturing

The results of the PARAFAC and the proposed models in terms of the explained variance per group for the three-way array X are summarized in Table 2. Consistent with the findings observed in the simulated datasets, the PARAFAC model exhibited a relatively higher representation of the first group compared to the subsequent two groups. This disparity in groups representation becomes more pronounced, particularly in the higher-order components. In contrast, the proposed model, resulted in a more uniform representation of all three groups.

These findings are further validated by the ratios between the explained variances of each group with respect to group 3, in Fig. 8.

The initial components show that the PARAFAC model exhibits a twofold ratio in representing the first group when compared to the third group. While the disparity diminishes with subsequent components, a preference for the first group persists, albeit to a lesser degree. This observation suggests that the PARAFAC model captures the variability of the first group more effectively than it does for the entire dataset. In contrast, the proposed model already explains similar degree of variances explained by all groups in the first few components. As more components are added, the model portrays an equal representation of all three groups. This demonstrates that the proposed model is more effective at capturing the common variability and fostering a broader understanding of the data.

Fig. 9 displays the scores of the first two components for both models. In the PARAFAC case, it is evident that the scores corresponding to group 1 exhibit greater separation than those of the remaining groups, in agreement with the above findings. Moreover, the direction of the

and and Part of the

Table 2

Cumulative % of Variance Explained by Groups: Case study X three-way matrix.

PARAFAC Algorithm

The month of the official										
	Comp 1	Comp 2	Comp 3	Comp 4	Comp 5	Comp 6	Comp 7	Comp 8	Comp 9	Comp 10
Group 1	39.00	39.25	39.18	53.57	53.76	65.48	66.00	66.22	67.93	71.53
Group2	26.85	26.90	27.49	33.15	35.74	39.58	42.55	43.15	48.79	52.57
Group3	20.66	21.42	21.70	27.17	44.48	46.98	47.69	48.47	52.79	55.54
Proposed Algorithm										
	Comp 1	Comp 2	Comp 3	Comp 4	Comp 5	Comp 6	Comp 7	Comp 8	Comp 9	Comp 10
Group 1	32.90	33.10	33.30	38.27	42.62	46.53	50.30	56.59	59.70	58.36
Group2	29.71	29.69	30.18	36.88	41.08	39.33	47.20	48.31	53.68	53.06
Group3	23.38	24.21	24.65	30.06	36.84	36.30	44.05	48.76	52.84	55.59



Fig. 8. Case study: Ratios of the explained variance for each group with respect to group 3.



Fig. 9. Case study: Three-way array X, first two dimensions scores representation with a 99% confidence interval ellipse.

scores of group 2 is different from the other two groups. Conversely, in the proposed model scenario, the scores of all three groups exhibit significant overlap within a common region. Notably, the groups also demonstrate shared patterns direction of variability. This outcome is consistent with the objective of the proposed model, which aims to recover loadings that express the shared patterns and variability across all groups.

The biplot of the loadings matrices **B** for both models are illustrated in Fig. 10. In the PARAFAC model, the variables V2, V4, V8, V11, V44 and V45 have the highest model coefficients. The proposed model showcases similar high coefficients for the variables V11, V44 and V45, albeit with slightly reduced magnitudes. Moreover, other variables such as V1, V2, V8, V9, and V10 are shown to have high model coefficients in the proposed model. However, in the PARAFAC model, the coefficients for these variables are lower in comparison. The disparity in the estimates of the model coefficients can likely be attributed to the unequal representation of the variability of the three-way array by the two models. As shown in the above findings, the PARAFAC model primarily captures the variability specific to the first group rather than the entirety of the dataset. It implies that the emphasis placed on the first group by the PARAFAC model may result in model coefficients representing this group rather than the entire three-way array. This difference between the two models may result in interpretations and conclusions that are different.

Following the analysis of the outcomes and consultation with expert engineers in the field, a consensus was reached regarding the results obtained from the proposed model. The variables V1 and V2 are associated with the pressure applied to each micro-layer, which determines its fusion with the main component. This step is highly important in the production process, as it significantly impacts the numerous mechanical properties of the final product. The variables V8 and V9are linked to the bulk temperature, representing the temperature of the main component of all the fused layers. This temperature plays a crucial role in material fusion and adhesion of the individual microlayers, so the fusion between all the micro-layers influences the printing process overall success and contributes to several quality aspects of the end products. Moreover, the variable V10 is linked to the creation of each individual micro-layer, underscoring its significance in the overall



Fig. 10. Case study: Biplot of the first two components of the loadings matrix B of the PARAFAC and the proposed model.

process. Lastly, both models highlighted the importance of the variables *V*44 and *V*45, which pertain to the sensors in the cooling step. Cooling temperatures are of noteworthy influence on warpage, a factor directly affecting various qualities of the final product. These findings reinforce the understanding that careful management of cooling temperatures is imperative for ensuring optimal product outcomes across multiple quality dimensions.

Discussion

Across all four data sets, we demonstrated that the proposed approach was able to provide a model that is representative of all the groups present in the data. This was supported by similar levels of variance explained across all groups compared to the variance discrepancies in the PARAFAC model. This is indicative of a model that is representative of the variability of the entire data set and all groups. Despite the fact that one or more groups may exhibit larger variability, the suggested model was able to reduce this variability and produce model coefficients that are common to all groups. The two simulation strategies depicted two ways of inducing a multi-group structure in the data and in both cases, the proposed model was able to regulate the dominance of the group with the highest variance. In the case study, we presented a real-world application in additive manufacturing, aiming to gain valuable process insights and understanding. To achieve this aim, it is necessary to consider a model that accurately represents all groups in the data set. The proposed model showcases a percentage of explained variance similar for all groups, indicating its ability to effectively capture the common variability of the three-way array and its multiple groups. This stands in contrast to the PARAFAC model, which exhibits a greater representation of the first group and its variability. The strength of the proposed model lies in its capability to mitigate the higher variability observed in the first group and successfully recover a common structure that encompasses all groups. An in-depth analysis of the loadings matrix **B** provides further insights for process understanding. The results of the PARAFAC model indicate a lack of homogeneity in variability among the groups, leading to different model coefficients. This aspect holds crucial implications for a deeper understanding of the underlying process, as in the PARAFAC case, the coefficients of the variables are more representative of the first group rather than the entire data set. In this context, the proposed model proves to be more suitable and advantageous for three-way arrays that present a multi-group structure. Furthermore, we recommend extending the analysis to include a meticulous examination of residuals, the threeway matrix E, when conducting further studies on differences between various groups. As the model represents the common structure shared

among all groups, analyzing the residuals can shed light on the different variations a specific group represents. However, it is essential to note that such an analysis lies beyond the scope of the proposed model in this paper. Instead, it serves as a natural direction for future studies, warranting careful consideration and examination on a case-by-case basis.

4. Conclusions

The extension of the PARAFAC model proposed in this paper makes it possible to analyze complex data that present a multi-grouping structure. This transfer of methodology from the PARAFAC model to the multi-group settings was defined by adding constraints to the objective function to consider the group structure for the model parameters. We provided the corresponding extension of the ALS algorithm to solve the proposed model. This algorithm was set to be also used in the absence of multi-group data as its initialization corresponds to the PARAFAC model.

Three simulation studies and a real case were used to illustrate the capability of the extended PARAFAC method to render a model that explains the common variability of all the groups. In all applications, using the model without a multi-group structure when the groups have large differences resulted in a model being representative of the most dominant group (i.e. the group with the largest variance). In the real case study in additive manufacturing, we illustrated the impact of considering the grouping structure to analyze the importance of process variables in this type of data. By employing the extended PARAFAC method, we achieved a more comprehensive understanding of the entire manufacturing process, avoiding a biased focus on individual groups. This perspective facilitated a more profound analysis of the data and enabled us to extract valuable insights that would have been overlooked in a model dominated by a single group.

Our results demonstrate that the proposed model performs well in capturing the common and explanatory loadings of the data, even when the variances of the groups are unequal, resulting in a more robust and reliable model. The proposed model effectively represents the entire dataset, and this is of utmost importance for gaining valuable insights and understanding of the process. When a model is heavily influenced by one dominant group, the loadings tend to explain more of that particular group's characteristics, neglecting the broader dynamics of the entire process. The proposed method offers a solution to this challenge, ensuring a more balanced representation of the data and fostering deeper process understanding across all groups.

CRediT authorship contribution statement

Marta Rotari: Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. Valeria Fonseca Diaz: Writing – original draft, Formal analysis, Conceptualization. Bart De Ketelaere: Writing – review & editing, Writing – original draft, Supervision, Conceptualization. Murat Kulahci: Supervision, Resources, Project administration, Conceptualization.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

Data will be made available on request.

References

- M. Spooner, D. Kold, M. Kulahci, Selecting local constraint for alignment of batch process data with dynamic time warping, Chemometr. Intell. Lab. Syst. 167 (2017) 161–170.
- [2] D. Louwerse, A.K. Smilde, Multivariate statistical process control of batch processes based on three-way models, Chem. Eng. Sci. 55 (7) (2000) 1225–1235.
- [3] M. Ryckewaert, G. Chaix, D. Héran, A. Zgouz, R. Bendoula, Evaluation of a combination of nir micro-spectrometers to predict chemical properties of sugarcane forage using a multi-block approach, Biosyst. Eng. 217 (2022) 18–25.
- [4] M. Dyrby, M. Petersen, A.K. Whittaker, L. Lambert, L. Nørgaard, R. Bro, S.B. Engelsen, Analysis of lipoproteins using 2d diffusion-edited nmr spectroscopy and multi-way chemometrics, Anal. Chim. Acta 531 (2) (2005) 209–216.
- [5] C.M. Rubingh, S. Bijlsma, R.H. Jellema, K.M. Overkamp, M.J. van der Werf, A.K. Smilde, Analyzing longitudinal microbial metabolomics data, J. Proteome Res. 8 (9) (2009) 4319–4327.
- [6] S. Wold, P. Geladi, K. Esbensen, J. Öhman, Multi-way principal components-and pls-analysis, J. Chemom. 1 (1) (1987) 41–56.
- [7] R. Bro, Multiway calibration. multilinear pls, J. Chemom. 10 (1) (1996) 47-61.
- [8] R. Harshman, Foundations of the parafac procedure: Models and conditions for an "explanatory" multi-modal factor analysis, UCLA Work. Pap. Phonetics 16 (1970).
- [9] L.R. Tucker, et al., The extension of factor analysis to three-dimensional matrices, Contrib. Math. Psychol. 110119 (1964).
- [10] R. Vitale, O.E. de Noord, J.A. Westerhuis, A.K. Smilde, A. Ferrer, Divide et impera: How disentangling common and distinctive variability in multiset data analysis can aid industrial process troubleshooting and understanding, J. Chemom. 35 (2) (2021) e3266.
- [11] Y. Ben-Shabat, M. Lindenbaum, A. Fischer, 3Dmfv: Three-dimensional point cloud classification in real-time using convolutional neural networks, IEEE Robot. Autom. Lett. 3 (4) (2018) 3145–3152.

- [12] C. Ma, Y. Guo, Y. Lei, W. An, Binary volumetric convolutional neural networks for 3-d object recognition, IEEE Trans. Instrum. Meas. 68 (1) (2018) 38–48.
- [13] J.D. Carroll, J.-J. Chang, Analysis of individual differences in multidimensional scaling via an n-way generalization of "eckart-young" decomposition, Psychometrika 35 (3) (1970) 283–319.
- [14] L.R. Tucker, Some mathematical notes on three-mode factor analysis, Psychometrika 31 (3) (1966) 279–311.
- [15] R. Bro, Parafac. tutorial and applications, Chemometr. Intell. Lab. Syst. 38 (2) (1997) 149–171.
- [16] J.A. Westerhuis, T. Kourti, J.F. MacGregor, Analysis of multiblock and hierarchical pca and pls models, J. Chemom. J. Chemom. Soc. 12 (5) (1998) 301–321.
- [17] D. Chessel, M. Hanafi, Analyses de la co-inertie de k nuages de points, Rev. Stat. Appl. 44 (2) (1996) 35–60.
- [18] S. Wold, N. Kettaneh, K. Tjessem, Hierarchical multiblock pls and pc models for easier model interpretation and as an alternative to variable selection, J. Chemom. 10 (5–6) (1996) 463–482.
- [19] A. Eslami, A. Kohler, M. El Qannari, S. Bougeard, General overview of methods of analysis of multi-group datasets, in: HDSDA, 2011, pp. 108–123.
- [20] S. Legleye, A. Eslami, S. Bougeard, Assessing the structure of the cast (cannabis abuse screening test) in 13 European countries using multigroup analyses, Int. J. Methods Psychiatr. Res. 26 (1) (2017) e1552.
- [21] A. Tenenhaus, M. Tenenhaus, Regularized generalized canonical correlation analysis for multiblock or multigroup data analysis, Eur. J. Oper. Res. 238 (2) (2014) 391–403.
- [22] M. Hanafi, E.M. Qannari, B. Jaillais, Multi-Block and Three-Way Data Analysis, 2019, http://dx.doi.org/10.1016/B978-0-12-409547-2.14717-1.
- [23] B.N. Flury, Common principal components in k groups, J. Amer. Statist. Assoc. 79 (388) (1984) 892–898.
- [24] W. Krzanowski, Principal component analysis in the presence of group structure, J. R. Stat. Soc. Ser. C. Appl. Stat. 33 (2) (1984) 164–168.
- [25] J. Kallus, P. Johansson, S. Nelander, R. Jörnsten, Mm-pca: integrative analysis of multi-group and multi-view data, 2019, arXiv preprint arXiv:1911.04927.
- [26] A. Eslami, E.M. Qannari, A. Kohler, S. Bougeard, Algorithms for multi-group pls, J. Chemom. 28 (3) (2014) 192–201.
- [27] A. Eslami, E. Qannari, A. Kohler, S. Bougeard, Multivariate analysis of multiblock and multigroup data, Chemometr. Intell. Lab. Syst. 133 (2014) 63–69.
- [28] A.C. Olivieri, G.M. Escandar, H.C. Goicoechea, A.M. de la Peña, Fundamentals and Analytical Applications of Multiway Calibration, Elsevier, 2015.
- [29] T.G. Kolda, B.W. Bader, Tensor decompositions and applications, SIAM Rev. 51 (3) (2009) 455–500.
- [30] J. h. Jiang, H. l. Wu, Y. Li, R. q. Yu, Three-way data resolution by alternating slice-wise diagonalization (asd) method, J. Chemom. J. Chemom. Soc. 14 (1) (2000) 15–36.
- [31] P. Paatero, A weighted non-negative least squares algorithm for three-way 'parafac'factor analysis, Chemometr. Intell. Lab. Syst. 38 (2) (1997) 223–242.
- [32] N. Faber, L. Buydens, G. Kateman, Generalized rank annihilation method. i: Derivation of eigenvalue problems, J. Chemom. 8 (2) (1994) 147–154.
- [33] E. Sanchez, B.R. Kowalski, Tensorial resolution: a direct trilinear decomposition, J. Chemom. 4 (1) (1990) 29–45.
- [34] H.-P. Yeh, K. Meinert, M. Bayat, J. Hattel, Part-scale thermo-mechanical modelling for the transfusion module in the selective thermoplastic electrophotographic process, in: WCCM-APCOM 2022, in: Manufacturing and Materials Processing, vol. 1000, 2022.
- [35] Y. Wu, K. He, Group normalization, in: Proceedings of the European Conference on Computer Vision, ECCV, 2018, pp. 3–19.