



HAL
open science

Estimation of high dimensional probability density functions with low rank-tensors models: application to flow cytometry

Philippe Flores

► **To cite this version:**

Philippe Flores. Estimation of high dimensional probability density functions with low rank-tensors models: application to flow cytometry. Automatic. Université de Lorraine, 2024. English. NNT : 2024LORR0021 . tel-04625735

HAL Id: tel-04625735

<https://hal.univ-lorraine.fr/tel-04625735>

Submitted on 26 Jun 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



**UNIVERSITÉ
DE LORRAINE**

**BIBLIOTHÈQUES
UNIVERSITAIRES**

AVERTISSEMENT

Ce document est le fruit d'un long travail approuvé par le jury de soutenance et mis à disposition de l'ensemble de la communauté universitaire élargie.

Il est soumis à la propriété intellectuelle de l'auteur. Ceci implique une obligation de citation et de référencement lors de l'utilisation de ce document.

D'autre part, toute contrefaçon, plagiat, reproduction illicite encourt une poursuite pénale.

Contact bibliothèque : ddoc-theses-contact@univ-lorraine.fr
(Cette adresse ne permet pas de contacter les auteurs)

LIENS

Code de la Propriété Intellectuelle. articles L 122. 4

Code de la Propriété Intellectuelle. articles L 335.2- L 335.10

http://www.cfcopies.com/V2/leg/leg_droi.php

<http://www.culture.gouv.fr/culture/infos-pratiques/droits/protection.htm>

Estimation of high dimensional probability density functions with low-rank tensor models: application to flow cytometry.

*Estimation de densités de probabilité en grandes dimensions par
modèles tensoriels de rang faible : application à la cytométrie en
flux.*

THÈSE

présentée et soutenue publiquement le 16 avril 2024

pour l'obtention du

Doctorat de l'Université de Lorraine

(Mention automatique, traitement du signal et génie informatique)

par

Philippe Flores

Composition du jury

<i>Rapporteurs :</i>	Olivier MICHEL	<i>Professeur des Universités, Université Grenoble Alpes</i>
	Vicente ZARZOSO	<i>Professeur des Universités, Université Côte d'Azur</i>
<i>Examinatrices :</i>	Marianne CLAUSEL	<i>Professeure des Universités, Université de Lorraine</i>
	Mariya ISHTEVA	<i>Professeure Assistante, KU Leuven</i>
<i>Directeurs de thèse :</i>	David BRIE	<i>Professeur des Universités, Université de Lorraine</i>
	Konstantin USEVICH	<i>Chargé de recherche, CNRS, Université de Lorraine</i>
<i>Présidente du jury :</i>	Marianne CLAUSEL	

Mis en page avec la classe thesul.

Acknowledgements

In these three years of thesis, I found myself with mixed feelings. First, this experience went extremely fast. In the meantime, I also look back at it thinking a lot of things happened during that time. It started with Covid, I moved 3 times in the process, and so much more. . . With those words, I would like to thank the people that were with me throughout my journey.

First and foremost, I would like to give my full gratitude to my supervisors, David and Konstantin. I know I was not the easiest PhD student and for that I want to thank you for bringing me back on track, giving me motivation, believing in me when I did not. Thank you for your scientific advices, always on point. I enjoyed our discussions with the most interests.

I would also like to thank Olivier Michel, Vicente Zarzoso, Marianne Clausel and Mariya Ishteva for accepting to be part of this PhD committee. Special thanks to the members of my thesis committee Pierre Comon and Marianne. I would also like to give a special thank to Olivier and Vicente for their understanding and thoughtful discussions along this manuscript writing.

During the realization of this project, I had the occasion and privilege to work in collaboration with a lot of people. I would like to give a big thank you to the biologists of CRAN, Guillaume Harlé and Stéphanie Grandemange, for their help on flow cytometry and everything. On the medical doctors part, I would like to thank Maud D'Aveni and Anne-Béatrice Notarantonio for their implication in the project. Special thanks to Anne-Béatrice who helped me tremendously on the flow cytometry part. I learnt a lot from you and I hope we will continue this collaboration. A warm thank you to Joseph Chege, Martin Haardt and Arie Yeredor for our collaboration and helpful discussions.

For the people of the 4th floor at CRAN, I would like to thank everyone of you for making me feel like I belong there. To Saulo, Narech, Faustine, Hadi, Sebastian, Jean-Christophe, Ricardo, thank you for your support along the way. Special thanks to Clémence and Julien, with whom I share most of my time. I have only good memories with you and I'm happy to count you as my friends.

To the members of the R&T department at IUT Nancy Brabois goes my gratitude for their trust and welcoming atmosphere. Special thank to you David for giving me the opportunity to gain that much experience as a teaching assistant.

I am also grateful for all the discussions I couldn't have had without the thesis: conferences, seminars or workshops. I cannot be exhaustive here, but all my gratitude comes to those who exchanged with me along the way.

On a more personal note, hence in french. Premièrement, je voudrais remercier les amis qui sont toujours autour de moi, même si je leur ai fait défaut à quelques occasions: Dylan, Maxime, Sandra, Éléonore, tous mes amis de prépa, etc. À ma famille que j'aime, je voudrais vous remercier pour votre soutien et tout ce que vous m'avez apporté. À mes soeurs, je sais que je peux compter sur vous n'importe quand et je sais qu'on restera toujours "Les Quatre". À mes parents, pleins de merci pour être toujours à l'écoute alors que je suis si loin de vous. Pour Nathan et Louis, j'aimerais vous dire à quel point je suis fier de vous, que je sais à quel point il m'a été difficile de ne pas vous voir grandir autant que je le voulais. Et enfin, à Elsa, je ne peux dire avec des mots à quel point tu m'as aidé pendant ces 3 années. Pour ton soutien, ta motivation, tes conseils, ta gentillesse, j'ai de la chance que tu fasses partie de ma vie.

Contents

List of Figures	vii
List of Tables	xi
List of notations	1
List of acronyms	3
<hr/>	
General introduction	5
Chapter 1 Flow cytometry	9
1.1 General principle and applications	10
1.1.1 Geometrical data acquisition	12
1.1.2 Fluorescence data acquisition	13
1.2 Flow cytometry data pre-processing	16
1.2.1 Flow cytometry dataset example	16
1.2.2 Compensation	18
1.2.3 Cleaning	20
1.2.4 Non-linear transformation	23
1.3 Flow cytometry data analysis	24
1.3.1 Manual data analysis: Gating	24
1.3.2 Computational flow cytometry	27
1.4 Conclusion of Chapter 1	31
Chapter 2 Tensor preliminaries	33
2.1 Basic definitions	34
2.1.1 Matrix operations	34
2.1.2 Tensors and basic definitions	35
2.2 Canonical Polyadic decomposition	37
2.2.1 Definition of the CP decomposition	37

2.2.2	Uniqueness of the CP decomposition	38
2.2.3	Generic uniqueness	40
2.3	Identifiability of polynomial models	41
2.3.1	Polynomial mapping and uniqueness	42
2.3.2	Uniqueness and recoverability of additive models	43
2.4	Conclusion of Chapter 2	45
Chapter 3 Histograms for probability density function estimation		47
3.1	Density estimation	48
3.1.1	Problem statement	48
3.1.2	Parametric approaches	49
3.1.3	Non-parametric approaches	50
3.2	Histograms	51
3.2.1	Definitions	51
3.2.2	Curse of dimensionality: choosing the number of bins	55
3.3	Naive Bayes model and CP decomposition	57
3.3.1	Graphical models	57
3.3.2	Naive Bayes model	59
3.3.3	Link between naive Bayes model and tensor decompositions	61
3.3.4	Curse of dimensionality: estimating the factors of the NBM	63
3.4	Fully coupled tensor factorization of 3D marginals	64
3.4.1	Marginalization of naive Bayes model	64
3.4.2	Fully coupled tensor factorization	66
3.4.3	Curse of dimensionality: number of marginals	67
3.5	Conclusion of Chapter 3	68
Chapter 4 Partially coupled tensor factorization for probability mass function estimation		71
4.1	Partially coupled tensor factorization	72
4.1.1	Principles	72
4.1.2	Optimization	73
4.1.3	An example in the case of 4 variables	74
4.2	Coupling strategies	77
4.2.1	Coupling strategies as hypergraphs	78
4.2.2	Examples of coupling strategies	80
4.2.3	An algorithm for generating balanced couplings	87
4.3	Numerical experiments	90
4.3.1	Random strategies and performance regarding the number of triplets	91
4.3.2	Comparison between random and balanced couplings.	91

4.3.3	Comparison with a KLD-based method	94
4.4	Conclusion of Chapter 4	97
Chapter 5 Identifiability of coupled tensor models		99
5.1	Recoverability results	100
5.1.1	Reduced re-parametrization of factor matrices	100
5.1.2	Jacobian structure and algorithm for recoverability bound search	102
5.1.3	Recoverability necessary condition and number of degrees of freedom	104
5.1.4	Recoverability results for fully coupled tensor factorization	105
5.1.5	Recoverability results for randomly selected triplets	105
5.1.6	Analysis of defective cases	106
5.1.7	Recoverability results for balanced couplings	110
5.2	Identifiability of Cartesian product coupling	111
5.2.1	Re-parametrization of the coupled model and identifiability results	113
5.2.2	Identifiability of the equivalent tensor model	114
5.2.3	An algebraic proof of identifiability	115
5.3	Application of Theorem 5.2.3	119
5.3.1	Identifiability result in the general case of variable partition coupling	119
5.3.2	Identifiability result for an even partition of variables	120
5.3.3	Identifiability result for the fully coupled case	121
5.4	Conclusion of Chapter 5	121
Chapter 6 <i>CTFlowHD</i> : a tensor-based method for flow cytometry data analysis		123
6.1	Presentation of the method <i>CTFlowHD</i>	124
6.1.1	Naive Bayes model and flow cytometry data	124
6.1.2	<i>CTFlowHD</i> workflow of analysis	125
6.1.3	Implementation	127
6.1.4	Application to a controlled flow cytometry dataset	128
6.2	Clustering and visualizations	131
6.2.1	Visualizations based on hierarchical clustering	132
6.2.2	Marginal visualizations	134
6.2.3	<i>CTFlowHD</i> : a versatile visualization tool	136
6.3	Application to real flow cytometry data	140
6.3.1	Application to 8-marker datasets	141
6.3.2	Application to 24-marker datasets	146
6.4	Conclusion of Chapter 6	148
General Conclusion		149

Contents

Appendix A Brochure of BD for FACSymphony flow cytometers	153
Appendix B PCTF3D applied to supervised classification on the MNIST dataset	163
Appendix C Presentation of PCTF3D in the context of flow cytometry	171
Bibliography	179

List of Figures

1.1	Flow cytometer principle.	11
1.2	Fluorescence flow cytometry principle.	12
1.3	Principle of FSC value acquisition.	13
1.4	Forward and side scatter for a cell [2].	14
1.5	Bivariate plot (FSC-SSC) human blood cells [104]. Area values are plotted here: FSC-A and SSC-A.	14
1.6	Screenshot featuring a table containing FCM variables information.	18
1.7	FCM data represented with 2-dimensional scatter plots.	19
1.8	Overlapping of fluorophore emission bandwidths.	20
1.9	Effect of compensation on raw FCM data.	21
1.10	Viable cell detection gate with SSC and propidium iodide bivariate plots [90].	22
1.11	Singlets cell detection gate with FSC-A and FSC-H bivariate plots [115].	22
1.12	Application of a "Logicle" scale to the controlled dataset.	25
1.13	Example of manual gating sequence.	26
1.14	Gating applied to the controlled dataset.	26
1.15	Example of SPADE visualization.	28
1.16	Example of viSNE visualization.	29
2.1	An order-3 tensor with its 3 sets of fibers.	36
3.1	Example of a Gaussian Mixture Model in a two variable case.	50
3.2	A theoretical histogram (Left plot) and a empirical histogram (Right plot) as an estimation of the scaled PDF.	53
3.3	Example of graphical models.	59
3.4	Naive Bayes model as a graphical model.	60

3.5	Link between CPD and NBM.	63
3.6	Example of marginalization of the NBM.	65
4.1	Representation of couplings for 6 variables.	78
4.2	Theoretical and empirical distributions of the degree sequence.	86
4.3	Empirical distributions of the degree sequence for random couplings.	86
4.4	Error regarding the number of triplets for different coupling strategies.	92
4.5	Error on 3D marginals for balanced and random coupling strategies.	93
4.6	Results for the sensitivity experiment.	95
4.7	Comparison of CTF3D, PCTF3D and SQUAREM-PMF in terms of runtime and FMS.	96
5.1	Recoverability bounds for different values of M and I in the case of the CTF3D.	106
5.2	Recoverability bounds in the case of random partial coupling.	107
5.3	Recoverability bounds in the case of balanced partial coupling.	110
5.4	Structure of \mathcal{Y} as a concatenation of 3D marginals.	112
5.5	The matrix \mathbf{Q}_1 . The matrix $\mathcal{J}_{\mathcal{P}} \in \mathbb{R}^{I \times (I-1)}$ is the Jacobian matrix of the mapping for truncated factors (see Equation (5.9)).	116
5.6	Visualization of the constraints which form a polyhedron in a 3D space.	118
5.7	Identifiability bounds in the fully coupled case along with recoverability results.	122
6.1	Naive Bayes model for flow cytometry data analysis.	125
6.2	<i>CTFlowHD</i> workflow that contains the algorithm PCTF3D.	126
6.3	Manual gating providing ground truth for the controlled experiment.	129
6.4	<i>CTFlowHD</i> output for a controlled experiment before the visualization step.	130
6.5	<i>CTFlowHD</i> output for a controlled experiment after a hierarchical clustering applied to rank-one terms (single linkage).	131
6.6	<i>CTFlowHD</i> output for a controlled experiment after a hierarchical clustering applied to rank-one terms (complete linkage).	133
6.7	Bivariate plot CFSE/CTV for the controlled experiment combined with <i>CTFlowHD</i> 2D marginals.	134
6.8	<i>CTFlowHD</i> 2D marginals combined with hierarchical clustering (complete linkage).	135
6.9	Controlled experiment visualized with a K-means clustering method.	136
6.10	<i>CTFlowHD</i> applied with a t-SNE visualization.	137

6.11	t-SNE applied to raw data from the controlled experiment.	138
6.12	SPADE applied to raw data from the controlled experiment.	138
6.13	<i>CTFlowHD</i> applied to the controlled experiment with a SPADE visualization. . .	139
6.14	Hematopoiesis tree.	141
6.15	<i>CTFlowHD</i> applied to a graft sample dataset with 8 markers (full coupling). . .	142
6.16	<i>CTFlowHD</i> applied to a graft sample dataset with 8 markers (partial coupling). .	144
6.17	Cell population signatures for the fully coupled strategy.	145
6.18	Cell population signatures for the partially coupled strategy.	145
6.19	<i>CTFlowHD</i> applied to a 24-marker FCM dataset.	147
B.1	Examples of digit images from the MNIST dataset.	164
B.2	Cropping and resizing of MNIST images.	164
B.3	Quantization of MNIST images.	165
B.4	Representation of triplets as a 3-neighbor mapping.	167
B.5	Feature images for the different likelihood estimators.	169

List of Figures

List of Tables

1.1	FCM variable information for the controlled FCM experiment.	16
1.2	Cell populations used in the controlled FCM experiment with their marker expressions.	17
3.1	Evolution of the required number of samples with the optimal I^* and the order M	57
3.2	Number of marginals exhibiting the third level of curse of dimensionality.	68
4.1	Examples of couplings in the case of 6 variables (see Figure 4.1).	81
4.2	Coupled strategies properties and associated complexities.	81
4.3	Examples of '+2' couplings associated with their incidence matrix and degree sequence.	83
4.4	Examples of '+1' couplings with their associated incidence matrix and degree sequence.	85
4.5	Coupling strategies for $M = 10$	92
5.1	Evolution over M of R_{\max} for most favorable coupling case with $d_1 = 1$	108
5.2	Evolution over M of R_{\max} for most favorable coupling case with $d_1 = d_2 = 1$	109
5.3	Number of trials regarding the number of triplets considered in \mathcal{T} for the balanced couplings experiment.	110
6.1	Properties of the 3 populations used in the controlled experiment.	128
6.2	Analysis run times for the controlled experiment.	140
6.3	Identification of cell populations obtained with <i>CTFlowHD</i> with hierarchical clustering.	143
B.1	Number of samples for each digit in the labelled MNIST dataset.	163
B.2	Accuracy of MNIST classification using PCTF3D and a Bayes classifier.	169

List of Tables

List of notations

\mathcal{I}	sets (calligraphic uppercase)
$\mathbb{1}_{\mathcal{I}}$	indicator function of the set \mathcal{I}
\mathbb{N}	set of non-negative integers
\mathbb{R}	set of real numbers
\mathbb{R}_+	set of real non-negative numbers
m, M	scalars (lowercase or uppercase)
$[[1, M]]$	set of integers $\{1, 2, \dots, M\}$
\mathbf{v}	vectors (bold lowercase)
v_i	i -th entry of \mathbf{v}
$\mathbb{1}_I$	vector of all-ones of size I
\mathbf{A}	matrices (bold uppercase)
a_{ij}	i -th row and j -th column entry of matrix \mathbf{A}
$\mathbf{A}_{:j}$	j -th column of \mathbf{A}
$\mathbf{A}_{i:}$	i -th row of \mathbf{A}
\mathbf{A}^{-1}	matrix inverse
\mathbf{A}^T	matrix transpose
\mathbf{A}^\dagger	matrix pseudo-inverse
$\text{Diag}(\mathbf{d})$	diagonal matrix with diagonal \mathbf{d}
$\text{diag}(\mathbf{A})$	vector of diagonal elements of \mathbf{A}
$\text{rank}(\mathbf{A})$	rank of \mathbf{A}
$\kappa(\mathbf{A})$	Kruskal rank of \mathbf{A}

List of notations

\mathbf{I}_M	Identity matrix of size $M \times M$
$\ \cdot\ _2$	2-norm
\mathcal{X}	high-order tensors (bold calligraphic uppercase)
$x_{i_1 \dots i_M}$	(i_1, \dots, i_M) -th entry of \mathcal{X}
$\mathbf{X}_{(m)}$	mode- m unfolding of \mathcal{X}
$\text{vec}(\mathcal{X})$	vectorization of \mathcal{X}
○	outer product
⊗	Kronecker product
⊙	Khatri-Rao product
*	Hadamard product
$\ \cdot\ _F$	Frobenius norm
\mathcal{H}	theoretical histogram tensor
$\tilde{\mathcal{H}}$	empirical histogram tensor
$\hat{\mathcal{H}}$	estimated histogram tensor
R	CP decomposition rank
$\mathbf{A}^{(m)}$	m -th factor matrix of the CPD of \mathcal{H}
$\hat{\mathbf{A}}^{(m)}$	estimated m -th factor matrix
$\mathbf{a}_r^{(m)}$	(m, r) -th factor of the CPD of \mathcal{H}
$\hat{\mathbf{a}}_r^{(m)}$	estimated (m, r) -th factor
$\{j, k, \ell\}$	triplet of variables
\mathcal{T}	set of triplets of variables
\mathcal{T}_m	set of triplets of variables that contains the m -th variable
\mathbf{V}	incidence matrix of a hypergraph
\mathbf{d}	degree sequence of a hypergraph
P	number of doublets of a hypergraph
T	number of triplets of a hypergraph
$\mathcal{H}^{(j k \ell)}$	marginal histogram tensor for variables in $\{j, k, \ell\}$
$\tilde{\mathcal{H}}^{(j k \ell)}$	empirical estimation of $\mathcal{H}^{(j k \ell)}$

List of acronyms

ADMM	Alternating Direction Method of Multipliers
AO	Alternating Optimization
ALL	Acute Lymphoblastic Leukemia
ALL	Acute Myeloid Leukemia
BW	Black and White
CD	Cluster of Differentiation
CFC	Computational Flow Cytometry
CP	Canonical Polyadic
CPD	Canonical Polyadic Decomposition
EM	Expectation-Maximization
FACS	Fluorescence Active Cell Sorting
FCM	Flow CytoMetry
FCS	Forward SCatter
FCS-A	Forward SCatter - Area
FCS-H	Forward SCatter - Height
CTF3D	Coupled Tensor Factorization of 3D marginals
FMS	Factor Match Score
GMM	Gaussian Mixture Model
GvHD	Graft versus Host Disease

List of acronyms

GvT	Graft versus Tumor
KLD	Kullback-Leibler Divergence
MAP principle	Maximal A Posteriori principle
MISE	Mean Integrated Squared Error
MSE	Mean Squared Error
NBM	Naive Bayes Model
PCTF3D	Partially Coupled Tensor Factorization of 3D marginals
PDF	Probability Density Function
PMF	Probability Mass Function
SSC	Side SCatter
SSC-A	Side SCatter - Area
SSC-H	Side SCatter - Height

General introduction

This work has been conducted at CRAN – previously Centre de Recherche en Automatique de Nancy – (CNRS UMR 7039) in the BioSiS department (Biologie, Signaux et Systèmes en Cancérologie et Neurosciences). It is part of a collaboration between biologists from CRAN, medical doctors from CHRU of Nancy and signal processing researchers from CRAN. The objective of this collaboration was to develop an automated unsupervised flow cytometry data analysis method.

Problem statement

Flow cytometry measures biological properties of a flow of individual cells. It is widely used in biological studies with both clinical and academic research applications. For example, a clinical application of flow cytometry is the diagnosis of auto-immune diseases such as leukemia. It also can be used to study the immune system or to apprehend why certain leukemia patients suffer from graft relapses while others do not. During the past decades, flow cytometry has permitted to measure more and more biological properties (up to 50 parameters) for more and more cells (up to millions) in a reduced time. Meanwhile, as manual data analysis proved its limitations, the interest for automated flow cytometry data analysis gave rise to the new field of computational flow cytometry. Most methods developed in this field are adaptations of already existing clustering and classification methods.

Biological cell properties measured by a flow cytometer are of two types. Geometrical (size or granularity of a cell) properties are mainly used in a pre-processing step to select certain types of cells (for example lymphocytes T). Fluorescence properties are obtained by dyeing cell genes of interests and measuring afterwards the fluorescence emitted by each cell. In terms of biological applications, the problem of flow cytometry data analysis includes three main tasks:

- Identifying a cell population of interest inside a pool of samples;

- Quantifying the proportion of a cell population of interest;
- Exploring a dataset to find and characterize new cell populations.

From a signal processing point of view, the problem of flow cytometry data analysis comes down to clustering and/or classification problems. In this work, this problem is addressed as a non-parametric (histogram) joint density estimation. This is a fundamental problem in data processing which has motivated numerous works over years. Indeed, having an estimate of the joint density is a kind of Grail in data processing since, for example, classification and clustering tasks can be made using the maximum *a posteriori* principle. However, density estimation is often considered impossible in practice due to the curse of dimensionality which states that the complexity of a problem increases exponentially with its number of dimensions. To give some figures, to estimate a histogram in dimension 10 requires more than a billion points, which is out of reach in most practical situations, including flow cytometry. In fact, in a wide range of situations, only 3-D and 4-D histograms can be estimated.

To circumvent the curse of dimensionality, we developed a high dimension joint density estimation approach which includes 3 main ingredients: (i) the use of a limited number of low dimension (3-D) histograms, (ii) a probabilistic model whose complexity remains linear with the number of dimensions, (iii) a tensor-based algorithmic framework with theoretical guarantees ensuring reproducible results. This algorithmic framework yields low-dimensional probabilistic features that can further be used in classification/clustering tasks.

Outline of the manuscript

The first three chapters of this manuscript present the flow cytometry data analysis problem and the two main mathematical tools that are used in this work, namely tensors and histograms for density estimation. The first chapter introduces the basic principles of flow cytometry data acquisition. Also, some representative methods from computational flow cytometry are presented together with remaining challenges that will be addressed in this work. The second chapter introduces the tensorial background used in this thesis. In particular, the canonical polyadic decomposition and associated classical uniqueness results are recalled. Then, a special emphasis is given to the less classical results relating to the uniqueness of polynomial additive models. They will be extensively used in the chapter dedicated to the identifiability of coupled tensor models. Finally, in chapter 3, the problem of probability density function estimation is introduced. The focus is put on histograms that are nothing but probability tensors whose estimation is hampered

by the curse of dimensionality. To tackle the curse of dimensionality, the joint density is assumed to follow the naive Bayes model. A key point is that it can be interpreted as a constrained canonical polyadic decomposition; two representative tensor-based methods for density estimation are presented and discussed.

The last three chapters present the main contributions of this thesis. Chapter 4 presents a reduced complexity approach which makes use of a limited (and controlled) number of 3-D histograms through a coupled tensor factorization algorithm. This method is termed as PCTF3D for Partially Coupled Tensor Factorization of 3D marginals. Choosing which histograms are coupled is crucial. Different coupling strategies are proposed and studied with the framework of hypergraphs. In particular, an algorithm for balanced coupling generation that ensures that variables are represented evenly in terms of occurrences, is proposed. This chapter concludes with numerical experiments aiming at evaluating and comparing the performances of PCTF3D to other state-of-the-art methods.

In chapter 5, uniqueness of the coupled model used in PCTF3D is studied in terms of model recoverability (*there exists only a finite number of decompositions*) thanks to algebraic geometry tools. This consists in the re-parametrization of the coupled model and the study of the rank of its Jacobian. An algorithm is proposed to search for the maximal rank for recoverability. This algorithm yields sufficient recoverability conditions that can be applied to a wide range of situations including fully and partially coupled tensor models. When applied to random and balanced couplings, the recoverability bounds show that random couplings may lead to the so-called defective cases with lower recoverability bounds. Those defective cases do not appear when balanced couplings are considered. Thanks to the polynomial additive model framework, this chapter also gives sufficient identifiability conditions (generic uniqueness) for a particular set of couplings. By examining the set of constraints of the corresponding model, it is also proved that the identifiability of the coupled model is implied by the identifiability of a non-negative tensor decomposition. This allows the derivation of improved sufficient identifiability conditions.

Finally, chapter 6 presents the computational package called *CTFlowHD* which stands for Coupled Tensor for Flow cytometry in High-Dimensions. The core method of this package is the algorithm PCTF3D which, when applied to flow cytometry data, results in a low complexity model of the joint density consisting in the factors of the naive Bayes model. These factors can be then post-processed using standard classification and clustering methods to provide biologically meaningful interpretations with a reduced computational burden. The versatility of *CTFlowHD*

is illustrated on a controlled dataset where multiple available visualization tools are presented. In a real biological context, *CTFlowHD* is able to identify blood cell populations in a 8-variable experiment. Lastly, the results of *CTFlowHD* are displayed for a 24-marker experiment whose interpretations are still in progress.

List of Publications

Conference proceedings:

- [43] **Philippe Flores**, Guillaume Harlé, Anne-Béatrice Notarantonio, Konstantin Usevich, Maud d’Aveni, Stéphanie Grandemange, Marie-Thérèse Rubio, and David Brie. *Coupled Tensor Factorization for Flow Cytometry Data Analysis*. IEEE 32nd International Workshop on Machine Learning for Signal Processing (MLSP), pages 1–6, 2022;
- **Philippe Flores**, Joseph K. Chege, Konstantin Usevich, Martin Haardt, Arie Yeredor and David Brie. *Probability Mass Function Estimation Approaches with Application to Flow Cytometry Data Analysis*. 2023 IEEE 9th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP), 2023.

National conference proceedings:

- [44] **Philippe Flores**, Guillaume Harlé, Anne-Béatrice Notarantonio, Konstantin Usevich, Maud d’Aveni, Stéphanie Grandemange, Marie-Thérèse Rubio, and David Brie. *Factorisation Couplée de Tenseurs pour l’Analyse de Données de Cytométrie en Flux*. GRETSI: XXVIIIème Colloque Francophone de Traitement du Signal et des Images, 2022;
- [45] **Philippe Flores**, Konstantin Usevich and David Brie. *Identifiabilité de Modèles Tensoriels Couplés pour l’Estimation de Loi de Probabilité Discrète*. GRETSI: XXIXème Colloque Francophone de Traitement du Signal et des Images, 2023.

Code

Simulations in this thesis were run on a MacBook Pro with 2.3 GHz Intel Core i9 and 32GB RAM. The code is implemented in MATLAB[®]. The n-way tensor toolbox [5] was used as a foundation for the code. The code for *CTFlowHD* is provided in a GitHub repository available at https://github.com/philippesflores/fcm_ctflowhd.

Chapter 1

Flow cytometry

Contents

1.1	General principle and applications	10
1.1.1	Geometrical data acquisition	12
1.1.2	Fluorescence data acquisition	13
1.2	Flow cytometry data pre-processing	16
1.2.1	Flow cytometry dataset example	16
1.2.2	Compensation	18
1.2.3	Cleaning	20
1.2.4	Non-linear transformation	23
1.3	Flow cytometry data analysis	24
1.3.1	Manual data analysis: Gating	24
1.3.2	Computational flow cytometry	27
1.4	Conclusion of Chapter 1	31

This chapter aims at introducing the problem of flow cytometry data analysis which is addressed in this thesis. Flow cytometry is a powerful tool that analyzes biological cells and is used in many applications.

First, the principles of flow cytometry will be introduced by presenting how flow cytometry data is acquired in practice. Flow cytometry measurements can be divided into two categories: geometrical measurements that gives insights into cell morphology and optical measurements used to detect the presence of particular biological properties.

Before analysis, flow cytometry data pre-processing includes 3 steps: *compensation*, *cleaning* and a *non-linear transformation*. These pre-processing steps are presented in Section 1.2 and illustrated on a controlled flow cytometry dataset designed especially for this work.

Finally, the problem of flow cytometry data analysis is introduced in Section 1.3. The first method presented is manual gating which shows limitations in practice, especially because of the growing number of variables. Bigger datasets have imposed to end users the development of automated tools that are presented in Section 1.3.2. A few representative examples are featured as well as remaining challenges for the field of computational flow cytometry.

1.1 General principle and applications

A flow cytometer is a device that measures cell biological properties. It takes a solution with cells in suspension as an input and outputs the biological properties of each cell. As its name suggests, a flow cytometer creates a flow of individual cells from the volume of cells (see Figure 1.1). What makes flow cytometry (FCM) powerful is that cytometers can analyze thousands of cells per second which permits to analyze large cell pools (up to millions) making FCM a more objective and faster method compared to biological studies carried out manually [99].

Originally, FCM was used to sort cells according to their geometrical properties. However, with the development of biomarkers, FCM has moved to Fluorescent Active Cell Sorting (FACS). A biomarker is a molecule or a gene that is present in a certain type of cell. The principle of FACS is then to identify a cell by measuring which biomarkers it contains and in what quantities. There exist numerous biomarkers that can be possibly chosen in FCM experiments. Human cell differentiation molecules are presented in [118]. For example, the marker CD3 is expressed by Lymphocytes T and therefore permits to identify this cell population while Macrophages are identified with the marker CD14.

In order to detect the presence of markers, there exists a specific antibody that can bound to

each marker. Antibodies are named after their associated markers. For example, the CD3 marker antibody is called Anti-CD3. Thus, it is common not to mention antibodies in FCM experiments and to only refer to biomarkers. Antibodies are not fluorescent molecules but they can be dyed with fluorophores. A fluorophore is a fluorescent molecule that can be bonded to an antibody. When lit by a laser, the group biomarker-antibody-fluorophore (if present in a cell) will emit light that will permit to detect the presence of this biomarker. A schematic example is given in Figure 1.2.

The principle of FCM (see Figure 1.1) is as follows: the flow of cells passes through one or more lasers. The light re-scattered by each cell is then separated thanks to dichroic mirrors. Lasers and dichroic mirrors wavelength properties are chosen according to the biological properties of interest for the FCM experiment. Flow cytometers output two types of data for each cell. The first type characterizes the cell morphology like size or granularity while the other type is related to the fluorescence properties of cells [55]. Together with advances in fluorophore design, FACS has evolved to increase the number of parameters that can be measured simultaneously on a single cell. While the first FACS were using 2 or 4 parameters, nowadays 18-parameter FACS

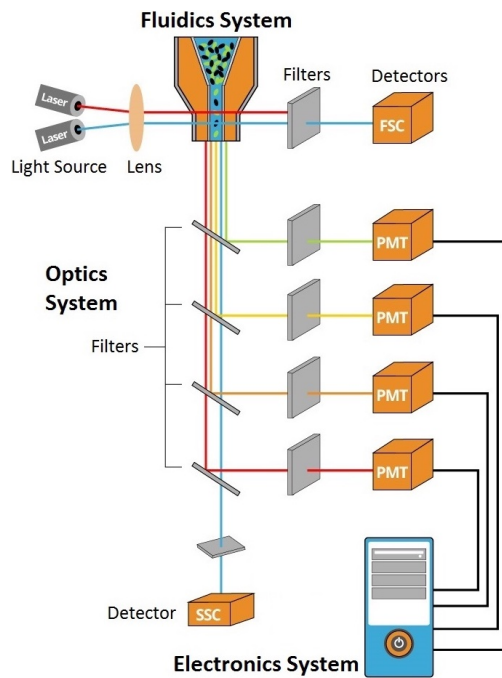


Figure 1.1: Flow cytometer principle. The fluidic system permits to create a flow of individual cells. Lasers emit light through the cell flow and the light scattered by each cell is separated thanks to the optics system. Image taken from [6].

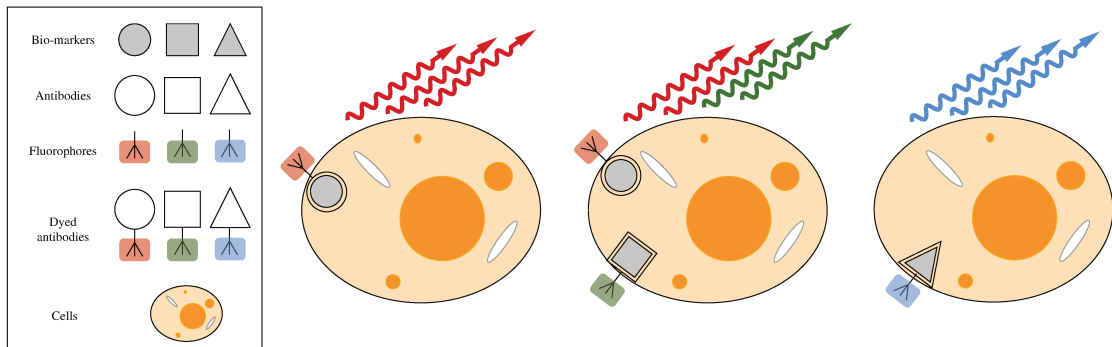


Figure 1.2: Fluorescence flow cytometry principle. When red blue and green light passes through cells, fluorophores associated with markers emit light in particular wavelengths so that biomarkers can be associated with cells.

are routinely used and 30 to 50-parameter FACS are commercially available (see Appendix A).

FCM has proved its worth in numerous applications: it has been used to count Lymphocytes T [52], and to study DNA and chromosomes [18]. Today, FACS can be considered as the reference method for the analysis of biological cells, either in oncology [50, 8, 112], environmental science [113] or agronomy [98]. However, the application where FACS has had the greatest impact is certainly immunology. One of the first applications of FACS was to monitor the immunological status of patients by counting T cells [99]. It is still widely used in HIV studies [22, 87, 85, 77]. In addition, with the increasing number of biomarkers, FACS has enabled to better study and understand the human immune system [82, 47, 21], with important implications for the diagnosis and treatment of leukemia [17, 32, 105].

1.1.1 Geometrical data acquisition

The first type of data measured by FCM is related to cell geometric properties. When a flow containing cells goes through a flow cytometer, the flow is lit by a laser. Two photo-detectors are placed after the laser (see Figure 1.1). The first detector measures Forward SCatter (FSC) and is placed directly in front of the laser. The second detector called Side SCatter (SSC) is placed with a 90 degrees angle with the laser beam and is intended to measure the light re-scattered by the cell. When the laser is lit without any cell, the FSC detector is measuring the intensity of the laser and the SSC detector does not receive any light. However, when a cell passes through, it is shadowing the FSC detector (see Figure 1.3). Moreover, cells scatter light to the SSC channel. Both detectors provide peak-like signal from which three values are determined: FSC-A and SSC-A which measure the area under the peak, FSC-W and SSC-W measuring the width of the peaks

and FSC-H and SSC-H measuring the height of the peaks. When combined with the flow rate of the fluidic system, FSC values permit to estimate the geometric properties of a cell [99] such as diameter and volume. Also, it is worth mentioning that SSC results from the light refracted or reflected at the interface between the laser and intra-cellular structures, such as granules and nucleus (see Figure 1.4). Thus, SSC provides information about the internal complexity (i.e. granularity) of a cell.

These geometric measurements (FSC and SSC) are used mainly in the pre-processing steps which will be presented in section Section 1.2. For example, it is possible from those geometric measurements to sort three populations of blood cells: Lymphocytes, Monocytes and Granulocytes (see Figure 1.5).

1.1.2 Fluorescence data acquisition

In addition to geometrical measurements, FACS aims at measuring the presence of specific biomarkers in a volume of cells. In FCM, a biomarker refers to a state or a component of cell.

Biomarkers must be chosen in accordance with the cell populations that need to be identified. For example, the biomarker TCR (T Cells Receptor) is a protein complex found on the surface of T cells. Antibodies also reflect the presence of specific biomarkers; they can themselves attach to a protein complex. Coming back to the example of Lymphocytes T, the antibody CD3 (which stands for Cluster of Differentiation number 3) can bind to the biomarkers TCR and

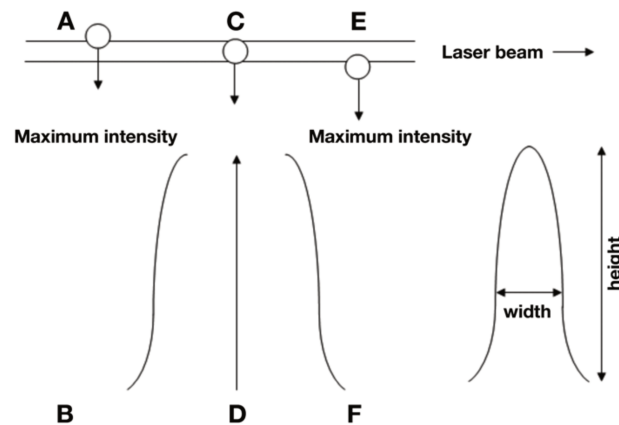


Figure 1.3: Principle of FSC value acquisition [9]. When a cell enters the beam (A), a signal is generated (B) and increases until the cell is in the center of the beam (C). Hence, the signal reaches maximum intensity (D). As the cell leaves the beam (E), the signal falls back to zero (F).

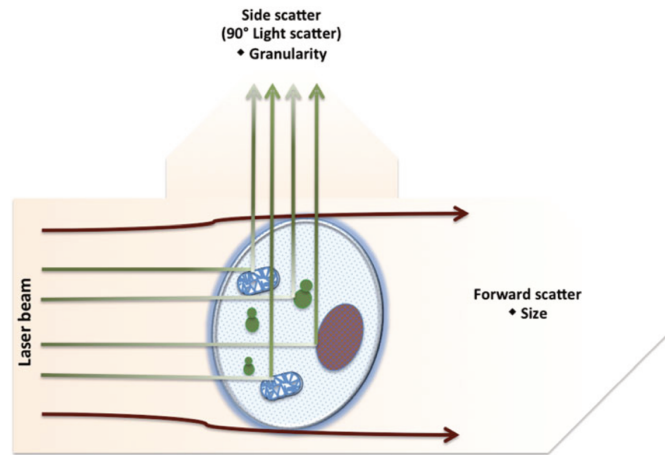


Figure 1.4: Forward and side scatter for a cell [2].

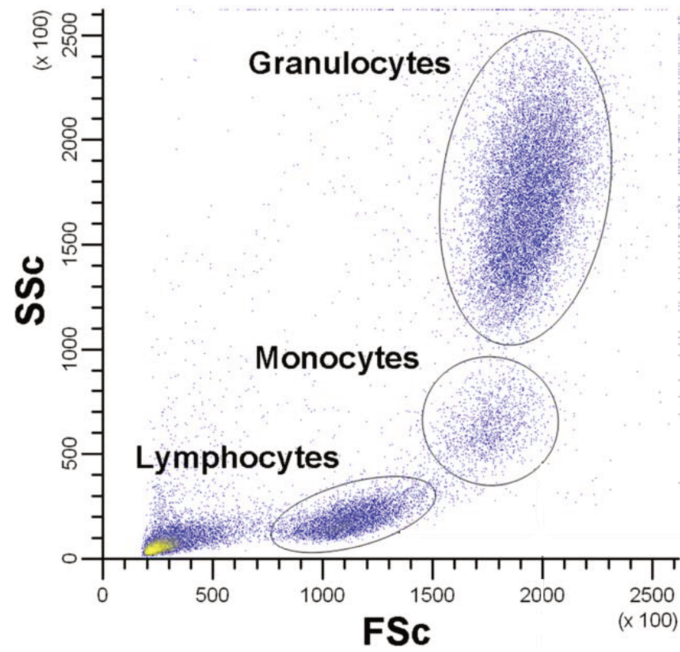


Figure 1.5: Bivariate plot (FSC-SSC) human blood cells [104]. Area values are plotted here: FSC-A and SSC-A.

form the complex TCR/CD3 that will be used to identify T-cell population. Throughout the years, antibodies were found for many biomarkers. For example, Lymphocytes B antibodies are CD19 and CD20 while CD14 is used to identify Monocytes.

Both markers and antibodies have a fluorescence property called auto-fluorescence. However, in most cases, auto-fluorescence is not significant and does not permit to identify a cell. To tackle this issue, each antibody is associated to a fluorophore, that is a fluorescent component that can be attached to an antibody. A fluorophore is characterized by an absorbance spectrum (also referred to as excitation spectrum since it reflects how the excitation light is absorbed by the fluorophore) and an emission spectrum. For a laser whose emission spectrum is compatible with the emission spectrum of the fluorophore, a light will be produced with the corresponding emission spectrum. By according a photo-detector with the emission spectrum of the fluorophore, it is possible to measure the presence of this fluorophore, hence its associated antibody, hence the marker of interest.

Let us summarize the identification of a cell population. First a marker that permits its characterization is chosen. Then, an antibody specific to this marker is dyed with a fluorophore. The volume of cells is then dyed with the pair antibody/fluorophore. At this point, there are two types of cells:

- Positive cells: the cells that express the marker are called positive cells. For example, T cells have TCRs hence express the marker CD3. T cells are called CD3-**positive** cells or CD3+ cells,
- Negative cells: the cells that do not express the marker are called negative cells. For example, B cells does not contain TCR hence they do not express the marker CD3. These cells are called CD3-**negative** cells or CD3- cells.

The volume of cells is then passed into a flow cytometer and the fluorescence of each cell is measured. In the following, the term marker will be used for the triple marker/antibody/fluorophore for conciseness and simplicity purposes. One of these terms will be used when it will be required to be specific.

For clinical applications, most markers are chosen a priori. However, for research applications, it is possible to select a set of markers for an exploratory analysis aiming at revealing new (unknown) populations.

Table 1.1: FCM variable information for the controlled FCM experiment. PE stands for PhycoErythrin and BV stands for Brilliant Violet™.

Variable	Full name	Fluorescence	Fluorophore	Emission maximum
CFSE	CarboxyFluorescein Succinimidyl Ester	Yes	–	488nm
TCR	T Cell Receptor	No	PE	573nm
CTV	CellTrace™ Violet	Yes	–	421nm
MHC-II	Major Histocompatibility Complex of class II	No	BV-510	510nm

1.2 Flow cytometry data pre-processing

From a data analysis point of view, a cytometer permits to obtain an observation matrix \mathbf{X}_{raw} containing N rows and M columns. First, M corresponds to the number of variables measured in a FCM experiment, including both geometrical and fluorescence variables. Concerning N , it corresponds to the number of FCM events registered by the flow cytometer. The observation matrix \mathbf{X}_{raw} is contained in a `.FCS` file which is a standardized file format for FCM [102]. Along with FCM data, a `.FCS` file can contain metadata that provides information on the conditions of data acquisition in a FCM experiment.

1.2.1 Flow cytometry dataset example

In order to show how a `.FCS` file is organized, this section presents a flow cytometry dataset that will be used as a framework example in the following.

Flow cytometry experiment

The goal of this FCM experiment was to provide a FCM controlled dataset. Thus, instead of using a sample with unknown cell populations, three cell populations were chosen and were mixed to create a volume of cells. The three populations were lymphocytes T or T cells, lymphocytes B or B cells and macrophages. In addition to geometrical variables, the experiment featured 4 variables: CFSE, TCR, CTV and MHC-II (see Table 1.1). The markers TCR and MHC-II are classical biomarkers that are associated with antibodies and fluorophores. However, CFSE and CTV are fluorescent complexes that bind to free amino acids. Thus, fluorophores are not needed for those markers.

Table 1.2: Cell populations used in the controlled FCM experiment with their marker expressions. For each marker, M denotes if the cell populations is marked with the marker while the E-column features the expected expression of cell populations.

Cell population	CFSE		TCR		CTV		MHC-II	
	M	E	M	E	M	E	M	E
Lymphocytes B	Yes	+	No	-	No	-	Yes	++
Lymphocytes T	No	-	Yes	+	No	-	No	-
Macrophages	No	-	No	-	Yes	+	Yes	+

Let us go through the chronological process of the experiment. First, each cell population was marked with either a marker (TCR) or a free amino acid detector (CFSE or CTV). T cells were marked with TCR, B cells with CFSE and macrophages with CTV. Then, the three cell populations were mixed in 3 different proportions, thus creating 3 FCM datasets. Each volume was then marked with the marker MHC-II¹. Finally, samples were passed individually in a flow cytometer. To give some figures, samples contain between 50,000 to 100,000 cells.

Flow cytometry metadata

First, let us present the metadata contained in a '.FCS' file. First, the list of variables measured during the experiment is available (see Figure 1.6). A mandatory information contained in the metadata is the compensation matrix. This matrix is needed to perform the compensation pre-processing step (see Section 1.2.2). Other information can be present in the metadata: operating system, cytometer reference, time of analysis, analysis duration, byte order, etc.

Flow cytometry data

FCM raw data can be plotted in bivariate scatter plots. For example, Figure 1.7a shows a way to represent geometrical data. This point cloud is plotting the SSC value regarding FSC. It is shown on this plot that the higher the FSC value is, the higher is the SSC value. This can also be interpreted by the fact that bigger cells re-scatter more light compared to smaller cells. On the other hand, it is possible to plot fluorescence raw data (see Figure 1.7b). On this figure, it is possible to distinguish 3 different populations. It seems that rightmost cells on this plot are B

¹In this experiment, cells were marked separately to ensure that only the cells marked express their associated marker. This was possible because – unlike classical FCM experiments – cells were separated and mixed in this controlled experiment.

	name	marker	range	bit
1	'FSC-H'	'FSC-H'	16777216	32
2	'FSC-A'	'FSC-A'	16777216	32
3	'SSC-H'	'SSC-H'	16777216	32
4	'SSC-A'	'SSC-A'	16777216	32
5	'FL1-H'	'CFSE-H'	16777216	32
6	'FL1-A'	'CFSE-A'	16777216	32
7	'FL2-H'	'TCR PE-H'	16777216	32
8	'FL2-A'	'TCR PE-A'	16777216	32
9	'FL3-H'	'CTV-H'	16777216	32
10	'FL3-A'	'CTV-A'	16777216	32
11	'FL4-H'	'MHCII BV510-H'	16777216	32
12	'FL4-A'	'MHCII BV510-A'	16777216	32
13	'FSC-W'	'FSC-Width'	10000	32
14	'Time'	'Time'	900000000	32

Figure 1.6: Screenshot featuring a table containing FCM variables information.

cells as they are CFSE+. With the same reasoning for the CTV marker, the upmost cells on this plot are CTV+ hence macrophages. This leaves the last cells to be T cells which do not express neither CFSE neither CTV. Note that Figure 1.7b is a bivariate plot with logarithmic scales (see Section 1.2.4)

1.2.2 Compensation

In order to identify a large number of cell populations, FCM experiments are conducted with more and more biomarkers, each biomarker being associated to a specific fluorophore characterized by its emission spectrum. When the number of fluorophore increases, the emission spectra are likely to overlap: light intensities measured by photo-detectors is a mix of different emission spectra. An example is given in Figure 1.8. This example features the fluorophores FITC (Fluorescein isothiocyanate) and PE [88]. In this figure, both fluorescence emission spectra are plotted. FITC emission is maximal around 515nm; typically, a filter centered on 530nm is used to collect the emitted light. The emission of PE is farther red, with a maximum around 575nm; typically, a filter centered on this emission maximum is used to collect PE. Note that FITC has some emission in the wavelength bands used to collect PE fluorescence (B); typically, the amount of light in the 575nm band is 15% of that in the 530nm band (A). The PE has very little emission in the 530nm band (C), usually less than 2% of the emission in the 575nm band (D). Compensation permits to unmix light intensities received by photo-detectors. By knowing the level of overlapping of each pair of fluorophores, it is possible to create a matrix $\mathbf{S} \in \mathbb{R}^{M \times M}$

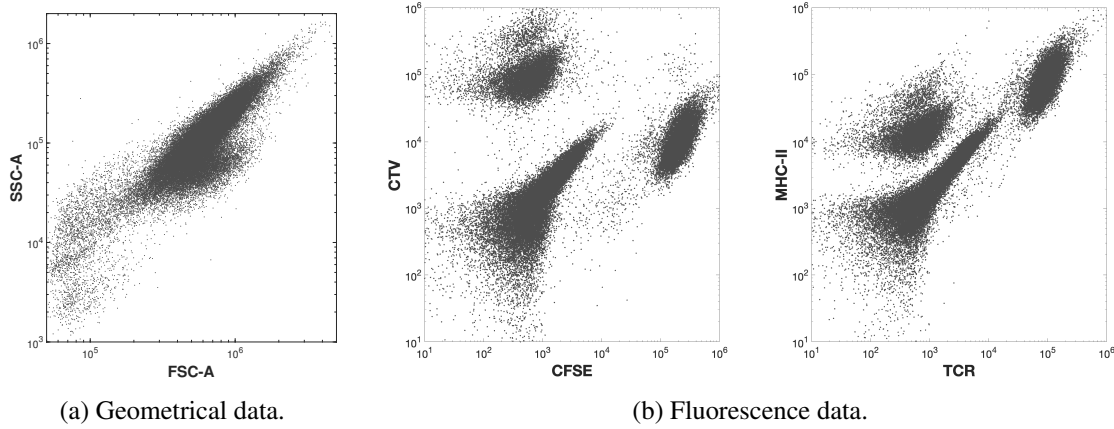


Figure 1.7: FCM data represented with 2-dimensional scatter plots.

storing the overlapping values. Because compensation is linear, applying compensation is equivalent to solving a set of linear equation which can be done by inverting the compensation matrix [88]:

$$\mathbf{X}_{\text{comp}} = \mathbf{S}^{-1} \mathbf{X}_{\text{raw}}.$$

Note that even if compensation permits to unmix fluorescence, FCM experiments are designed such that the spectrum overlapping between fluorophores is minimal.

Compensation applied to the controlled dataset

In the controlled experiment, 4 variables are subject to compensation. The compensation matrix is contained in the metadata and is equal to:

$$\mathbf{S} = \begin{bmatrix} 1 & 0.4361 & 0.0095 & 0.0733 \\ 0.513 & 1 & 0.5184 & 0.8 \\ 0 & 0.0004 & 1 & 0.0878 \\ 0 & 0.0057 & 0.1307 & 1 \end{bmatrix}$$

Figure 1.9 shows the effect of the compensation pre-processing step. The three cell populations observed on raw data are still visible after compensation. Because compensation unmix the data provided by each channel, compensated data seem to be distributed more independently across variables. Both plots of Figure 1.9 are shown with a logarithmic scale which is not ideal but convenient for the moment (see Section 1.2.4).

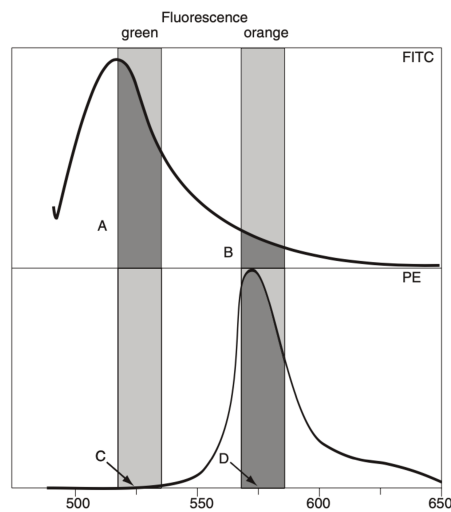


Figure 1.8: Overlapping of fluorophore emission bandwidths for fluorophores FITC and PE [88]. Emission spectra are shown in black. Photodetector reception bandwidths are shown with shaded regions.

1.2.3 Cleaning

In the raw data acquired through an FCM, not all events do correspond to actual cells; this is why cleaning needs to be applied to extract cells from outlying events. Three types of outliers are presented in this section as well as possible methods to detect and remove them.

Time inconsistencies: It is possible that the flow of individual cells is not consistent due to clogs or in transitional periods (beginning and end of the experiment) for example. To prevent that, only events in the stationary flow rate zone are kept for further analysis [9].

Dead cells and debris: Despite all precautions, it is not possible to guarantee that the cells present in the suspension volume are still alive at the time of analysis. Moreover, when cell apoptosis (cell death) occurs, cell debris are scattered throughout the cell suspension and may be counted as events, creating outliers in the analysis. To eliminate both cell debris and apoptotic cells, bivariate SSC and propidium iodide diagrams allow only living cells to be selected. Propidium iodide is a laboratory reagent for the fluorescent staining of nucleic acids; it is used in flow cytometry for staining apoptotic cells and nuclei. As the membrane of a living cell is impermeable to propidium iodide, the state of the cell determines the ability of propidium iodide to stain viable and apoptotic cells (see Figure 1.10) [9].

Doublets: The simultaneous passing of two cells through the laser beam results in an outlier event. To eliminate those events called doublets and to keep only single events called

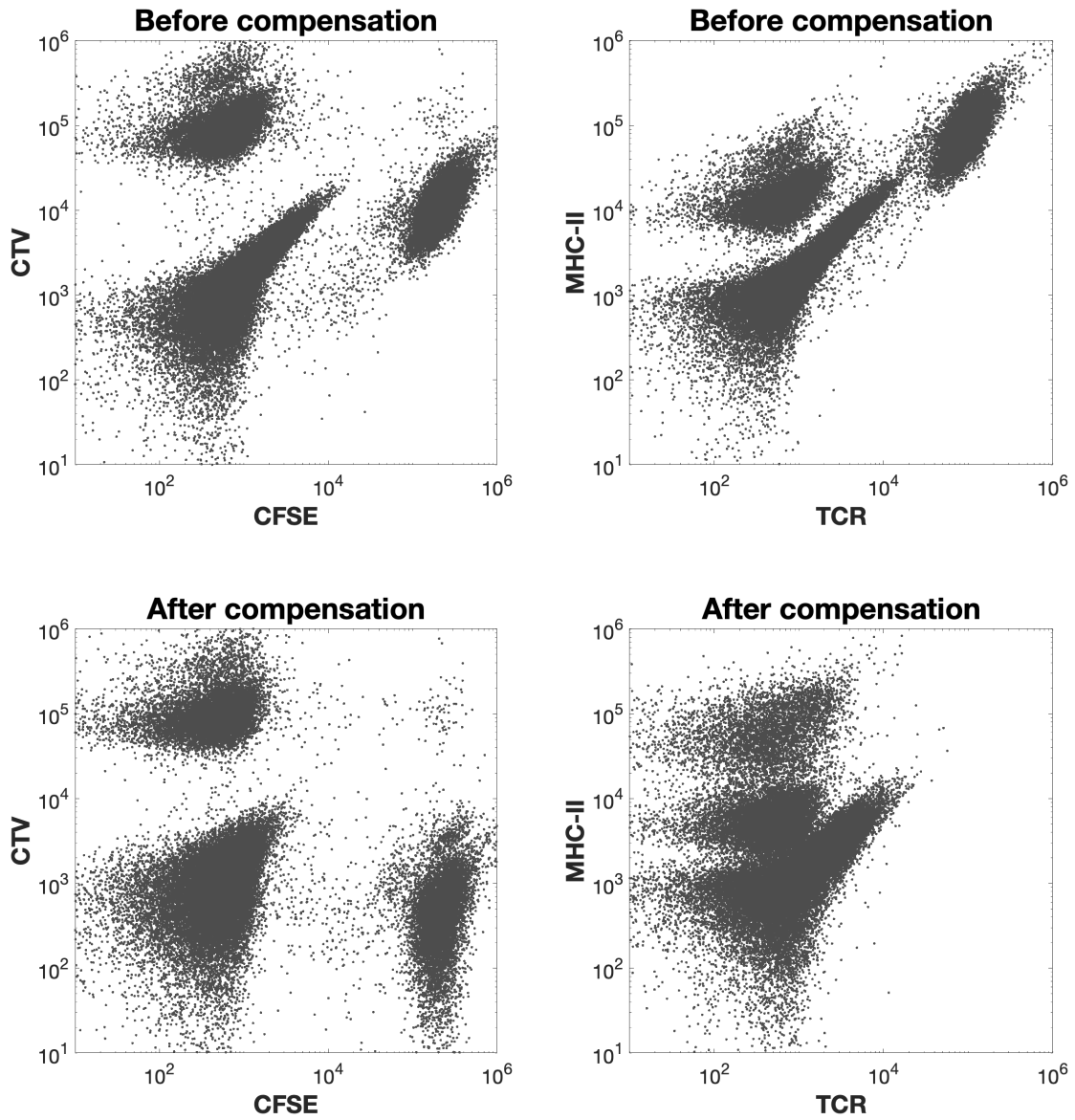


Figure 1.9: Effect of compensation on raw FCM data.

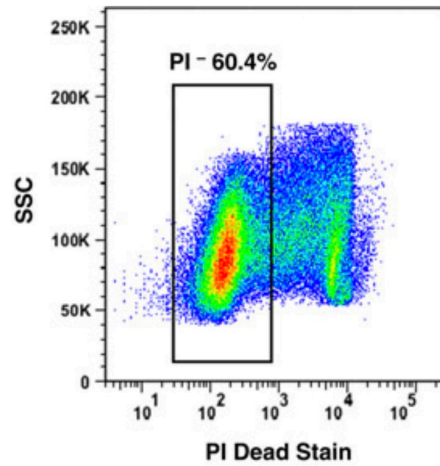


Figure 1.10: Viable cell detection gate with SSC and propidium iodide bivariate plots [90].

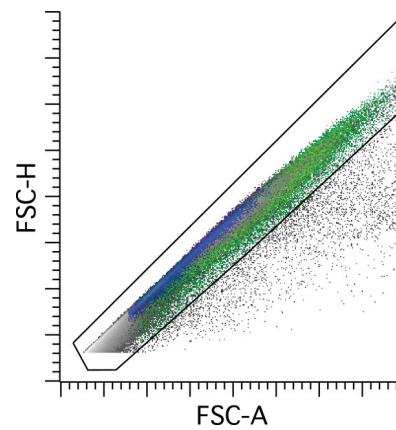


Figure 1.11: Singlets cell detection gate with FSC-A and FSC-H bivariate plots [115].

singlets, the bivariate graph containing the height of FSC (FSC-H) and the area of FSC (FSC-A) is plotted [24]. For example, let us consider a case where two cells are passing through the laser at once side by side. In this case, the ratio between the width of the signal (see Figure 1.11) and its height is doubled. Another extreme case is when one cell is eclipsed behind another cell. In this case, the FSC area value will remain the same. However, two cells will result in twice more light received by SSC detectors. In summary, by interpreting the FSC and SSC values, it is possible to detect doublets.

Implementation

Cleaning can be performed manually by selecting cells on specific bivariate scatter plots. However, applying cleaning manually introduced subjectivity in the FCM data analysis process. There exist several packages that perform FCM pre-processing steps. In this work, we used the package FlowAI [73]. This package is developed in R and permit to clean FCM datasets from the previously presented outliers. We provide R scripts at github.com/philippflores/fcm_ctflowhd/tree/main/preprocessingR.

1.2.4 Non-linear transformation

The data provided by flow cytometers is a unit-less real positive number ranging between 0 and 10^6 . While some cells may respond strongly (positive cells), some others may only respond slightly which results in a high dispersion within the FCM fluorescence data. To visualize and analyze FCM fluorescence data, a non-linear transformation is generally applied along each fluorescence parameter. Due to the wide dispersion, a logarithmic scale can be applied but it often results in artifacts especially for negative cell populations. Indeed, compensation may result in data containing slightly negative values. In addition, a pure logarithmic scale over-amplifies small variations in the neighborhood of 0. A commonly used transformation is the "Logicle" scale [80] that can handle both negative and positive values. It is based on the arcsinh function and behaves as a linear scale around zero values and like a logarithmic scale for large values. Alternative log-like transforms are also available in the literature [7, 106].

To conclude this section, it is worth mentioning that the design of non-linear transformations used in flow cytometry shares strong similarities with non-uniform quantization since the goal of these non-linear transformation is to flatten the probability density of the considered one dimensional variable.

Implementation

Non-linear transformations can be applied manually by choosing two parameters: a parameter to set the range of FCM data and a parameter that fixes the threshold between the linear behavior of the scale and the logarithmic behavior of the scale. In this work, we used the package flowTrans [42]. This package is developed in R and permits to set the parameters of the "Logicle" scale and apply this transformation to FCM datasets. We provide R scripts at github.com/philippflores/fcm_ctflowhd/tree/main/preprocessingR.

Non-linear transformation applied to the controlled dataset

In the controlled experiment, a non-linear transformation must be applied to every fluorescence variable. Figure 1.12 shows the effect of the non-linear transformation "Logicle" applied to this dataset.

1.3 Flow cytometry data analysis

1.3.1 Manual data analysis: Gating

To analyze flow cytometry datasets, events are displayed in two ways: one-dimensional histograms and a bivariate point clouds if two or more markers are considered [78, 89]. Cell subpopulations in a bivariate zone of interest are then selected for further analysis. These sequential analysis are denoted *Gating* and are performed manually. An example of a gating sequence on scatter plots is shown in Figure 1.13.

Gating applied to the controlled dataset

In the framework of the controlled experiment, gating is applied to characterize the 3 cell populations. Figure 1.14 shows how the 3 cell populations can be identified with a bivariate scatter plot. First, the point cloud plotted on this figure permits to confirm that 3 distinct populations are contained in the sample. Then, with the information of Table 1.2, it is possible to identify each population:

- blue cells are CTV+ hence the blue population corresponds to macrophages,
- green cells are CFSE+ hence the green population corresponds to lymphocytes B,
- By elimination, red cells which are CFSE-/CTV- corresponds to lymphocytes T.

Finally, by counting how many cells are in each gate, it is possible to quantify the proportion of each cell population: T cells represents 56.1% while B cells and macrophages represent respectively 21.9% and 20.8%. Note that the sum of these proportion is not equal to 1, as there are "ungated" cells which do not fall into any of gates.

Limitations of gating

Even if gating is the most used method to analyze FCM datasets [25], this technique shows practical limitations that pushed the development of automated FCM data analysis method (see

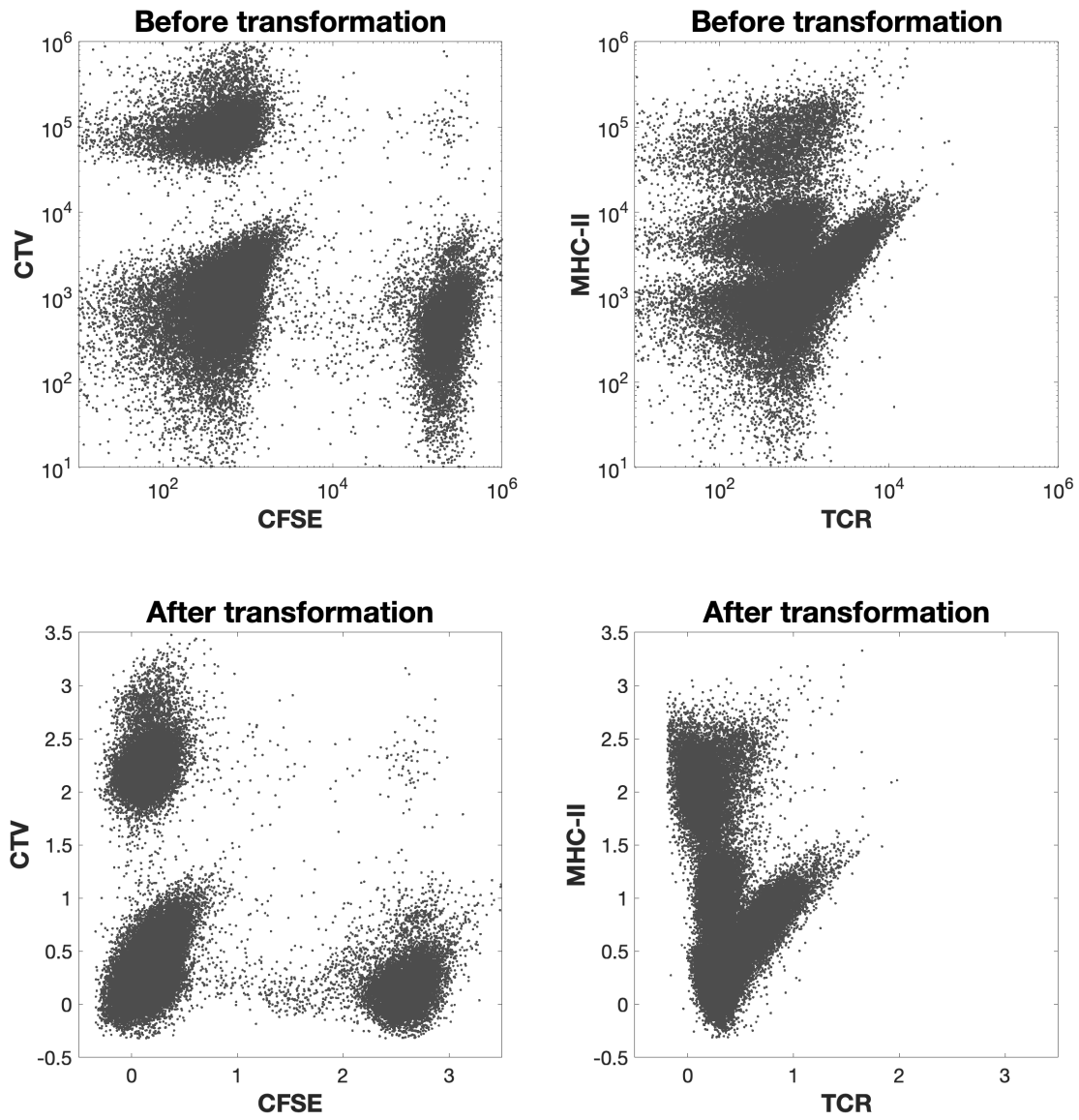


Figure 1.12: Application of a "Logicle" scale to the controlled dataset.

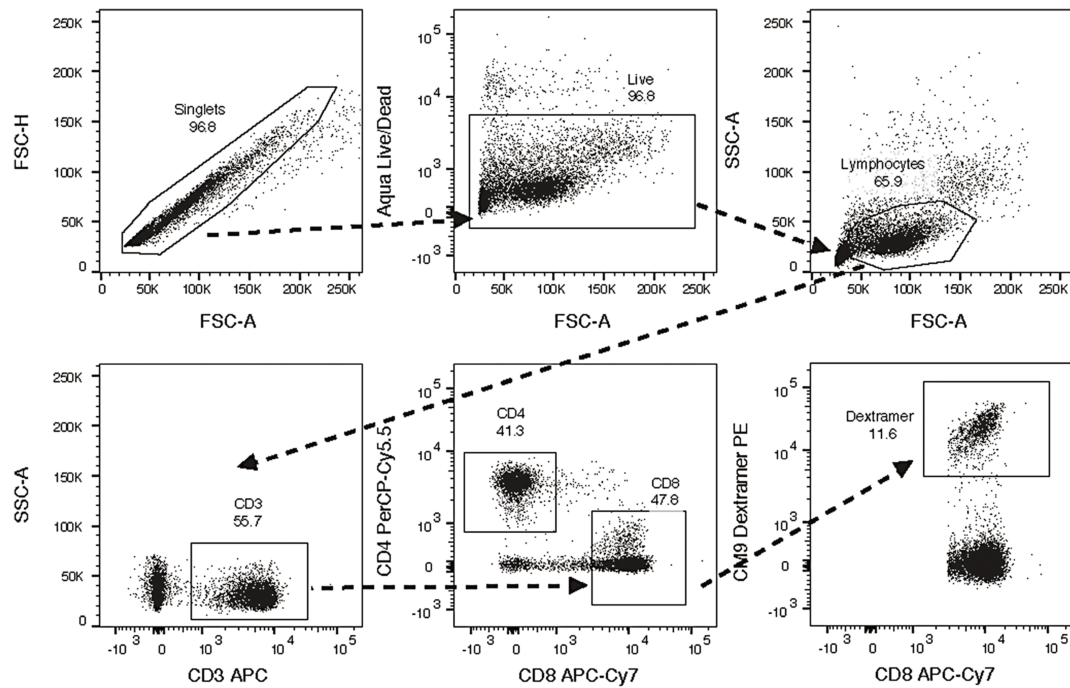


Figure 1.13: Example of manual gating sequence. The first three steps (upper plots) features physical characteristics (size and diameter). The last 3 steps (lower plots) permit to separate cells according to the fluorescence/biological properties. For each fluorescence plot, axis represent the amount of fluoresced light emitted by cells. Image taken from [72].

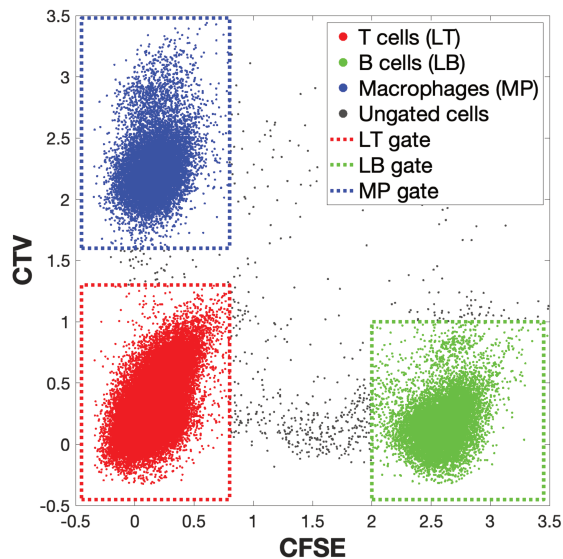


Figure 1.14: Gating applied to the controlled dataset.

Section 1.3.2).

Because gating is carried out manually by practitioners, gates introduce subjectivity into FCM data analysis. Indeed, some practitioners are choosing more inclusive gates than others which may create difficulties for comparison and reproducibility of results among studies. Here, we would like to emphasize that the controlled FCM experiment (see Section 1.2.1) was designed to have 3 distinct cell populations in terms of fluorescence properties. Thus, gating can be easily applied without introducing subjectivity in the process. That way, the gating results for this experiment can be considered as ground-truth which cannot be the case for real FCM experiments.

Another limitation of gating comes with the augmentation of the number of biomarkers being analyzed jointly. Indeed, when the number of FCM variables M increases, the number of possible bivariate scatter plots increases quadratically with M [92]. This leads to difficulties in the exploration of complex datasets [82]. As an example, for a $M = 10$, there are $\binom{M}{2} = 45$ possible 2D scatter plots. Even if end users have biological insights for the choice of gating variables, exploring datasets with such a "funnel-type" analysis can be problematic. Indeed, gating is focusing on specific cells while removing others hence preventing exhaustive exploration of FCM datasets.

Finally, gating is time-consuming [3]. For FCM experiments where several samples are analyzed (e.g. for multiple patients or group of patients), repeating the same gating method for each sample may be impossible to apply in practice. Moreover, when experts search for multiple cell populations, the whole gating process must be performed again from scratch.

1.3.2 Computational flow cytometry

Because of the gating limitations (see Section 1.3.1), methods for automated flow cytometry data analysis have appeared over the last 15 years [92]. These methods falls to the recent field of Computational Flow Cytometry (CFC). CFC methods have the advantage of being less time-consuming than gating and more objective. In this subsection, a few representative CFC methods are presented as well as remaining challenges around FCM data analysis.

In clinical use, cell populations of interests are known before-hand. Therefore, end users tend to choose supervised classification to identify and quantify cell populations of interests. One supervised CFC method that can be mentioned is CITRUS [16] which uses supervised learning after a hierarchical clustering of cells. Research-oriented FCM studies often aim at exploring datasets either to find new cell populations of interests, or to find new properties

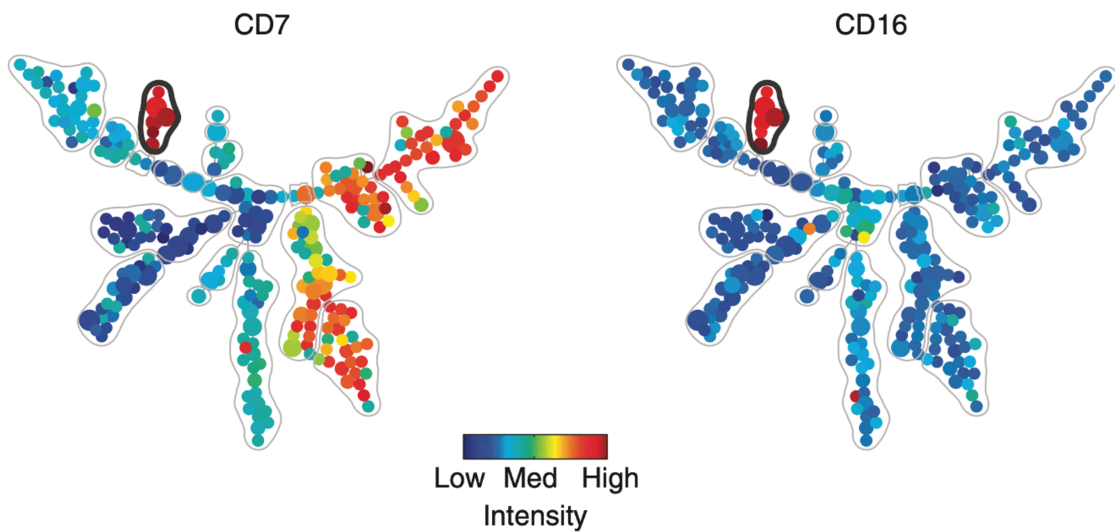


Figure 1.15: Example of SPADE visualization. For each of the 2 markers, the minimum spanning tree is plotted and colored with mean fluorescence intensity.

of already-identified populations [25]. In this case, CFC methods use unsupervised clustering methods to analyze datasets.

According to [25], the most-cited CFC methods are viSNE¹ [4] and SPADE [84]. Those two methods are based on already-existing machine learning methods. SPADE uses k-means [39] and plots the resulting clusters as a minimum spanning tree [49]. An example is plotted in Figure 1.15. On the other hand, viSNE is an adaption of t-Stochastic Neighbor Embedding (t-SNE) [109]. The principle of t-SNE is to map samples on a two-dimensional space. The difference between viSNE and t-SNE is that viSNE performs a down-sampling before applying t-SNE to reduce the computational load of the method. An example of t-SNE is plotted in Figure 1.16. The new mapping provided by t-SNE is chosen such that the distance in between samples in the new mapping is approximately the same as in the high-dimensional space.

Combined with clustering algorithms, CFC methods use dimension reduction to visualize high-dimensional data. For example, SPADE uses a two-dimensional tree to map k-means nodes. It is the same for t-SNE/viSNE, as single cells are mapped on a two-dimensional space. In a following step, end users use their expertise to cluster cells or nodes which have same properties directly on the new visualizations with their expertise.

¹Unfortunately, the code for viSNE is available on demand but we did not receive any answers from the authors.

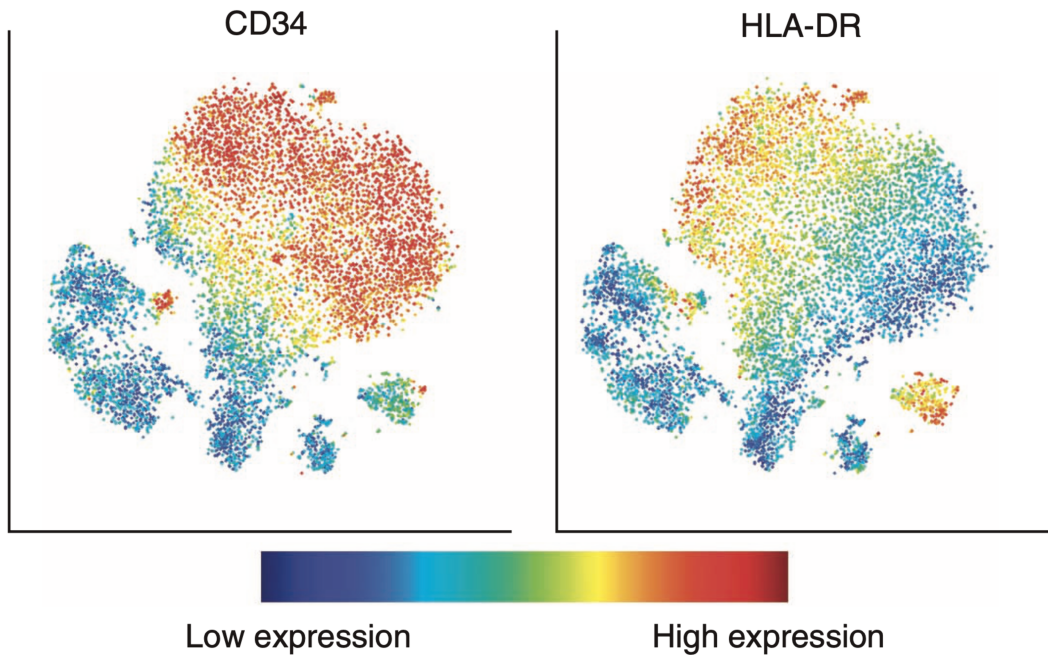


Figure 1.16: Example of viSNE visualization [4, Figure 3]. viSNE uses single cell visualizations because all cells are plotted on a new two dimensions space where color represents the expression of each of the 2 markers considered.

Limitations and remaining challenges in computational flow cytometry

Scaling to more cells and more markers With the development of the recent years, flow cytometers analyze more cells with more variables in a short amount of time. CFC data analysis methods have still not scale to these dataset scales [69]. Indeed, some methods can take hours to run (see [92, Table 2]) which prevent their use in clinical conditions. Most CFC methods do not handle datasets with large number of cells [69]. This is the case of t-SNE hence viSNE that can become computationally prohibitive [10]. Methods that can handle high numbers of cells like SPADE does not scale very well to high number of variables. Indeed, SPADE's complexity relies on k-means whose complexity increases with the number of variables [53]. An ongoing challenge is then to reduce the computational burden of the methods used in CFC.

This problem is aggravated by the fact that end users may need to apply multiple tools to fully interpret FCM datasets [3]. As an example, let us study the two methods SPADE and t-SNE for the two following setups in terms of computational load:

- Situation **S1**: $N = 10^4$ cells and $M = 20$ markers,
- Situation **S2**: $N = 10^6$ cells and $M = 6$ markers.

The situation **S1** features a lot of markers but not a high number of cells while **S2** is featuring more cells for fewer marker. SPADE is performing well in **S2** while t-SNE performs well in **S1**. However, if one wants to apply both methods, it will be computationally expensive in both situations.

Note that the problem of run times is strongly linked to the curse of dimensionality [11] which states that the complexity of a problem increases exponentially with the number of considered variables.

Enhancing interpretability and visualization Because of the numerous and various FCM applications, the development of an efficient CFC method for all applications is impossible. Therefore, CFC methods have been developed to perform well in particular setups. For example, CITRUS works very well to compare group of patients but it can be difficult to globally visualize FCM datasets with this method. Indeed, CITRUS does not give an overview of FCM datasets as it focuses on the differences between FCM samples. Methods like t-SNE/viSNE or SPADE give an overview of FCM datasets but show limitations to detect rare cell populations. As end users choose the level of granularity of analysis [69], an ongoing CFC challenge remains in the development of methods that can adapt to different levels of granularity, especially in terms of visualizations.

Detecting rare cell populations Another ongoing challenge concerns the detection of rare cell populations (less than 0.1%). Because more markers are available in FCM experiments, more cell heterogeneity can be revealed by FCM datasets [82]. Cell heterogeneity means that inside a cell populations, there may be intra-groups of cells that behave different from one to another. Characterizing cell heterogeneity is very useful to interpret FCM datasets and find rare cell populations. A potential solution to this problem is to determine cell populations profiles. Instead of determining the properties of a cell population, it may be more informative to assign a range of properties hence a fluorescence profile to each cell population.

This challenge is linked to the computational load of CFC methods. Indeed, a classical workaround to reduce run times is to use down-sampling which consists in analyzing a subset of cells. This may result in the overlooking of rare cell populations.

1.4 Conclusion of Chapter 1

In this chapter, the principles of flow cytometry were presented as well as flow cytometry data analysis methods. The widely used approach for data analysis is manual gating which showed limitations as the number of variables increases. This motivated the development of computational flow cytometry but challenges remain to be solved. This work addresses some of these challenges by developing a probabilistic approach associated with low-rank tensor decompositions.

Chapter 2

Tensor preliminaries

Contents

2.1 Basic definitions	34
2.1.1 Matrix operations	34
2.1.2 Tensors and basic definitions	35
2.2 Canonical Polyadic decomposition	37
2.2.1 Definition of the CP decomposition	37
2.2.2 Uniqueness of the CP decomposition	38
2.2.3 Generic uniqueness	40
2.3 Identifiability of polynomial models	41
2.3.1 Polynomial mapping and uniqueness	42
2.3.2 Uniqueness and recoverability of additive models	43
2.4 Conclusion of Chapter 2	45

In this chapter, we introduce the tensorial background used throughout the manuscript. In Section 2.1, we define tensors as multidimensional arrays. Along with the definition of tensors, we recall a set of classical operations and tools used in the tensor framework. Then, in Section 2.2, we present the Canonical Polyadic Decomposition (CPD) which is the tensor decomposition used in this thesis. We recall the uniqueness properties of the CPD, which is a distinctive feature of the higher-order (tensor) case and does not hold in the matrix case. Finally, we introduce in Section 2.3 the framework of polynomial additive maps that will be used to study the identifiability of coupled tensor decompositions.

2.1 Basic definitions

2.1.1 Matrix operations

First, some useful matrix definitions that will be used throughout the manuscript are recalled.

Definition 2.1.1. Matrix vectorization – For $\mathbf{A} \in \mathbb{R}^{I \times J}$, the column-major vectorization of \mathbf{A} is a vector of size IJ denoted $\text{vec}(\mathbf{A})$ and defined as the stacking of \mathbf{A} columns in the natural order:

$$\text{vec}(\mathbf{A}) := \left[\mathbf{A}_{:,1}^T \quad \cdots \quad \mathbf{A}_{:,J}^T \right]^T.$$

Definition 2.1.2. Kronecker product – For two matrices $\mathbf{A} \in \mathbb{R}^{I \times J}$ and $\mathbf{B} \in \mathbb{R}^{K \times L}$, their Kronecker product is a matrix of size $(IK) \times (JL)$ denoted $\mathbf{A} \otimes \mathbf{B}$ and defined as:

$$\mathbf{A} \otimes \mathbf{B} := \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1J}\mathbf{B} \\ \vdots & & \vdots \\ a_{I1}\mathbf{B} & \cdots & a_{IJ}\mathbf{B} \end{bmatrix}.$$

For details and properties on the Kronecker product, please refer to [67].

Definition 2.1.3. Khatri-Rao product – For two matrices $\mathbf{A} \in \mathbb{R}^{I \times K}$ and $\mathbf{B} \in \mathbb{R}^{J \times K}$, their Khatri-Rao product (also defined as the column-wise Kronecker product) is a matrix of size $(IJ) \times K$ denoted $\mathbf{A} \odot \mathbf{B}$. The resulting matrix is defined with the Kronecker product between columns of both matrices:

$$\mathbf{A} \odot \mathbf{B} := \left[\mathbf{A}_{:,1} \otimes \mathbf{B}_{:,1} \quad \cdots \quad \mathbf{A}_{:,K} \otimes \mathbf{B}_{:,K} \right].$$

Note that for columns vectors $\mathbf{a} \in \mathbb{R}^I$ and $\mathbf{b} \in \mathbb{R}^J$, Kronecker and Khatri-Rao products are identical: $\mathbf{a} \otimes \mathbf{b} = \mathbf{a} \odot \mathbf{b}$. More details on the Khatri-Rao product are available in [59].

Definition 2.1.4. Kruskal rank [61] – For a matrix \mathbf{A} , its Kruskal rank is denoted $\kappa(\mathbf{A})$ and is defined by the maximum integer such that any combination of k columns of \mathbf{A} are linearly independent.

Note that in general, the Kruskal rank of a matrix cannot exceed its rank but these two ranks can be different. Indeed, if there exists a combination of k columns linearly independent, it holds that $\text{rank}(\mathbf{A}) \geq k$ but this does not mean that any combination of columns of \mathbf{A} are linearly independent.

2.1.2 Tensors and basic definitions

Tensors are classically defined as multilinear operators [30] for a set of vector spaces. When the basis of the vector spaces are fixed, tensors are represented by M -way arrays similarly to linear operators for matrices. In this thesis, we do not make distinctions between tensors and M -way arrays.

Definition 2.1.5. Tensor – An order- M tensor is an M -way array of elements in \mathbb{R} .

$$\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_M}.$$

The elements of \mathcal{X} are denoted $x_{i_1 \dots i_M}$ where $i_m \in \llbracket 1, I_m \rrbracket$, $m \in \llbracket 1, M \rrbracket$. A tensor is said cubic if it has the same number of elements along all dimensions ($I_1 = \dots = I_M$).

Note that scalars, vectors and matrices are tensors of order respectively $M = 0$, $M = 1$ and $M = 2$. For $M \geq 3$, M -order tensors are called higher-order tensors. In the following, “higher-order tensors” will be referred to as “tensors”.

For an order- M tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_M}$, fibers are vectors containing entries of \mathcal{X} by fixing all indices except one. For $m \in \llbracket 1, M \rrbracket$, there are $I_1 \dots I_{m-1} I_{m+1} \dots I_M$ m -th mode fibers which can be written

$$\mathcal{X}_{i_1 \dots i_{m-1} : i_{m+1} \dots i_M} := \left[x_{i_1 \dots i_{m-1} 1 i_{m+1} \dots i_M} \quad \dots \quad x_{i_1 \dots i_{m-1} I_m i_{m+1} \dots i_M} \right]^T.$$

Examples of fibers are given in Figure 2.1 for an order-3 tensor. As with rows and columns for matrices, there exist fibers for each mode of \mathcal{X} . The column-major vectorization of \mathcal{X} is denoted

by $\text{vec}(\mathcal{X})$ and is the vector of size $\prod_{m=1}^M I_m$ that stores elements of \mathcal{X} in the order defined in [60]. Note that the tensor vectorization definition is consistent with the matrix vectorization (see Definition 2.1.1).

Definition 2.1.6. Tensor unfolding – A tensor can be transformed into a matrix through tensor matricization. For an order- M tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_M}$, the mode- m unfolding of \mathcal{X} is a matrix of size

$$I_m \times \left(\prod_{\substack{\ell=1 \\ \ell \neq m}}^M I_\ell \right)$$

and is denoted $\mathbf{X}_{(m)}$. The m -th unfolding is built by stacking m -th mode fibers as columns of $\mathbf{X}_{(m)}$ in the column-major vectorization order $(1, \dots, m-1, m+1, \dots, M)$.

Definition 2.1.7. Tensor contraction – Let \mathcal{X} be an order- M tensor of size $I_1 \times \dots \times I_M$ and $\boldsymbol{\mu} \in \mathbb{R}^{I_m}$ a vector of size I_m . The operation of tensor contraction on the m -th mode is denoted \bullet_m . If $\mathcal{Y} = \mathcal{X} \bullet_m \boldsymbol{\mu}$, the tensor \mathcal{Y} is the order- $(M-1)$ tensor of size $I_1 \times \dots \times I_{m-1} \times I_{m+1} \times \dots \times I_M$ such that:

$$y_{i_1 \dots i_{m-1} i_{m+1} \dots i_M} = \sum_{i_m=1}^{I_m} \mu_{i_m} x_{i_1 \dots i_m \dots i_M}.$$

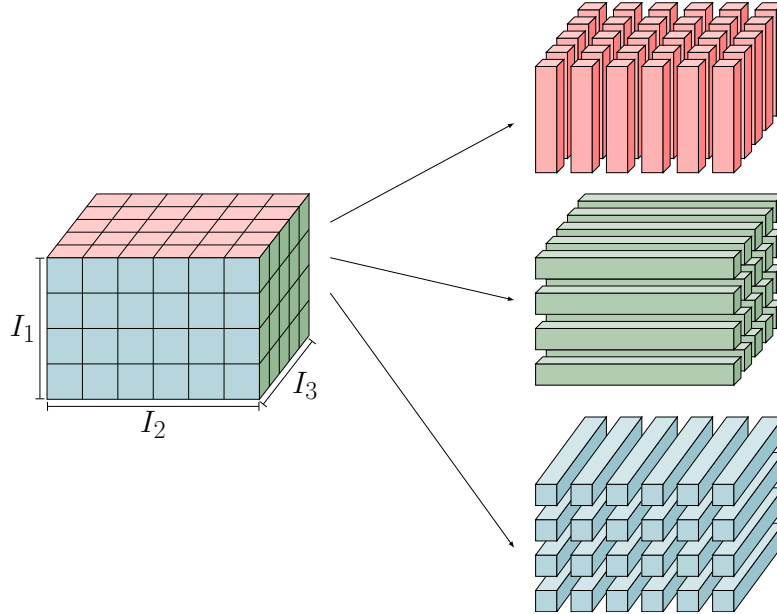


Figure 2.1: An order-3 tensor with its 3 sets of fibers.

2.2 Canonical Polyadic decomposition

2.2.1 Definition of the CP decomposition

Similarly to matrices, higher-order tensors can also be decomposed into low-rank components as shown in this section.

Definition 2.2.1. Rank-one tensor and outer product – A tensor \mathcal{X} is said to be rank-one if it can be written as the outer product of a set of vectors $\{\mathbf{a}^{(m)} \in \mathbb{R}^{I_m}\}_{m=1}^M$:

$$\mathcal{X} = \mathbf{a}^{(1)} \circ \dots \circ \mathbf{a}^{(M)},$$

where \circ denotes the outer product. The outer product of a set of M vectors is an order- M tensor \mathcal{X} of size $I_1 \times \dots \times I_M$ whose entries are defined by:

$$x_{i_1 \dots i_M} = \prod_{m=1}^M a_{i_m}^{(m)}.$$

Definition 2.2.2. Canonical Polyadic Decomposition [54] – The Canonical Polyadic Decomposition (CP decomposition or CPD) of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_M}$ represents \mathcal{X} has a sum of R rank-one tensors:

$$\mathcal{X} := \sum_{r=1}^R \lambda_r \mathbf{a}_r^{(1)} \circ \dots \circ \mathbf{a}_r^{(M)}, \quad (2.1)$$

where $\mathbf{a}_r^{(m)} \in \mathbb{R}^{I_m}$ for $m \in \llbracket 1, M \rrbracket$ and $r \in \llbracket 1, R \rrbracket$ are called the factors of the decomposition and $\boldsymbol{\lambda} := \begin{bmatrix} \lambda_1 & \dots & \lambda_R \end{bmatrix}^\top$ is called the loading vector of size R .

Definition 2.2.3. Factor matrix – Let \mathcal{X} an order- M tensor that admits a CPD (2.1) of rank R . For each mode $m \in \llbracket 1, M \rrbracket$, the m -th factor matrix $\mathbf{A}^{(m)}$ is the matrix of size $I_m \times R$ that stores the factors of the CPD of \mathcal{X} :

$$\mathbf{A}^{(m)} := \begin{bmatrix} \mathbf{a}_1^{(m)} & \dots & \mathbf{a}_R^{(m)} \end{bmatrix}$$

To compress the expression (2.1), the CPD of \mathcal{X} is written with factor matrices with the following notation:

$$\mathcal{X} = \llbracket \boldsymbol{\lambda}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)} \rrbracket. \quad (2.2)$$

A shorthand notation will be also used for the CPD with $\boldsymbol{\lambda} = \mathbb{1}_R$:

$$\left[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)} \right] = \left[\mathbb{1}_R; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)} \right].$$

Other tensor decompositions exist like the Tucker decomposition [107] or the Block Term Decomposition (BTD) [35]. In this thesis, we will focus only on the CPD. Indeed, we will see in Chapter 3 that the naive Bayes model for the multivariate distributions is tightly linked to the CPD.

Property 2.2.4. Unfoldings and CPD – For an order- M tensor \mathcal{X} that admits a CPD (2.1) of rank R , unfoldings of \mathcal{X} can be expressed in terms of factor matrices:

$$\mathbf{X}_{(m)} = \mathbf{A}^{(m)} \text{Diag}(\boldsymbol{\lambda}) (\mathbf{A}^{(M)} \odot \dots \odot \mathbf{A}^{(m+1)} \odot \mathbf{A}^{(m-1)} \odot \dots \odot \mathbf{A}^{(1)})^T,$$

where $\text{Diag}(\boldsymbol{\lambda})$ is the diagonal matrix which has $\boldsymbol{\lambda}$ as its diagonal.

Definition 2.2.5. Rank of a tensor – For a higher-order tensor \mathcal{X} , the rank of \mathcal{X} (or also the CP rank of \mathcal{X}) is defined as the smallest integer such that eq. (2.1) holds and will be denoted $\text{rank}_{CP}(\mathcal{X})$.

There exist other notions of rank [59] (multilinear, block-term) which are different from the CP-rank. However, in the rest of the thesis only the CPD will be used and the term *rank* will only refer to the CP-rank. For non-negative tensors (tensors with non-negative entries), the following constrained version of the CPD can be introduced.

Definition 2.2.6. Non-negative rank of a tensor – For a higher-order non-negative tensor $\mathcal{X} \in \mathbb{R}_+^{I_1 \times \dots \times I_M}$, the non-negative rank of \mathcal{X} , denoted as $\text{rank}_{CP+}(\mathcal{X})$, is the smallest integer R such that \mathcal{X} admits a non-negative CPD, i.e., the CPD eq. (2.1) with non-negative factor matrices $\mathbf{A}^{(m)} \in \mathbb{R}_+^{I_m \times R}$ (or $\mathbf{A}^{(m)} \geq 0$) and a non-negative loading vector $\boldsymbol{\lambda} \geq 0$.

In general, the non-negative rank can be strictly greater than the rank but there are many cases where the two ranks coincide [83].

2.2.2 Uniqueness of the CP decomposition

The biggest advantage of tensor decompositions over their matrix counterparts is their uniqueness under relatively mild conditions. This is not true for matrices, as the matrix factorization

(without additional constraints) is unique up to a change of basis. For higher-order tensors, however, the uniqueness of the CPD can hold up to trivial ambiguities only [29], which are described below.

First, the permutation ambiguity means that it is possible to change the order of the sum in (2.1) without changing the sum of factors. Indeed, for a tensor \mathcal{X} that admits a CP decomposition (2.2), if we consider a permutation matrix $\mathbf{\Pi} \in \mathbb{R}^{R \times R}$ then \mathcal{X} admits two distinct decompositions:

$$\mathcal{X} = \left[\boldsymbol{\lambda}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)} \right] = \left[\mathbf{\Pi} \boldsymbol{\lambda}; \mathbf{A}^{(1)} \mathbf{\Pi}, \dots, \mathbf{A}^{(M)} \mathbf{\Pi} \right].$$

Second, scaling ambiguities expose that if we consider $M + 1$ scalars $(\alpha_0, \alpha_1, \dots, \alpha_M)$ such that $\prod_{m=0}^M \alpha_m = 1$ (not all equal to 1), it is possible to obtain another CPD for the tensor \mathcal{X} :

$$\mathcal{X} = \sum_{r=1}^R \lambda_r \mathbf{a}_r^{(1)} \circ \dots \circ \mathbf{a}_r^{(M)} = \sum_{r=1}^R (\alpha_0 \lambda_r) (\alpha_1 \mathbf{a}_r^{(1)}) \circ \dots \circ (\alpha_M \mathbf{a}_r^{(M)}).$$

To remove the scaling ambiguities for factors, one may divide each factor by its norm (for example, the 1-norm $\|\cdot\|_1$). By doing that, all factors become norm one and the norms of each factor are transferred to the loading vector $\boldsymbol{\lambda}$. Let $\sigma_r^{(m)}$ denote the norm of the factor $\mathbf{a}_r^{(m)}$, then we have:

$$\mathcal{X} = \sum_{r=1}^R \lambda_r \mathbf{a}_r^{(1)} \circ \dots \circ \mathbf{a}_r^{(M)} = \sum_{r=1}^R \prod_{m=1}^M \sigma_r^{(m)} \lambda_r \frac{1}{\sigma_r^{(1)}} \mathbf{a}_r^{(1)} \circ \dots \circ \frac{1}{\sigma_r^{(M)}} \mathbf{a}_r^{(M)},$$

and the second decomposition does not contain scaling ambiguities for factors (except a possible sign change). We will see in Chapter 3 that, in the case of our study, the factors of the CPD will be already in the normalized form (due to sum-to-one constraints) and the sign ambiguity will be absent (due to non-negativity).

Definition 2.2.7. Uniqueness – For an order- M tensor, the decomposition (2.2) is said unique if there exists no other decompositions up to scaling and permutation ambiguities.

The most classical result is the Kruskal sufficient uniqueness condition for order-3 tensors [61].

Proposition 2.2.8. Kruskal condition [61] – If the matrices $\mathbf{A}^{(1)} \in \mathbb{R}^{I_1 \times R}$, $\mathbf{A}^{(2)} \in \mathbb{R}^{I_2 \times R}$ and $\mathbf{A}^{(3)} \in \mathbb{R}^{I_3 \times R}$ satisfy

$$\kappa(\mathbf{A}^{(1)}) + \kappa(\mathbf{A}^{(2)}) + \kappa(\mathbf{A}^{(3)}) \geq 2R + 2,$$

then, the tensor $\mathcal{X} = \llbracket \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)} \rrbracket$ is of rank R and has a unique CP decomposition.

A generalization of Kruskal's sufficient condition to order- M tensors was provided in [100] and states that a CPD is unique if:

$$2R + (M - 1) \leq \sum_{m=1}^M \kappa(\mathbf{A}^{(m)}).$$

2.2.3 Generic uniqueness

In general, a property is said *generic* if it holds almost everywhere [31]. In this subsection, we introduce the notion of generic uniqueness in the case of tensor decompositions and some classical results around this notion.

Definition 2.2.9. Generic uniqueness – *The CP model of rank R (2.2) is said to be generically unique (or identifiable) if the CPD is unique for all possible factors except a set of factors of Lebesgue measure zero. Equivalently, the CPD is unique (with probability 1) for factor matrices drawn from an absolutely continuous distribution.*

In other words, if generic uniqueness is ensured, the probability of drawing random factors such that the CPD is not unique is zero.

By using the generalized Kruskal condition (2.2.2) in the case of cubic tensors, the CPD is generically unique if

$$2R + (M - 1) \leq \min(I, R)M.$$

This expression was obtained because $\kappa(\mathbf{A}^{(m)}) = \min(I, R)$ for factor matrices drawn from an absolutely continuous distribution. In [26], milder sufficient conditions were proved in the case of tensors of order 3, stating that the CPD is generically unique if

$$R \leq 4^{\log_2(I)-1}.$$

For order- M tensors, the result of [26] can be extended by grouping variables in 3 disjoint groups [74, Theorem 2]. However, there are more powerful results which are related to the notion of generic rank.

Definition 2.2.10. Typical ranks – *For a tensor space $\mathbb{R}^{I_1 \times \dots \times I_M}$, a typical rank is a rank that appears with non-zero probability when the entries are drawn randomly from an absolutely continuous distribution. This means that the set of tensors of non-typical ranks is of Lebesgue measure zero.*

Remark 2.2.11. *For tensors with complex entries, there always exists one typical rank called generic rank which is equal to the smallest real typical rank [12] for real tensors of same dimensions.*

There are many recent results on generic rank and generic uniqueness and in many cases, the CP models for ranks below the generic ranks are identifiable [27, 83]. In particular, the following result will be useful.

Proposition 2.2.12 ([27, Theorem 1.1][83, Corollary 30]). *Let $I_1 \geq I_2 \geq I_3$ and*

$$R_0 = \left\lceil \frac{I_1 I_2 I_3}{I_1 + I_2 + I_3 - 2} \right\rceil.$$

Then the CP model for tensors of size $I_1 \times I_2 \times I_3$ is identifiable for $R < R_0$ if $I_1 I_2 I_3 \leq 15000$, and if (I_1, I_2, I_3, R) do not fall into the following exceptional cases:

(I_1, I_2, I_3)	R
$(4, 4, 3)$	5
$(4, 4, 4)$	6
$(6, 6, 3)$	8
$I_1 > I_2 I_3 - (I_2 - 1) + (I_3 - 1)$	$R \geq I_2 I_3 - (I_2 - 1) - (I_3 - 1)$

Note that in the cubic case ($I_1 = I_2 = I_3 = I$), the identifiability results go back to Strassen [103], and identifiability is proved without constraints on I for a slightly weaker condition than in Proposition 2.2.12 (see discussion in [27]).

Remark 2.2.13. *Many identifiability results in the literature are proved for complex-valued tensors. However, as shown in [83], complex identifiability implies real identifiability (and also non-negative identifiability) for ranks below the generic rank.*

2.3 Identifiability of polynomial models

In this thesis, we consider CP models which are coupled and whose factors are constrained to be non-negative and sum-to-one (see Chapters 3 and 4). For these reasons, the uniqueness results presented in the previous section cannot be directly applied. In this section, we recall results on identifiability of polynomial models from [15] which use the tools of algebraic and semi-algebraic geometry [83]. The results presented in this section will be used later in Chapter 5 to study the identifiability of the coupled models considered in the thesis.

2.3.1 Polynomial mapping and uniqueness

Let μ be a polynomial mapping from a space of n parameters $\Theta \subset \mathbb{R}^n$ to a space of S observations (or measurements) denoted $\mathcal{S} \subset \mathbb{R}^S$. For simplicity, we will refer to (μ, Θ) as the model. In the following, we assume that Θ is of positive measure and contains an open (Euclidean) subset.

Example 2.3.1. *The CP model can be viewed as a special case of polynomial models. In this case,*

$$\mu(\boldsymbol{\theta}) = \left[\boldsymbol{\lambda}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)} \right],$$

where for example $\boldsymbol{\theta}$ is defined as

$$\boldsymbol{\theta} = (\boldsymbol{\lambda}, \text{vec}(\mathbf{A}^{(1)}), \dots, \text{vec}(\mathbf{A}^{(M)})).$$

The choice of $\Theta = \mathbb{R}^{(IM+1)R}$ corresponds to the real CP model while $\Theta = \mathbb{R}_+^{(IM+1)R}$ corresponds to the non-negative CP model (see Definition 2.2.6).

One of the key notions for studying model identifiability is the notion of recoverability [15].

Definition 2.3.2. Recoverability – *The model (μ, Θ) is said recoverable at a given $\boldsymbol{\theta} \in \Theta$ if there exists a finite number of elements in the pre-image of $\mu^{-1}(\mu(\boldsymbol{\theta}))$. The model (μ, Θ) is generically recoverable if it is recoverable for a generic $\boldsymbol{\theta} \in \Theta$ (that is for all $\boldsymbol{\theta}$ except a subset of Lebesgue measure 0).*

To study the recoverability, the following tool is crucial:

Definition 2.3.3. Jacobian matrix of a parametrization – *For a polynomial model (μ, Θ) , the Jacobian of the parametrization at a given $\boldsymbol{\theta} = [\theta_1 \ \dots \ \theta_n]^\top \in \Theta$, is a $S \times n$ matrix defined by:*

$$\mathcal{J}_\mu(\boldsymbol{\theta}) := \left(\frac{\partial y_s}{\partial \theta_i} \right)_{s=1, i=1}^{S, n},$$

where \mathbf{y} is the vector such that $\mathbf{y} = \mu(\boldsymbol{\theta}) = [y_1(\boldsymbol{\theta}) \ \dots \ y_S(\boldsymbol{\theta})]^\top$.

The following proposition shows that the generic recoverability is determined by the rank of the Jacobian.

Proposition 2.3.4. *Special case of [15, Theorem 4.9] – The following two statements are equivalent:*

1. There exists a $\boldsymbol{\theta}^* \in \mathbb{R}^n$ such that the $\mathcal{J}_\mu(\boldsymbol{\theta})$ is full column rank.
2. The model (μ, \mathbb{R}^n) is generically recoverable.

Remark 2.3.5. In the previous proposition, if at a single point $\boldsymbol{\theta}^*$ the Jacobian is full column rank, then it is full rank for a generic point $\boldsymbol{\theta} \in \mathbb{R}^n$ (almost everywhere on \mathbb{R}^n). This follows from [83, Lemma 1].

Note that the Proposition 2.3.4 cannot be directly applied to tensor decomposition models due to scaling ambiguities (which implies the rank-deficiency of the Jacobian matrix). To apply this proposition, a re-parametrization with a smaller parameter space can be used (see Section 5.1.2).

2.3.2 Uniqueness and recoverability of additive models

The CP model (2.1), as well as the coupled models considered in the thesis, falls into a special case of additive polynomial models.

Definition 2.3.6. Additive model – Let us consider a polynomial model (μ_1, Θ_1) where $\Theta_1 \subset \mathbb{R}^{n_1}$. Then, the R -term additive model (μ, Θ) is defined as

$$\mu(\boldsymbol{\theta}) = \mu_1(\boldsymbol{\theta}_1) + \dots + \mu_1(\boldsymbol{\theta}_R), \quad (2.3)$$

where $\Theta = \Theta_1^R$ and $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_R)$, $\boldsymbol{\theta}_k \in \Theta_1$.

Example 2.3.7. The CP model of Example 2.3.1 can be seen as an additive polynomial model after reordering the parameters such that

$$\boldsymbol{\theta}_r = (\lambda_r, \mathbf{a}_r^{(1)}, \dots, \mathbf{a}_r^{(M)}),$$

and with μ_1 the mapping is given as

$$\mu_1(\boldsymbol{\theta}_r) = \lambda_r \mathbf{a}_r^{(1)} \circ \dots \circ \mathbf{a}_r^{(M)}. \quad (2.4)$$

Additive polynomial models generalize tensor decompositions and can be studied in the framework of X-rank decompositions [83]. Similarly to tensor decompositions, uniqueness of additive models can be defined.

Definition 2.3.8. Essential uniqueness of additive models – Let $\mathbf{y} \in \mathcal{S}$ admitting an R -term decomposition (2.3) such that $\mathbf{y} = \mu(\boldsymbol{\theta})$. Then, the model (μ, Θ) is said to be essentially unique at a given $\boldsymbol{\theta} \in \Theta$ if elements in μ^{-1} are equal to $\boldsymbol{\theta}$ up to a permutation of terms in (2.3), and change of parameters $\boldsymbol{\theta}_r$ that do not change the one-term result $\mu_1(\boldsymbol{\theta}_r)$.

The model is called identifiable if it is essentially unique for a generic $\boldsymbol{\theta} \in \Theta$.

Note that the essential uniqueness in sense of Definition 2.3.8 corresponds exactly to uniqueness in the case of tensor decompositions (see Definition 2.2.7), as scaling is the only ambiguity that can result in the same rank-one tensor in Equation (2.4).

Remark 2.3.9. Recoverability is a necessary condition for identifiability. Indeed, for an identifiable additive model (μ, Θ) (see Definition 2.3.8) with recoverable (μ_1, Θ_1) , the model is recoverable in the sense of Definition 2.3.2.

Thanks to Remark 2.3.9, the Jacobian can be used to study identifiability and recoverability of additive models. Indeed, if the rank of Jacobian is maximal at a given point $\boldsymbol{\theta} \in \Theta$ and equal to

$$\text{rank}(\mathcal{J}_\mu(\boldsymbol{\theta})) = Rn_1, \quad (2.5)$$

then the model is recoverable by Proposition 2.3.4.

Note that the Jacobian matrix has the following block structure:

$$\mathcal{J}_\mu(\boldsymbol{\theta}) = \begin{bmatrix} \mathcal{J}_{\mu_1}(\boldsymbol{\theta}_1) & \cdots & \mathcal{J}_{\mu_1}(\boldsymbol{\theta}_R) \end{bmatrix} \quad (2.6)$$

where $\mathcal{J}_{\mu_1}(\boldsymbol{\theta}_r)$ represents the Jacobian of the rank-1 parametrization μ_1 applied to the r -th block of parameters $\boldsymbol{\theta}_r$. If Equation (2.5) holds, then any subsets of block columns of the Jacobian is also full column rank. This means that the same R' -term models $R' < R$ are also generically recoverable. Let us denote R_{\max} the integer such that (2.5) holds for $R \leq R_{\max}$ and does not hold for $R > R_{\max}$. We call such R_{\max} the recoverability bound because the model (μ, \mathbb{R}^n) is necessary not recoverable if the Jacobian is not full column rank. Note that in order to apply the Equation (2.5) to tensor CP models, a re-parametrization that removes the scaling ambiguities of factors is needed and will be presented in Section 5.1.1. The following lemma will be also useful.

Lemma 2.3.10. Let (μ, Θ) an identifiable (respectively recoverable) additive model and let $\tilde{\Theta}$ be a subset of Θ . If $\tilde{\Theta}$ contains an open subset, then $(\mu, \tilde{\Theta})$ is identifiable (respectively recoverable).

Proof. We provide a proof for the case of identifiability. Let \mathcal{Y} be the subset of Θ that contains the non-identifiable parameters. Because (μ, Θ) is identifiable, \mathcal{Y} is of Lebesgue measure zero. Then, $\mathcal{Y} \cap \tilde{\Theta}$ is also of Lebesgue measure zero which implies identifiability of $(\mu, \tilde{\Theta})$. \square

2.4 Conclusion of Chapter 2

In this chapter, the necessary background on tensor decompositions was presented. The CP decomposition which permits to express a tensor as a sum of rank-one factors was introduced. The key properties of the CPD such as uniqueness were recalled. This decomposition is strongly linked to the naive Bayes model, a model for probability density functions introduced in Chapter 3 and later applied to the flow cytometry data.

Finally, tensor decompositions also rely on powerful uniqueness results. Some classical results for the CP decomposition were provided in this chapter. In addition to that, polynomial and additive models were introduced. This framework will be used in the context of low rank tensor decompositions in Chapter 5, both in terms of identifiability and recoverability.

Chapter 3

Histograms for probability density function estimation

Contents

3.1	Density estimation	48
3.1.1	Problem statement	48
3.1.2	Parametric approaches	49
3.1.3	Non-parametric approaches	50
3.2	Histograms	51
3.2.1	Definitions	51
3.2.2	Curse of dimensionality: choosing the number of bins	55
3.3	Naive Bayes model and CP decomposition	57
3.3.1	Graphical models	57
3.3.2	Naive Bayes model	59
3.3.3	Link between naive Bayes model and tensor decompositions	61
3.3.4	Curse of dimensionality: estimating the factors of the NBM	63
3.4	Fully coupled tensor factorization of 3D marginals	64
3.4.1	Marginalization of naive Bayes model	64
3.4.2	Fully coupled tensor factorization	66
3.4.3	Curse of dimensionality: number of marginals	67
3.5	Conclusion of Chapter 3	68

This chapter aims at describing the main methodological problems and challenges of coupled decompositions of marginals for density estimation which are addressed in the thesis. First, an overview of classical parametric and non-parametric methods is provided. A focus is given to histograms as non-parametric estimators which suffers from the curse of dimensionality, as the complexity of histogram estimation grows exponentially with dimensions.

To tackle this problem, we resort to a particular distribution model, the so-called naive Bayes model, which is a special instance of graphical models. The naive Bayes model results in a mixture model that can be seen as a constrained low-rank tensor decomposition whose complexity remains linear with dimensions. We also present the 3D-marginal-based approach initiated in [74] that further reduces the complexity of the histogram model by a coupled CP decomposition of 3D marginals.

3.1 Density estimation

3.1.1 Problem statement

Let us consider M random variables grouped in a random vector $\mathbf{x} = (X^{(1)}, \dots, X^{(M)})$ having an absolutely continuous distribution. For $m \in \llbracket 1, M \rrbracket$, the variable $X^{(m)}$ takes values in a subset $\mathcal{I}^{(m)} \subset \mathbb{R}$. The problem of joint distribution estimation is to estimate the joint PDF of all variables

$$p_{\mathbf{x}} : \begin{cases} \mathcal{I}^{(1)} \times \dots \times \mathcal{I}^{(M)} & \longrightarrow & \mathbb{R} \\ (x_1, \dots, x_M) & \longmapsto & p_{\mathbf{x}}(x_1, \dots, x_M) \end{cases}$$

from N samples which are supposed independent and identically distributed. Each sample is a realization of the random vector \mathbf{x} and can be stored as the n -th row of an observation matrix $\mathbf{X} \in \mathbb{R}^{N \times M}$.

Definition 3.1.1. Multivariate Probability Density Function – For a random vector $\mathbf{x} = (X^{(1)}, \dots, X^{(M)})$ taking values in \mathbb{R}^M and having an absolutely continuous distribution (with respect to Lebesgue measure), its density (denoted $p_{\mathbf{x}}$) is the M -dimensional function such that for all $m \in \llbracket 1, M \rrbracket$, $[a_m, b_m] \subset \mathbb{R}$:

$$\Pr\left(X^{(1)} \in [a_1, b_1], \dots, X^{(M)} \in [a_M, b_M]\right) = \int_{a_1}^{b_1} \dots \int_{a_M}^{b_M} p_{\mathbf{x}}(x_1, \dots, x_M) dx_1 \dots dx_M.$$

The problem of probability density function (PDF) estimation arises in numerous signal processing applications. For example, knowing the joint PDF of M random variables permits

to cluster samples thanks to the Maximum a Posteriori (MAP) principle [39, 110]. Other data mining and signal processing application examples can be found with sensor networks [75], biology [51] or social network studies [76], to name a few.

Without any assumptions, the task of PDF estimation is considered impossible in practice. There are many approaches to density estimation which will be briefly reviewed in the following subsections. For most of these approaches, the biggest challenge is the curse of dimensionality which will be explained with the framework of histograms.

3.1.2 Parametric approaches

One of the most common parametric approaches is the Gaussian Mixture Model (GMM), introduced in 1894 by Pearson in [81]. The idea was to model a univariate density into a sum of 2 Gaussians. The model was then described by the components means μ_1 and μ_2 , variances σ_1 and σ_2 and the proportion of each component λ_1 and λ_2 such that $\lambda_1 + \lambda_2 = 1$. In the general case, the PDF is modeled as a mixture of R multivariate Gaussian densities $\{f_r\}_{r=1}^R$:

$$p_{\mathbf{x}}(x_1, \dots, x_M) = \sum_{r=1}^R \lambda_r f_r(x_1, \dots, x_M), \quad (3.1)$$

where $\lambda_r > 0$ for all $r \in \llbracket 1, R \rrbracket$ and

$$\sum_{r=1}^R \lambda_r = 1.$$

The component densities $\{f_r\}_{r=1}^R$ can be written

$$f_r(x_1, \dots, x_M) = \frac{1}{\sqrt{\det(\Sigma_r)(2\pi)^M}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_r)^\top \Sigma_r^{-1}(\mathbf{y} - \boldsymbol{\mu}_r)\right), \quad (3.2)$$

where $\Sigma_r \in \mathbb{R}^{M \times M}$ denotes the r -th covariance matrix and $\boldsymbol{\mu}_r \in \mathbb{R}^{M \times 1}$ denotes the r -th vector of means and \mathbf{y} stands for the vector $[x_1 \dots x_M]^\top$ ². An example is given in Figure 3.1 where a bivariate ($M = 2$) distribution is modeled with a GMM of $R = 4$ components which can be plotted to show the mixture of the different components.

To obtain the parameters $(\{\boldsymbol{\mu}_r\}_{r=1}^R, \{\Sigma_r\}_{r=1}^R, \boldsymbol{\lambda})$ of a GMM, one common approach is the maximum likelihood estimation (for example, using the Expectation-Maximization (EM) algorithm [36]). In general, finite mixtures (3.1) of other density functions (like exponential or

²It cannot be called \mathbf{x} as it is already the random vector of interest. \mathbf{x} stores the random variables while \mathbf{y} contains the values x_m for which the density is computed.

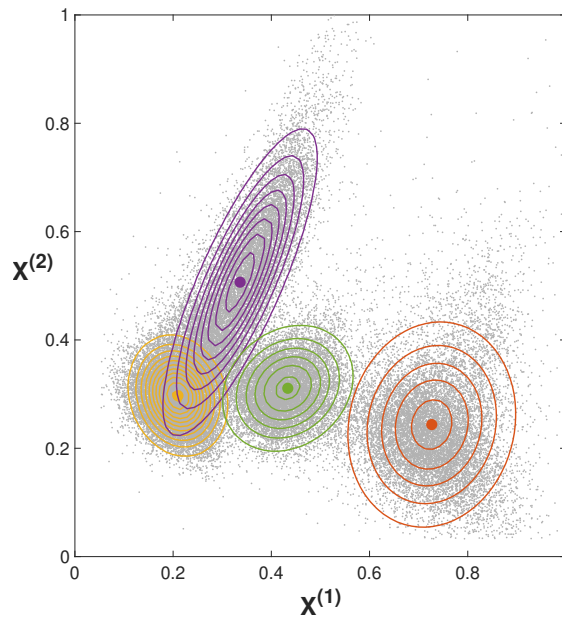


Figure 3.1: Example of a GMM for a bivariate distribution featuring the contour distribution of $R = 4$ components.

gamma distributions for example) can be considered and estimated with this framework. It can also be mentioned that parametric distribution estimation can be performed via Markov Chain Monte Carlo (MCMC) methods such as Gibbs sampling (see [71] for more details).

The first drawback of parametric approaches is that a parametric model for the density is a strong assumption in practice. For example, it has been shown that the GMM suffers from a lack of uniqueness of components in the mixture, as it is possible to obtain the same sum (3.1) with different sets of parameters [28]. Another drawback is the hardness to solve the estimation problems to obtain the parameters of a model. For example, when EM is used, it is very dependent on the initialization and can be stuck easily in local maxima [116]. EM is also known to have a slow convergence rate. Even if there exist acceleration methods for the EM algorithm such as SQUAREM [111], those approaches still does not scale for datasets with a large number of samples (see Section 4.3.3). In our work, we focus on non-parametric approaches which prevents us to consider a too restrictive model on our data.

3.1.3 Non-parametric approaches

One of the oldest non-parametric approaches is histogram estimation which tries to estimate the PDF as a piece-wise constant function by counting the samples that falls into bins (multidi-

mensional intervals) [94]. Histograms can be viewed as a special case of more general kernel density estimation methods [94]. The principle of kernel density estimation [97, 101] is to consider the density as a sum of weighted functions called kernels. Also introduced by Scott in [95], frequency polygons estimate continuous distributions "by connecting with straight lines the mid-bin values of a histogram". There exist other non-parametric approaches to density estimation such as log-concave modelling proposed [33], but it can model only unimodal distributions and its complexity is prohibitive for large datasets. Finally, let us mention recent density estimation techniques based on neural networks [62, 108], but these methods are still poorly understood.

In this thesis, we focus our work on histograms for multiple reasons. First, as we will see in the following sections, multivariate histograms can be represented as tensors. Thus, it is possible to use the tensor framework to tackle the problem of density estimation. As we shall explain in detail, low-rank tensor decompositions permit to break the curse of dimensionality arising in multivariate density estimation. Moreover, histograms are used by FCM end users on a daily basis. Therefore, this ensures that the proposed methodology can be easily plugged into state-of-the-art FCM analysis pipelines.

3.2 Histograms

In the literature, histograms are defined in two possible ways: from data science and density estimation points of view. From the first viewpoint, histograms are described just a descriptive statistics or a visualization technique for a dataset. We first use this perspective in Section 3.2.1 to define histograms in univariate and multivariate cases. The density estimation viewpoint will be used in Section 3.2.2 to explain how the number of bins can be chosen and to show that histogram-based density estimation also suffers from the curse of dimensionality.

3.2.1 Definitions

Univariate histogram

Let us first consider a real random variable X and let $\{X_n\}_{n=1}^N$ be a set of N realizations of X . The histogram is defined on an interval $\mathcal{I} = [\delta_0, \delta_I]$ partitioned into I disjoint intervals which are defined by $I + 1$ endpoints:

$$\delta_0 < \delta_1 < \dots < \delta_{I-1} < \delta_I.$$

Remark 3.2.1. From a data science perspective [114], the edges are often chosen as $\delta_0 = \min_{n \in \llbracket 1, N \rrbracket} \{X_n\}$ and $\delta_I = \max_{n \in \llbracket 1, N \rrbracket} \{X_n\}$.

In the following, we will consider only the case of the uniform discretization³ of \mathcal{I} . Note that in the FCM data analysis, due of the non-linear transformation presented in Section 1.2.4, the uniform discretization of \mathcal{I} can be viewed as a non-uniform discretization in the non-transformed space. The endpoints permit to define the intervals $\Delta_i, i \in \llbracket 1, I \rrbracket$ which are also called *bins*:

$$\Delta_i := \begin{cases} [\delta_{i-1}, \delta_i), & i < I, \\ [\delta_{I-1}, \delta_I], & \text{otherwise.} \end{cases} \quad (3.3)$$

The principle of density estimation with a histogram is to estimate the so-called *bin probabilities*:

$$h_i := \int_{\delta_{i-1}}^{\delta_i} p_X(x) dx = \Pr(X \in [\delta_{i-1}, \delta_i)).$$

Definition 3.2.2. Theoretical histogram – For a random variable X and a set of bins (3.3), the theoretical histogram is defined as the set of bin probabilities:

$$\{h_i \mid i \in \llbracket 1, I \rrbracket\}.$$

This set of bin probabilities can be stored in a 1-dimensional array (a vector) which will be denoted \mathbf{h} :

$$\mathbf{h} = \begin{bmatrix} h_1 & \dots & h_I \end{bmatrix}^\top.$$

Remark 3.2.3. In the following, we will consider that the sum of bin probabilities is 1 which means that $\mathbb{1}_I^\top \mathbf{h} = 1$. This is equivalent to assuming that the support of the random variable is contained on \mathcal{I} . If it is not the case, we can consider the random variable Y distributed according to the conditional density $p_{(X \mid X \in \mathcal{I})}$, then the sum of bin probabilities for the variable Y is 1.

Second, estimates of $h_i, i \in \llbracket 1, I \rrbracket$, are obtained by counting the number of samples that fall into each bin:

$$\tilde{h}_i = \frac{1}{N} \text{Card}\{n \in \llbracket 1, N \rrbracket \mid X_n \in \Delta_i\}.$$

³It is also possible to consider non-uniform discretization. Indeed, similarly to quantization in signal processing, non-uniform discretization can be chosen such that more bins are placed where more data points are present.

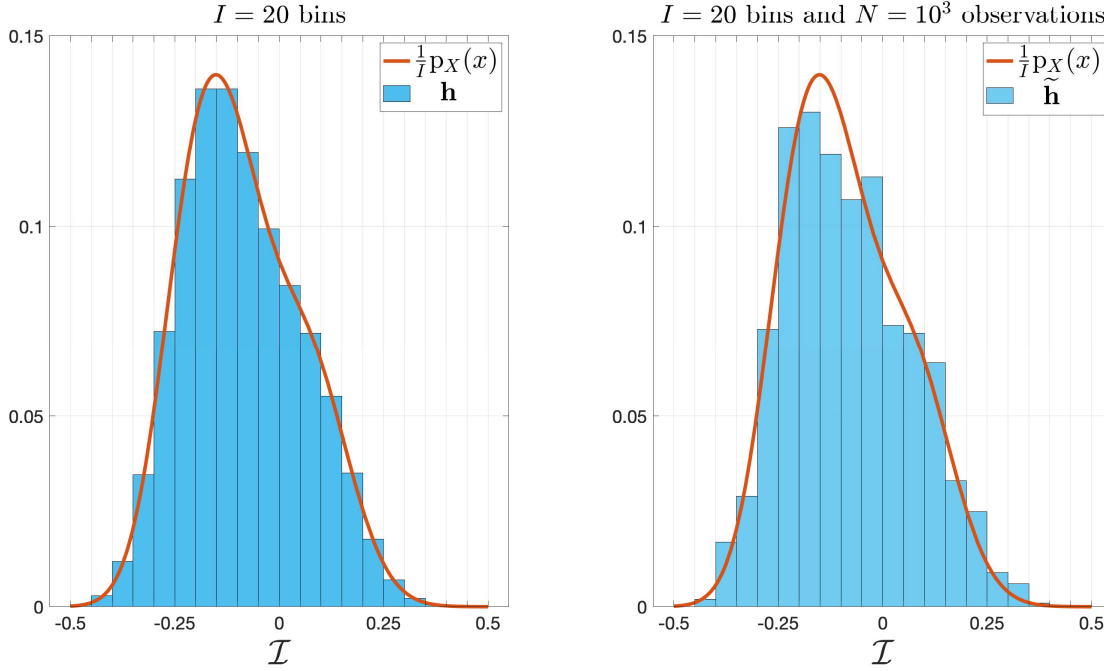


Figure 3.2: A theoretical histogram (**Left plot**) and an empirical histogram (**Right plot**) as an estimation of the scaled PDF.

Definition 3.2.4. Empirical histogram – For a set of N observations $\{X_n\}_{n \in \llbracket 1, N \rrbracket}$ of a random variable X and for a set of bins (3.3), the empirical histogram is defined as the set of estimates of bin probabilities :

$$\{\tilde{h}_i \mid i \in \llbracket 1, I \rrbracket\},$$

which can be stored in a 1-dimensional array (a vector) which will be denoted $\tilde{\mathbf{h}}$:

$$\tilde{\mathbf{h}} = [\tilde{h}_1 \quad \dots \quad \tilde{h}_I]^\top.$$

Figure 3.2 shows for a univariate density both a theoretical histogram (for $I = 20$ bins) and the associated empirical histogram obtained for $N = 10^3$ observations.

Multivariate histogram

For multivariate distributions, histograms are defined in a similar way. Let us consider a vector \mathbf{x} of M random variables $X^{(m)}$ and let $\mathbf{X} \in \mathbb{R}^{N \times M}$ be an observation matrix whose rows contains a set of N realizations of \mathbf{x} . First, we define for each variable an interval $\mathcal{I}^{(m)}$ that is divided in

I 1-dimensional bins denoted as $\Delta_{i_m}^{(m)}$, $i_m \in \llbracket 1, I \rrbracket$. With this notation, the bin probabilities are denoted as $h_{i_1 \dots i_M}$ and are defined as

$$\begin{aligned} h_{i_1 \dots i_M} &:= \int_{\Delta_{i_1}^{(1)}} \cdots \int_{\Delta_{i_M}^{(M)}} p_{\mathbf{x}}(x_1, \dots, x_M) dx_1 \dots dx_M \\ &= \Pr(X^{(1)} \in \Delta_{i_1}^{(1)}, \dots, X^{(M)} \in \Delta_{i_M}^{(M)}). \end{aligned} \quad (3.4)$$

Note that in the multivariate case, the notion of *bin* denotes the Cartesian product of M one-dimensional bins $\Delta_{i_1}^{(1)} \times \cdots \times \Delta_{i_M}^{(M)}$.

Definition 3.2.5. Theoretical histogram tensor – For a random vector \mathbf{x} and a set of intervals $\Delta_{i_m}^{(m)}$, the theoretical histogram tensor is defined as the set of bin probabilities:

$$\{h_{i_1 \dots i_M} \mid \forall m \in \llbracket 1, M \rrbracket, i_m \in \llbracket 1, I \rrbracket\}.$$

This set of bin probabilities can be stored in a M -dimensional array – an order- M tensor – which will be denoted \mathcal{H} :

$$\mathcal{H} = (h_{i_1 \dots i_M})_{i_1=1, \dots, i_M=1}^{I, \dots, I}.$$

To estimate the $h_{i_1 \dots i_M}$ from a set of N realizations stored in an observation matrix $\mathbf{X} \in \mathbb{R}^{N \times M}$, as in the 1D case we count the samples that fall into each bin

$$\tilde{h}_{i_1 \dots i_M} = \frac{1}{N} \text{Card} \left\{ n \in \llbracket 1, N \rrbracket \mid \mathbf{X}_n \in \Delta_{i_1}^{(1)} \times \cdots \times \Delta_{i_M}^{(M)} \right\}. \quad (3.5)$$

Definition 3.2.6. Empirical histogram tensor – For a set of N observations $\{\mathbf{X}_n\}_{n \in \llbracket 1, N \rrbracket}$ of a random vector \mathbf{x} and a set of intervals $\Delta_{i_m}^{(m)}$, the empirical histogram tensor is defined as the set of estimates of bin probabilities:

$$\{\tilde{h}_{i_1 \dots i_M} \mid \forall m \in \llbracket 1, M \rrbracket, i_m \in \llbracket 1, I \rrbracket\},$$

which can be stored in a M -dimensional array – an order- M tensor – which will be denoted $\tilde{\mathcal{H}}$:

$$\tilde{\mathcal{H}} = (\tilde{h}_{i_1 \dots i_M})_{i_1=1, \dots, i_M=1}^{I, \dots, I}.$$

Similarly to univariate histograms, we will assume that the sum of all entries of \mathcal{H} is equal to 1. In the following, multivariate histograms \mathcal{H} and $\tilde{\mathcal{H}}$ are cubic tensors because the same number of bins is chosen for all variables ($I_m = I$ for all $m \in \llbracket 1, M \rrbracket$).

3.2.2 Curse of dimensionality: choosing the number of bins

A crucial question in histogram estimation is the choice of the number of bins I with respect to the number of samples N . In this subsection, the performance of the histogram estimator is studied from a density estimation perspective. In the case of a random variable X taking values in \mathbb{R} , the density p_X is estimated with the *empirical histogram estimator* defined as

$$\tilde{f}_X(x) = \sum_{i \in \mathbb{Z}} \frac{\tilde{h}_i}{w} \mathbb{1}_{[\delta_{i-1}, \delta_i)}(x), \quad (3.6)$$

where w is the bin width and where $\mathbb{1}_{[\delta_{i-1}, \delta_i)}$ is the indicator function of the set $[\delta_{i-1}, \delta_i)$:

$$\mathbb{1}_{[\delta_{i-1}, \delta_i)}(x) = \begin{cases} 1 & \text{if } x \in [\delta_{i-1}, \delta_i), \\ 0 & \text{otherwise} \end{cases}.$$

In the following, bins have same width hence the bin width is a constant equal to $w = \delta_i - \delta_{i-1}$ for all $i \in \mathbb{Z}$. A typical criterion consists in minimizing the Mean Integrated Squared Error (MISE) defined as

$$\text{MISE} = \int_{-\infty}^{+\infty} \mathbb{E}[(\tilde{f}_X(x) - p_X(x))^2] dx. \quad (3.7)$$

In the Gaussian case ($X \sim \mathcal{N}(\mu, \sigma^2)$), Scott [96] established the following optimal bin width rule

$$w^* = \frac{3.5\sigma}{\sqrt[3]{N}}. \quad (3.8)$$

Equation (3.8) is usually considered as a generic rule, i.e., it provides a reasonable trade-off for arbitrary densities beyond the Gaussian case.

However, in practice, histograms are computed on a finite number of bins. Therefore, before picking the width of a histogram, one must determine the finite support on which the histogram will be computed. This interval will be denoted by \mathcal{I} . Consider the restriction of the histogram estimator \tilde{f}_X to \mathcal{I} given by

$$\tilde{f}_{X,\mathcal{I}}(x) = \tilde{f}_X(x) \mathbb{1}_{\mathcal{I}}(x).$$

Then, for any $\varepsilon > 0$, it is always possible to find an interval $\mathcal{I} = [\delta_-, \delta_+]$ such that

$$\int_{-\infty}^{+\infty} \mathbb{E}[(\tilde{f}_{X,\mathcal{I}}(x) - \tilde{f}_X(x))^2] dx < \varepsilon, \quad (3.9)$$

i.e., the truncated histogram $\tilde{f}_{X,\mathcal{I}}$ approximates arbitrary well the histogram \tilde{f}_X for an interval \mathcal{I}

“large” enough. For instance, in the zero-mean Gaussian case, a typical choice of $\mathcal{I} = [-5\sigma, 5\sigma]$ ensures that the condition (3.9) is satisfied for a small value of ε . More generally, if one selects a sufficiently large interval \mathcal{I} with length $S = |\mathcal{I}|$, then it is possible to adapt the optimal bin width rule (3.8) to select the optimal number of bins I^* for the interval \mathcal{I} . Since $I^*w^* = S$, plugging (3.8) yields

$$I^* = \left\lceil \frac{S\sqrt[3]{N}}{3.5\sigma} \right\rceil. \quad (3.10)$$

The rule (3.10) can be extended to the case of multivariate histograms. Indeed, Scott [96] also provides an optimal bin width rule in the multivariate Gaussian case. For $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$ where $\text{diag}(\Sigma) = (\sigma_1^2, \dots, \sigma_M^2)$, the optimal bin width for each variable is defined by

$$w_m^* = \frac{3.5\sigma_m}{2+M\sqrt{N}}. \quad (3.11)$$

Analogously to the univariate case, this rule can be adapted to the number of bins by considering a sufficiently large (regarding (3.9)) interval $\mathcal{I}^{(m)}$ of length S_m :

$$I_m^* = \left\lceil \frac{S_m^{2+M}\sqrt{N}}{3.5\sigma_m} \right\rceil. \quad (3.12)$$

Remark 3.2.7. *Although it seems that the rule (3.12) provides different numbers of bins per dimension, it is not the case in practice. Indeed, if one chooses an interval $\mathcal{I}^{(m)} = [-5\sigma_m, 5\sigma_m]$ for all $m \in \llbracket 1, M \rrbracket$, the ratio between S_m and σ_m is a constant (strongly linked to the tolerance of Equation (3.9)). Therefore, it leads to $I^* = I_m^*$ for $m \in \llbracket 1, M \rrbracket$. In the case, if one want to estimate the histogram for I bins per dimension, it requires the following optimal number of samples N^* :*

$$N^* = \lceil (0.35I)^{M+2} \rceil.$$

Curse of dimensionality: number of samples

One of the biggest challenges around density estimation is the curse of dimensionality [11] which denotes the fact that the complexity of the problem increases exponentially with M . Although this effect occurs with all approaches previously presented, the curse of dimensionality will be explained for the case of the histogram-based density estimation framework. Table 3.1 shows that estimating a multivariate histogram requires a number of samples growing exponentially with the order of the histogram. Indeed, if we take the example of a joint density of $M = 10$ variables, it is difficult in practice to have datasets containing billions of samples (following the

Table 3.1: Optimal number of samples N (see Equation (3.12)) required to estimate a multivariate normal distribution ($\mu_m = 0$ and $\sigma_m = 0.1$ for all $m \in \llbracket 1, M \rrbracket$) regarding M and I^* .

	$I^* = 5$	$I^* = 10$	$I^* = 20$
$M = 1$	6	43	343
$M = 2$	10	151	2401
$M = 3$	17	526	16807
$M = 5$	51	6434	823543
$M = 10$	826	3.3×10^6	1.3×10^{10}

rule (3.12)). In extreme cases, the required number of samples can increase to a point where $N < I^M$; meaning that only a subset of the estimated histogram values will be non-zeros. In addition to the lack of samples, visualizing and interpreting a higher-order histogram is challenging in practice.

3.3 Naive Bayes model and CP decomposition

3.3.1 Graphical models

Graphical models are a class of probabilistic models which encode the dependencies between variables by a graph [57]. In a graphical model, nodes of a graph (random variables) are connected by edges that represent statistical dependencies between variables. By exploiting the statistical independencies between variables, graphical models permit to reduce the model complexity [37], hence they can be used to address the curse of dimensionality. The two main types of graphical models are undirected and directed graphical models, based respectively on undirected and directed acyclic graphs. In this subsection, we provide expressions for absolutely continuous distributions. However, these expressions can be easily transposed to the case of random discrete variables.

In the undirected case, edges represent statistical dependencies between variables. In this case, the density can be factorized as a product of functions depending on smaller number of variables. For example, for the graph presented in Figure 3.3b, the density is factorized as

$$p_{\mathbf{x}}(x_1, x_2, x_3, x_4) = p_{(1,2)}(x_1, x_2)p_{(3,4)}(x_3, x_4),$$

where $p_{(1,2)}$ and $p_{(3,4)}$ are marginal distributions densities. The graph of Figure 3.3a corre-

sponds to an independence model hence its density can be factorized as

$$p_{\mathbf{x}}(x_1, x_2, x_3, x_4) = p_1(x_1)p_2(x_2)p_3(x_3)p_4(x_4).$$

More details and definitions on undirected graphical models are available in [57, 63]. Note that for discrete random variables, undirected graphical models are linked to tensor hyper-networks [86].

In the directed case, the density admits a factorization as a product of conditional densities based on the structure of an directed acyclic graph.

Definition 3.3.1. Directed graph – For a set of vertices $\mathcal{V} = \llbracket 1, M \rrbracket$, a directed graph is a tuple $(\mathcal{V}, \mathcal{E})$ where \mathcal{E} is a set of directed edges. An edge $\varepsilon = (m_1, m_2) \in \mathcal{E}$ contains a starting vertex m_1 and an ending vertex m_2 .

A directed graph can be represented similarly to undirected graphs, with the difference of edges are plotted as arrows pointing the ending vertex. A directed graph is called *acyclic* if it is not possible to find a sequence of vertices $(m_1, \dots, m_N) \in \mathcal{V}^N$ with edges that form a loop:

$$\{(m_1, m_2), (m_2, m_3), \dots, (m_{N-1}, m_N), (m_N, m_1)\} \subset \mathcal{E}.$$

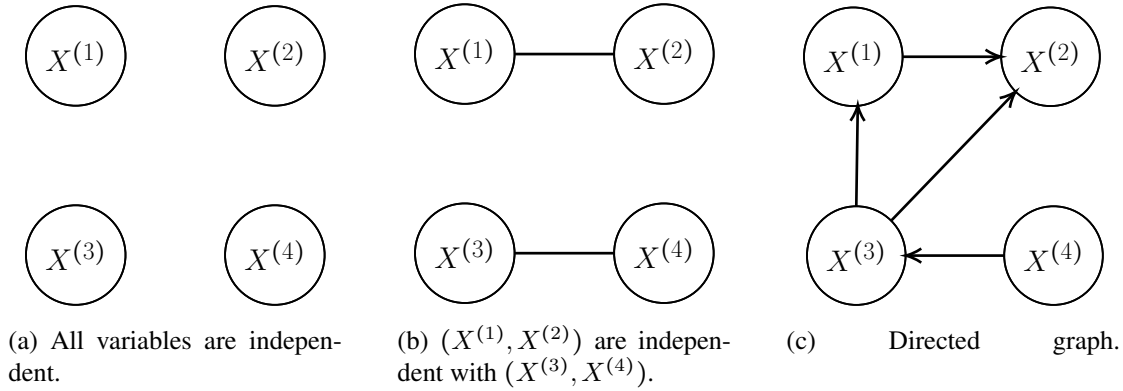
In an acyclic directed graph, a starting vertex of an edge is also called *parent* for the ending vertex. An example of an acyclic directed graphical model is given in Figure 3.3c. In this example, the 4 variables $\{X^{(m)}\}_{m=1}^M$ are represented as vertices $\mathcal{V} = \{1, 2, 3, 4\}$ linked by the edges $\mathcal{E} = \{(4, 3), (3, 1), (3, 2), (1, 2)\}$. An undirected graphical model permits to factorize the distribution with conditional densities:

$$p_{\mathbf{x}} = \prod_{m=1}^M P(X^{(m)} | \pi(m)), \quad (3.13)$$

where $\pi(m)$ is the set of parents of the m -th vertex. For example, the factorization (3.13) applied to the example of Figure 3.3c is the following:

$$p_{\mathbf{x}} = P_{X^{(4)}}P_{(X^{(3)} | X^{(4)})}P_{(X^{(1)} | X^{(3)})}P_{(X^{(2)} | X^{(3)}, X^{(1)})}.$$

More details on graphical models and density factorizations are available in [57, 63]. In the following, we will focus on the directed case which can also be interpreted as low-rank tensor decompositions [56].

Figure 3.3: Example of graphical models for $M = 4$ variables.

While graphical models are very helpful to model high-dimensional distributions, they are not flexible enough especially for multimodal distributions for example. To obtain a more flexible model, we may assume that there are some latent variables that are not observed (see in Figure 3.4). The distribution of such a model is obtained by marginalizing the joint distribution with respect to the latent variables. Latent variable models are linked to tensor decompositions [56] and this link will be explained in the next subsection for the so-called naive Bayes model.

3.3.2 Naive Bayes model

Continuous case

A Naive Bayes Model (NBM) is a particular case of latent variable models. In the NBM, there is a single latent variable L which is the parent of every random variable (see Figure 3.4), making each variable $X^{(m)}$ independent conditionally to L . The variable L is a discrete random variable with R outcomes (or states) $\{1, \dots, R\}$. For the graph of Figure 3.4, the joint “density”⁴ of \mathbf{x} and L denoted as $p_{\mathbf{x},L}$ is factorized with (3.13) as:

$$p_{\mathbf{x},L}(x_1, \dots, x_M, r) = \Pr(L = r) P_{(X^{(1)} | L=r)}(x_1) \cdots P_{(X^{(M)} | L=r)}(x_M).$$

⁴The function $p_{\mathbf{x},L}$ defines the joint distribution of the continuous random vector \mathbf{x} and the discrete random variable L such that $\Pr(\mathbf{x} \in \mathcal{A}, L = r) = \int_{\mathcal{A}} p_{\mathbf{x},L}(x_1, \dots, x_M, r) dx_1 \cdots dx_M$.

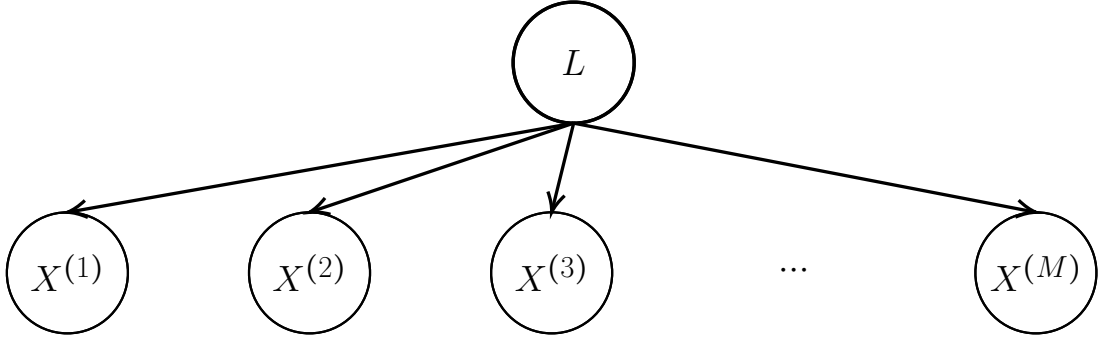


Figure 3.4: Naive Bayes model as a graphical model.

By marginalizing out the latent variable L over its R states, the naive Bayes model for the density of \mathbf{x} is obtained:

$$p_{\mathbf{x}}(x_1, \dots, x_M) = \sum_{r=1}^R \Pr(L = r) p_{(X^{(1)} | L=r)}(x_1) \cdots p_{(X^{(M)} | L=r)}(x_M). \quad (3.14)$$

Remark 3.3.2. The NBM (3.14) can be interpreted as a non-parametric finite mixture model [64] (see eq. (3.1)) where each f_r is an independent model hence a product of M densities

$$f_r(x_1, \dots, x_M) = p_{(X^{(1)} | L=r)}(x_1) \cdots p_{(X^{(M)} | L=r)}(x_M),$$

and the mixture weights are represented by the probabilities of the latent variable states ($\lambda_r = \Pr(L = r)$). This makes the NBM particularly interesting for flow cytometry data analysis, since a sample of cells can be seen as a mixture of cell populations (see Section 6.1.1).

Histogram tensor and discrete naive Bayes model

Let us consider a random vector \mathbf{x} that follows the continuous model (3.14) and let $\mathcal{I}^{(m)}$ be an interval partitioned into I bins as in Section 3.2.1. Then, because the bin probabilities (3.4) can be expressed as

$$h_{i_1 \dots i_M} = \int_{\Delta_{i_1}^{(1)}} \cdots \int_{\Delta_{i_M}^{(M)}} p_{\mathbf{x}}(x_1, \dots, x_M) dx_1 \dots dx_M,$$

and because the density is conditionally independent to L , it holds that:

$$\begin{aligned} h_{i_1 \dots i_M} &= \sum_{r=1}^R \Pr(L = r) \prod_{m=1}^M \int_{\Delta_{i_m}^{(m)}} \mathbb{P}_{(X^{(m)} | L=r)}(x_m) dx_m, \\ \implies h_{i_1 \dots i_M} &= \sum_{r=1}^R \Pr(L = r) \prod_{m=1}^M \Pr(X^{(m)} \in \Delta_{i_m}^{(m)} | L = r). \end{aligned} \quad (3.15)$$

Then, assuming that the support of $X^{(m)}$ is contained in the interval⁵ $\mathcal{I}^{(m)}$ (see Remark 3.2.3), let us define the random vector $\mathbf{z} = (Z^{(1)}, \dots, Z^{(M)})$ taking values in $\llbracket 1, I \rrbracket^M$ such that

$$\Pr(Z^{(1)} = i_1, \dots, Z^{(M)} = i_M) = \Pr(\mathbf{x} \in \Delta_{i_1}^{(1)} \times \dots \times \Delta_{i_M}^{(M)}).$$

Then it is easy to see that the distribution of \mathbf{z} follows a discrete version of the NBM:

$$\Pr(Z^{(1)} = i_1, \dots, Z^{(M)} = i_M) = \sum_{r=1}^R \Pr(L = r) \prod_{m=1}^M \Pr(Z^{(m)} = i_m | L = r). \quad (3.16)$$

Because we focus on estimation of histogram tensors, only the discrete NBM will be considered in the rest of the chapter.

3.3.3 Link between naive Bayes model and tensor decompositions

The NBM has a strong link to the CP decomposition (2.1). Considering a discrete NBM (3.16) where the M random variables $Z^{(m)}$ taking values in $\mathcal{I} = \llbracket 1, I \rrbracket$, then the probabilities can be arranged into a tensor \mathcal{H}

$$h_{i_1 \dots i_M} = \Pr(Z^{(1)} = i_1, \dots, Z^{(M)} = i_M).$$

Note that if \mathbf{z} is a quantized version of a continuous random vector (see previous subsection), then \mathcal{H} is exactly the histogram PMF tensor from Section 3.2.1. For simplicity, we refer to \mathcal{H} as the histogram tensor.

Let us define $\boldsymbol{\lambda}$ the vector of size R whose elements λ_r are equal to the latent state probability $\Pr(L = r)$. Moreover, we define $a_{i_m r}^{(m)}$ the conditional probability of the r -th state of L for the

⁵If it is not the case, we can consider, as in Remark 3.2.3, the conditional density supported on the product of $\mathcal{I}^{(m)}$ which also follows a continuous NBM.

m -th variable

$$a_{i_m r}^{(m)} = \Pr(Z^{(m)} = i_m \mid L = r),$$

which gives the following NBM expression:

$$h_{i_1, \dots, i_M} = \sum_{r=1}^R \lambda_r \prod_{m=1}^M a_{i_m r}^{(m)}. \quad (3.17)$$

Finally, by storing the values for each variable in *factor matrices* $\mathbf{A}^{(m)}$ of size $I \times R$

$$\mathbf{A}^{(m)} = \begin{bmatrix} a_{11}^{(m)} & \dots & a_{1R}^{(m)} \\ \vdots & & \vdots \\ a_{I1}^{(m)} & \dots & a_{IR}^{(m)} \end{bmatrix},$$

the NBM expression (3.17) becomes a CP decomposition

$$\mathcal{H} = \llbracket \boldsymbol{\lambda}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)} \rrbracket.$$

Note that the number of possible latent states R is also the decomposition rank of \mathcal{H} in section 3.3.3. In the following, we call a *factor* the vector $\mathbf{a}_r^{(m)}$ of size I which is defined by the r -th column⁶ of the factor matrix $\mathbf{A}^{(m)}$

$$\mathbf{a}_r^{(m)} = \begin{bmatrix} a_{1r}^{(m)} & \dots & a_{Ir}^{(m)} \end{bmatrix}^\top.$$

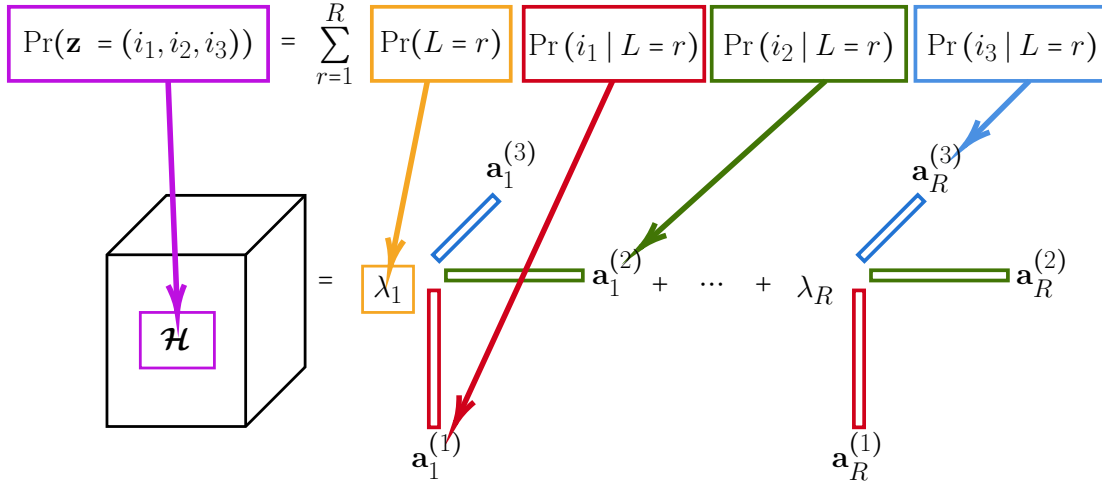
In Figure 3.5, the link between the CP decomposition and the NBM is visually represented.

Because the NBM is a probabilistic model, it imposes constraints on the induced CP model elements. Therefore, the CP factors and the loading vector must follow the two constraints below:

- Non-negativity: $\lambda_r \geq 0$ and $\mathbf{a}_r^{(m)} \geq 0$ for $m \in \llbracket 1, M \rrbracket$ and $r \in \llbracket 1, R \rrbracket$;
- Sum-to-one: $\sum_{r=1}^R \lambda_r = 1$ and $\sum_{i=1}^I a_{ir}^{(m)} = 1$ for $m \in \llbracket 1, M \rrbracket$ and $r \in \llbracket 1, R \rrbracket$.

These two conditions can be summarized with the *simplex* constraints and shortened by using

⁶Because $\mathbf{a}_r^{(m)}$ is defined as the r -th column of $\mathbf{A}^{(m)}$, the i_m -th element of $\mathbf{a}_r^{(m)}$ is the element of $\mathbf{A}^{(m)}$ of indices (i_m, r) , that is, equal to $a_{i_m r}^{(m)}$.


 Figure 3.5: Link between CPD and NBM for a discrete distribution with $M = 3$ variables.

factor matrices:

$$\begin{aligned} \boldsymbol{\lambda} &\geq 0, \quad \text{and} \quad \mathbb{1}_R^\top \boldsymbol{\lambda} = 1, \\ \mathbf{A}^{(m)} &\geq 0, \quad \text{and} \quad \mathbb{1}_I^\top \mathbf{A}^{(m)} = \mathbb{1}_R^\top. \end{aligned}$$

3.3.4 Curse of dimensionality: estimating the factors of the NBM

Regarding the complexity of the model, we saw that estimating the joint density of M variables is considered impossible in practice because it needs the estimation of a M -dimensional object. The NBM permits to reduce drastically the complexity of PDF estimation. By assuming a NBM for the joint density, the PDF is then described by R latent states and therefore RM one-dimensional densities. The complexity of the NBM is then linear with the number of dimensions M , which breaks the curse of dimensionality on this level.

If we assume that a joint discrete distribution follows a NBM, estimating this NBM comes to the estimation of the CP decomposition of \mathcal{H} . This problem can be formalized with the following optimization problem:

$$\begin{aligned} \widehat{\boldsymbol{\lambda}}, \widehat{\mathbf{A}}^{(1)}, \dots, \widehat{\mathbf{A}}^{(M)} &= \underset{\boldsymbol{\lambda}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)}}{\operatorname{argmin}} \left\| \widetilde{\mathcal{H}} - \llbracket \boldsymbol{\lambda}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)} \rrbracket \right\|_F^2, \\ \text{s.t.} \quad \boldsymbol{\lambda} &\geq 0, \mathbf{A}^{(m)} \geq 0, \\ \text{s.t.} \quad \mathbb{1}^\top \boldsymbol{\lambda} &= 1, \mathbb{1}^\top \mathbf{A}^{(m)} = \mathbb{1}^\top, \end{aligned} \tag{3.18}$$

where $\|\cdot\|_F$ denotes the Frobenius norm of a tensor and where $\tilde{\mathcal{H}}$ is the empirical histogram (see eq. (3.5)). We saw in Section 3.2.2 that a proper empirical estimation of \mathcal{H} requires a quantity of data not available in practice. We call this the first level of the curse of dimensionality. Even if both the NBM and the CPD are models whose complexities remain linear with M , the model itself cannot break fully the curse of dimensionality but only the first level. Now, because it is not possible to approximate well \mathcal{H} with the $\tilde{\mathcal{H}}$ given the amount of data available, the performance of an optimization algorithm, such as the AO-ADMM algorithm (Alternating Optimization solved with Alternating Direction Method of Multipliers [58]) is expected to be poor.

Workarounds exist like the method SQUAREM-PMF [23] that estimates the factors of the NBM directly from the data points. The complexity of this method relies on the number of samples and thus it permits to circumvent the curse of dimensionality on that matter. However, this comes at the expense of the convergence of the algorithm, as an increased number of missing values induces slower convergence rates [111]. Another workaround is the method of [74] presented in the following subsection that considers coupled tensor factorization of lower-order marginals. For order-3 marginals, the complexity now depends on $\binom{M}{3}I^3$ which is independent of the number of samples (the cost of computing histograms is often negligible). As it will be shown in Section 4.3.3, the complexities of both methods can be similar in particular situations. In the following subsection, it will be shown that the coupled tensor factorization method of [74] still has a prohibitive complexity regarding the number of variables.

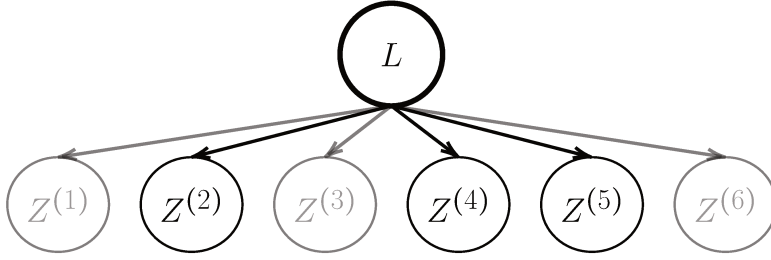
3.4 Fully coupled tensor factorization of 3D marginals

3.4.1 Marginalization of naive Bayes model

The idea of [74] to tackle the curse of dimensionality was to replace \mathcal{H} with a set of its projections (3D marginal tensors). For a subset of variables $(Z^{(j)}, Z^{(k)}, Z^{(l)})$, $\mathcal{H}^{(jkl)}$ denotes the 3D marginal tensor of size $I \times I \times I$ whose entries are

$$h_{i_j i_k i_\ell}^{(jkl)} = \Pr(Z^{(j)} = i_j, Z^{(k)} = i_k, Z^{(l)} = i_\ell).$$

As shown in [74], if \mathbf{z} follows a NBM then $\mathcal{H}^{(jkl)}$ has a also a CP decomposition which we will briefly explain below. Indeed, if \mathbf{z} follows a NBM (see Figure 3.4), then the marginal distribution of $(Z^{(j)}, Z^{(k)}, Z^{(l)})$ follows a NBM that corresponds to the subgraph of Figure 3.6.


 Figure 3.6: Example of marginalization of the NBM for $M = 6$.

More precisely the joint distribution of $(Z^{(j)}, Z^{(k)}, Z^{(l)})$ can be expressed as

$$\Pr(Z^{(j)} = i_j, Z^{(k)} = i_k, Z^{(l)} = i_\ell) = \quad (3.19)$$

$$\sum_{r=1}^R \Pr(L = r) \Pr(Z^{(j)} = i_j | L = r) \Pr(Z^{(k)} = i_k | L = r) \Pr(Z^{(l)} = i_\ell | L = r). \quad (3.20)$$

In terms of tensor notation, the eq. (3.20) can be written

$$\mathcal{H}^{(jkl)} = \llbracket \boldsymbol{\lambda}; \mathbf{A}^{(j)}, \mathbf{A}^{(k)}, \mathbf{A}^{(l)} \rrbracket. \quad (3.21)$$

Note that the expression (3.21) can be also obtained by using tensor contraction. Indeed, marginalizing out the m -th variable corresponds to contracting the full tensor \mathcal{H} on the m -th mode with the vector of ones $\mathbb{1}_I$. Therefore, it holds that the entries of $\mathcal{H}^{(jkl)}$ are

$$h_{i_j i_k i_\ell}^{(jkl)} = \underbrace{\sum_{i_1}^I \cdots \sum_{i_m}^I \cdots \sum_{i_M}^I}_{m \in \llbracket 1, M \rrbracket \setminus \{j, k, \ell\}} h_{i_1 \cdots i_M}.$$

As an example, Figure 3.6 features a NBM for $M = 6$ variables. For this example, the marginal distribution of variables $(Z^{(2)}, Z^{(4)}, Z^{(5)})$ is obtained by contracting the tensor \mathcal{H} on modes $(1, 3, 6)$ (which is equivalent to marginalizing out variables $(Z^{(1)}, Z^{(3)}, Z^{(6)})$)

$$\mathcal{H}^{(245)} = (((\mathcal{H} \bullet_6 \mathbb{1}_I) \bullet_3 \mathbb{1}_I) \bullet_1 \mathbb{1}_I).$$

3.4.2 Fully coupled tensor factorization

Thanks to the marginalization properties presented in Section 3.4.1, the set of 3D marginals $\{\mathcal{H}^{(jkl)}\}_{1 \leq j < k < \ell \leq M}$ admits a CP decomposition:

$$\mathcal{H}^{(jkl)} = \llbracket \boldsymbol{\lambda}; \mathbf{A}^{(j)}, \mathbf{A}^{(k)}, \mathbf{A}^{(\ell)} \rrbracket, \quad \text{for all } 1 \leq j < k < \ell \leq M, \quad (3.22)$$

where the factors of the CPD are shared and satisfy the non-negative and sum-to-one constraints.

In the estimation setup, [74] proposed to consider empirical histograms $\tilde{\mathcal{H}}^{(jkl)}$ as estimations of $\mathcal{H}^{(jkl)}$.

$$\tilde{h}_{i_j i_k i_\ell}^{(jkl)} = \text{Card}\{n \in \llbracket 1, N \rrbracket \mid (x_{nj}, x_{nk}, x_{n\ell}) = (i_j, i_k, i_\ell)\}, \quad (3.23)$$

While it is not possible to estimate \mathcal{H} , the 3D marginals $\{\mathcal{H}^{(jkl)}\}_{1 \leq j < k < \ell \leq M}$ can be obtained with a reasonable quantity of data points (see Section 3.2.2) by counting the number of occurrences inside each 3-dimensional bin. These empirical marginal histograms are then denoted $\tilde{\mathcal{H}}^{(jkl)}$ and are computed for all $\{j, k, \ell\} \subset \llbracket 1, M \rrbracket$ such that $1 \leq j < k < \ell \leq M$: where $\mathbf{x}_n \in \mathbb{R}^{1 \times M}$ is a row of the observation matrix \mathbf{X} containing N samples.

Therefore, the idea of [74] consisted in coupling estimated marginals $\{\tilde{\mathcal{H}}^{(jkl)}\}_{1 \leq j < k < \ell \leq M}$ to obtain the NBM of the full histogram. To do this, [74] proposed to solve the following optimization problem:

$$\begin{aligned} \hat{\boldsymbol{\lambda}}, \hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(M)} &= \underset{\boldsymbol{\lambda}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)}}{\text{argmin}} \sum_{j=1}^{M-2} \sum_{k=j+1}^{M-1} \sum_{\ell=k+1}^M \left\| \tilde{\mathcal{H}}^{(jkl)} - \llbracket \boldsymbol{\lambda}; \mathbf{A}^{(j)}, \mathbf{A}^{(k)}, \mathbf{A}^{(\ell)} \rrbracket \right\|_F^2, \\ \text{s.t. } \boldsymbol{\lambda} &\geq 0, \quad \mathbf{A}^{(1)} \geq 0, \quad \dots, \quad \mathbf{A}^{(M)} \geq 0, \\ \mathbb{1}_R^\top \boldsymbol{\lambda} &= 1, \quad \mathbb{1}_I^\top \mathbf{A}^{(1)} = \mathbb{1}_R^T, \quad \dots, \quad \mathbb{1}_I^\top \mathbf{A}^{(M)} = \mathbb{1}_R^T. \end{aligned} \quad (3.24)$$

This optimization problem (3.24) is solved via an alternating optimization method using ADMM to take into account the constraints. In the following, CTF3D (Coupled Tensor Factorization of 3D marginals) will denote the algorithm proposed by [74] that couples 3D marginals to obtain the factors of the CP decomposition.

Uniqueness of the fully coupled CP model

One of the key ideas behind the approach of [74] is to use the uniqueness properties of the CPD (see Section 2.2.2) of order-3 marginals to guarantee the recovery of the whole probability

tensor. Note that similarly it is possible to couple order-2 marginals to obtain the joint density [117] but at the expense of more restrictive uniqueness conditions.

The following generic uniqueness results for the coupled decomposition of marginals (3.22) have been proved in [74].

Theorem 3.4.1. [74, Theorem 1] – Let \mathcal{H} an order- M probability tensor that can be represented using a naive Bayes model of rank R . If $M \leq I$, then \mathcal{H} is identifiable from all marginals $\{\mathcal{H}^{(j k \ell)} \mid \{j, k, \ell\} \in \mathcal{T}_{all}\}$ if:

$$R \leq M(I - 2). \quad (3.25)$$

If $M > I$, then \mathcal{H} is identifiable from its marginals if:

$$R \leq \left(\left\lfloor \frac{\sqrt{MI - 1}}{I} \right\rfloor I - 1 \right)^2. \quad (3.26)$$

Theorem 3.4.2. [74, Theorem 2] – Let \mathcal{H} an order- M probability tensor that can be represented using a naive Bayes model of rank R . Let α be the maximal integer such that $2^\alpha \leq \lfloor \frac{M}{3} \rfloor I$. Then, \mathcal{H} is identifiable from all marginals $\{\mathcal{H}^{(j k \ell)} \mid \{j, k, \ell\} \in \mathcal{T}_{all}\}$ if:

$$R \leq 4^{\alpha-1}. \quad (3.27)$$

The idea of the proofs of Theorems 3.4.1 and 3.4.2 was to reduce the order- M coupled model into an order-3 model by grouping factor matrices into 3 groups. Thus, for example, the generic uniqueness results described in Section 2.2.3 can be used. However, the sum-to-one constraints seem to have been ignored in those theorem proofs. In Section 5.2, we show why and how the constraints should be properly taken care of, and that we can obtain stronger identifiability results than those of [74].

3.4.3 Curse of dimensionality: number of marginals

We saw that CTF3D permits to circumvent both the first two levels of the curse of dimensionality. The first one presented in Section 3.2.2 is solved by considering an NBM for the joint PDF of the variables. It narrowed down the problem of histogram estimation from I^M values to $R(IM + 1)$. Secondly, the principle of marginal factorization permits to circumvent the second curse of dimensionality (see Section 3.3.4). Indeed, obtaining the NBM factors using CTF3D requires just the estimation of $\binom{M}{3}$ empirical marginal 3D histograms (see Section 3.2.1), instead of the full M -dimensional empirical histogram.

Table 3.2: Number of marginals exhibiting the third level of curse of dimensionality, for $I = 15$ bins per dimension.

	I^M (see (3.5))	required N (see (3.12))	3D marginals $\binom{M}{3}$	Marginals values $\binom{M}{3}I^3$
$M = 3$	3375	4000	1	3375
$M = 4$	5.0×10^4	2.1×10^4	4	1.3×10^4
$M = 8$	2.6×10^9	1.6×10^7	56	1.9×10^5
$M = 10$	5.8×10^{11}	4.4×10^8	120	4.0×10^5
$M = 28$	8.5×10^{32}	4.0×10^{21}	3276	1.1×10^7

Therefore, the complexity now relies on the computation of 3D marginals which is no longer a problem with the quantity of data available in practice. However, the number of marginals to consider is $\binom{M}{3}$ and this number is increasing cubically with M (see Table 3.2). This means that there still remains a level of the curse of dimensionality which resides in the number of marginals to consider in the coupling. The first column denotes the number of values that describes the empirical estimator (see Equation (3.5)) for $I = 15$. The second column features the number of values required to estimate a M -dimensional histogram with $I = 15$, $\sigma = 0.1$ and $S = 10\sigma$ (see Equation (3.12)). Finally, the last column is plotting the third and last curse of dimensionality studied in this manuscript which is given by the number of triplets $\binom{M}{3}$. Despite the fact that the number of marginals values does not increase with M as much as I^M , it still represents a computational burden in terms of storage and handling of data. The next chapter will focus on a novel method that reduces again the complexity for histogram estimation.

3.5 Conclusion of Chapter 3

This chapter presented the density estimation problem. Histograms were described in detail, as well as their connection to tensors. The problem of the curse of dimensionality was explained, as well as the need for a model to reduce the complexity of the problem.

In the particular case of the naive Bayes model, a histogram is described with a number of values that remains linear with the number of dimensions. Moreover, the naive Bayes model is a particular case of graphical models, models that can be seen as low-rank tensor decompositions. Thus, it was shown that the problem of density estimation can be recast as a CP (Canonical Polyadic) approximation.

Estimating the factors of the decomposition remains a challenge, as it is still hampered by

the curse of dimensionality. To overcome this issue, the marginalization approach of [74] was presented, where the set of all 3D marginal histograms are estimated instead. However, this approach still has a number of drawbacks, such as the complexity growing cubically with the number of dimensions. In Chapter 4, a novel approach based on partially coupled tensor factorization will be presented, that aims at reducing the computational complexity of the coupled estimation problem. In Chapter 5, a detailed identifiability analysis is provided for the partially coupled and fully coupled cases, which improves and refines the results of [74]. Finally, it is shown in Chapter 6, how the coupled factorization can be effectively used for the high-dimensional FCM data analysis and visualisation.

Chapter 4

Partially coupled tensor factorization for probability mass function estimation

Contents

4.1	Partially coupled tensor factorization	72
4.1.1	Principles	72
4.1.2	Optimization	73
4.1.3	An example in the case of 4 variables	74
4.2	Coupling strategies	77
4.2.1	Coupling strategies as hypergraphs	78
4.2.2	Examples of coupling strategies	80
4.2.3	An algorithm for generating balanced couplings	87
4.3	Numerical experiments	90
4.3.1	Random strategies and performance regarding the number of triplets	91
4.3.2	Comparison between random and balanced couplings.	91
4.3.3	Comparison with a KLD-based method	94
4.4	Conclusion of Chapter 4	97

In this chapter, Section 4.1 introduces a novel approach called partially coupled tensor factorization which is included in [43, 44]. This method resides on considering only a subset of all possible marginals to obtain the CP decomposition of a dataset.

In Section 4.2, different choices for the set of triplets to consider in the coupling are presented. These so-called coupling strategies are introduced and formalized as hypergraphs and more particularly 3-graphs. Finally, balanced couplings are introduced as couplings for which variables are represented evenly in terms of occurrences. An algorithm to generate those strategies which relies on Lyndon words is presented.

Before concluding this chapter, Section 4.3 proposes numerical studies on the different characteristics of our method compared to the literature. Two studies were conducted on the estimation performance regarding the choice of triplets considered in the coupling. A third experiment compares thoroughly our novel approach with an algorithm using a Kullback-Leibler divergence metric. This chapter does not present numerical experiments on flow cytometry data which are reserved for Chapter 6. However, synthetic datasets and real datasets are presented in this section.

4.1 Partially coupled tensor factorization

4.1.1 Principles

To reduce the complexity of the method presented in Section 3.4, we proposed in [43] a new method for PMF estimation. This method is called Partially Coupled Tensor Factorization of 3D marginals or PCTF3D. The principle of PCTF3D is simple: couple only a subset of marginals instead of $\binom{M}{3}$. This permits to reduce the complexity of the method as there are less marginals to compute, store and couple to obtain the higher-order CPD. PCTF3D solves a similar problem presented in (3.24) where the sum over all possible triplets of variables is reduced to a subset of triplets:

$$\begin{aligned}
 \widehat{\lambda}, \widehat{\mathbf{A}}^{(1)}, \dots, \widehat{\mathbf{A}}^{(M)} &= \underset{\lambda, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)}}{\operatorname{argmin}} \sum_{\{j,k,\ell\} \in \mathcal{T}} \left\| \widetilde{\mathcal{H}}^{(jk\ell)} - \llbracket \lambda; \mathbf{A}^{(j)}, \mathbf{A}^{(k)}, \mathbf{A}^{(\ell)} \rrbracket \right\|_F^2, \\
 \text{s.t. } \lambda &\geq 0, \quad \mathbf{A}^{(1)} \geq 0, \quad \dots, \quad \mathbf{A}^{(M)} \geq 0, \\
 \mathbb{1}_R^\top \lambda &= 1, \quad \mathbb{1}_I^\top \mathbf{A}^{(1)} = \mathbb{1}_R^\top, \quad \dots, \quad \mathbb{1}_I^\top \mathbf{A}^{(M)} = \mathbb{1}_R^\top.
 \end{aligned} \tag{4.1}$$

In Equation (4.1), \mathcal{T} denotes the set of triplets that are considered in the coupling and thus will be referred to as the coupling in the following. Note that CTF3D is a particular case of PCTF3D where the coupling \mathcal{T} contains all possible triplets of variables in $\llbracket 1, M \rrbracket$:

$$\mathcal{T}_{\text{all}} := \{ \{j, k, \ell\} \subset \llbracket 1, M \rrbracket \mid 1 \leq j < k < \ell \leq M \}.$$

The conditions that must be satisfied by the couplings are discussed in the following subsection on coupling strategies.

4.1.2 Optimization

In PCTF3D, like in CTF3D the so-called AO-ADMM framework [58] is applied. In this framework, the variables are updated in an alternative fashion (hence Alternate Optimization). Each variable is updated by solving the problem (4.1) where all variables are fixed (and equal to their last update) except the one that is updated. This procedure defines $M + 1$ optimization sub-problems: one for the update of each factor matrix $\mathbf{A}^{(m)}$ and one for the update of the loading vector $\boldsymbol{\lambda}$. For $\mathbf{A}^{(m_0)}$, the $(t + 1)$ -th update denoted $\mathbf{A}_{t+1}^{(m_0)}$ is obtained via the current state of each factor and loading vector:

$$\mathbf{A}_{t+1}^{(1)}, \dots, \mathbf{A}_{t+1}^{(m_0-1)}, \mathbf{A}_t^{(m_0+1)}, \dots, \mathbf{A}_t^{(M)}, \boldsymbol{\lambda}_t.$$

For convenience, let us denote σ the function of m and t that returns the state of update of the m -th factor matrix when the factor m_0 is updated:

$$\sigma(m, t) := \begin{cases} t + 1 & \text{if } m \leq m_0 \\ t & \text{otherwise} \end{cases}$$

When updating the factor m_0 , the marginals $\tilde{\mathcal{H}}^{(jk\ell)}$ such that $m_0 \notin \{j, k, \ell\}$ are not used. Therefore, let \mathcal{T}_{m_0} the subset of triplets of variables of \mathcal{T} that contains the m_0 -th variable:

$$\mathcal{T}_{m_0} := \{ \{j, k, \ell\} \in \mathcal{T} \mid m_0 \in \{j, k, \ell\} \}.$$

The update $\mathbf{A}_{t+1}^{(m_0)}$ is obtained by solving the following optimization problem:

$$\begin{aligned} \mathbf{A}_{t+1}^{(m_0)} = \operatorname{argmin}_{\mathbf{A}^{(m_0)}} \sum_{\substack{k,\ell \\ \{m_0,k,\ell\} \in \mathcal{T}_{m_0}}} \left\| \tilde{\mathcal{H}}_{(1)}^{(m_0 k \ell)} - \mathbf{A}^{(m_0)} \operatorname{diag}(\boldsymbol{\lambda}_t) \left(\mathbf{A}_{\sigma(\ell,t)}^{(\ell)} \odot \mathbf{A}_{\sigma(k,t)}^{(k)} \right)^\top \right\|_F^2 \quad (4.2) \\ \text{s.t. } \mathbf{A}^{(m_0)} \geq 0, \quad \mathbb{1}_I^\top \mathbf{A}^{(m_0)} = \mathbb{1}_R^\top. \end{aligned}$$

The cost function (4.2) represents the sum of the errors over all marginals that contains the variable m_0 . Concerning $\boldsymbol{\lambda}$, its update is made after the update of all factor matrices and is obtained by solving the following optimization problem:

$$\begin{aligned} \boldsymbol{\lambda}_{t+1} = \operatorname{argmin}_{\boldsymbol{\lambda}} \sum_{\{j,k,\ell\} \in \mathcal{T}} \left\| \operatorname{vec}(\tilde{\mathcal{H}}^{(j k \ell)}) - (\mathbf{A}_{t+1}^{(\ell)} \odot \mathbf{A}_{t+1}^{(k)} \odot \mathbf{A}_{t+1}^{(j)}) \boldsymbol{\lambda} \right\|_F^2. \quad (4.3) \\ \text{s.t. } \boldsymbol{\lambda} \geq 0, \quad \mathbb{1}_R^\top \boldsymbol{\lambda} = 1. \end{aligned}$$

An overview of the procedure is presented in Algorithm 1. There are several possible convergence criteria that can be chosen to stop PCTF3D. First, it is possible to define a maximal number of outer iterations T_1 . It is also possible to stop PCTF3D when the difference between two updates is below a certain tolerance. Another criterion may be that the sum of the errors over all marginals is below a tolerance value. This last criterion may not be applied in certain cases, as the number of 3D marginals may make this estimation of this error computationally prohibitive.

The sub-problems (4.2) and (4.3) are solved via an Alternating Direction Method of Multipliers (ADMM) [14]. The two iterative procedures are presented in Algorithms 2 and 3. Because the complexity of ADMM is in $\mathcal{O}(IR^2)$ [58], PCTF3D has a complexity $\mathcal{O}(IR^2T)$ which depends on the number of triplets T . When all triplets are considered, the algorithm CTF3D obtains a complexity in $\mathcal{O}(IR^2 \binom{M}{3}) = \mathcal{O}(IR^2M^3)$.

4.1.3 An example in the case of 4 variables

For $M = 4$ variables, the goal is to estimate the histogram \mathcal{H} which can also be written $\mathcal{H}^{(1234)}$. Let us consider the set of triplets $\mathcal{T} = \{\{1, 3, 4\}, \{2, 3, 4\}\}$. Hence the two empirical histograms $\tilde{\mathcal{H}}^{(134)}$ and $\tilde{\mathcal{H}}^{(234)}$ are estimated with Equation (3.23).

To obtain an estimation of the four factor matrices $\{\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)}, \mathbf{A}^{(4)}\}$ and the loading

Algorithm 1: Algorithm PCTF3D that solves (4.1).

Input: $R, I, \mathbf{X}, \mathcal{T}$.

Output: $\{\mathbf{A}^{(m)}\}_{m=1}^M$ and $\boldsymbol{\lambda}$.

Initialization : $\{\mathbf{A}^{(m)} \in \mathbb{R}^{I \times R}\}_{m=1}^M$ and $\boldsymbol{\lambda} \in \mathbb{R}^R$ random such that non-negativity and sum-to-one constraints are satisfied.

for all (j, k, ℓ) in \mathcal{T} **do**

 | Estimate with \mathbf{X} the 3D histogram $\tilde{\mathcal{H}}^{(j k \ell)} \in \mathbb{R}^{I \times I \times I}$ with (3.23)

end

while convergence criterion is not satisfied **do**

for $m = 1$ to M **do**

 | Update of $\mathbf{A}^{(m)}$ by solving optimization problem (4.2),

end

 Update of $\boldsymbol{\lambda}$ by solving optimization problem (4.3)

end

Algorithm 2: Algorithm for updating factor matrices (solves (4.2)).

Input: $\{\mathbf{A}_{\sigma(m,t)}^{(m)}\}_{m=1}^M, \boldsymbol{\lambda}_t, \{\tilde{\mathcal{H}}^{(m_0 k \ell)}\}_{\{m_0, k, \ell\} \in \mathcal{T}_{m_0}}, T_2$.

Output: $\mathbf{A}_{t+1}^{(m_0)} = \mathbf{A}_{T_2}$.

Initialization : $\mathbf{A}_{t_2=0} = \mathbf{A}_t^{(m_0)}, \mathbf{U}_{t_2=0} \in \mathbb{R}^{I \times R}$ and $\mathbf{B}_{t_2=0} \in \mathbb{R}^{R \times I}$ with zeros.

for all triplets of variables $\{m_0, k, \ell\} \in \mathcal{T}_{m_0}$ **do**

 | $\mathbf{Q}^{(\ell k)} = \mathbf{A}_{\sigma(\ell,t)}^{(\ell)} \odot \mathbf{A}_{\sigma(k,t)}^{(k)}$,

end

$\mathbf{G} = (\boldsymbol{\lambda}_t \boldsymbol{\lambda}_t^\top) * \sum_{\substack{k, \ell \\ \{m_0, k, \ell\} \in \mathcal{T}_{m_0}}} \mathbf{Q}^{(\ell k) \top} \mathbf{Q}^{(\ell k)}$,

$\mathbf{W} = \text{diag}(\boldsymbol{\lambda}_t) \sum_{\substack{k, \ell \\ \{m_0, k, \ell\} \in \mathcal{T}_{m_0}}} \mathbf{Q}^{(\ell k) \top} \tilde{\mathcal{H}}_{(1)}^{(m_0 k \ell) \top}$,

$\rho = \frac{1}{R} \text{tr}(\mathbf{G})$,

for at most $t_2 \in \llbracket 1, T_2 \rrbracket$ iterations **do**

 | $\mathbf{B}_{t_2} = (\mathbf{G} + \rho \mathbf{I}_R)^{-1} (\mathbf{W} + \rho (\mathbf{A}_{t_2-1} + \mathbf{U}_{t_2-1})^\top)$,

 | $\mathbf{A}_{t_2} = \mathbf{A}_{t_2-1} - \mathbf{B}_{t_2}^\top + \mathbf{U}_{t_2-1}$,

 | Projection of \mathbf{A}_{t_2} onto the simplex constraints [38],

 | $\mathbf{U}_{t_2} = \mathbf{U}_{t_2-1} + \mathbf{A}_{t_2} - \mathbf{B}_{t_2}^\top$.

end

Algorithm 3: Algorithm for loading vector update (solves (4.3)).

Input: $\{\mathbf{A}_{t+1}^{(m)}\}_{m=1}^M$, $\boldsymbol{\lambda}_t$, $\{\tilde{\mathcal{H}}^{(jkl)}\}_{\{j,k,\ell\} \in \mathcal{T}}$, T_2 .

Output: $\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_{T_2}$.

Initialization : $\boldsymbol{\lambda}_{t_2=0} = \boldsymbol{\lambda}_t$, $\mathbf{u}_{t_2=0} \in \mathbb{R}^R$ and $\boldsymbol{\mu}_{t_2=0} \in \mathbb{R}^R$ with zeros.

for all triplets of variables $\{j, k, \ell\}$ in \mathcal{T} **do**

$\mathbf{Q}^{(\ell kj)} = \mathbf{A}_{t+1}^{(\ell)} \odot \mathbf{A}_{t+1}^{(k)} \odot \mathbf{A}_{t+1}^{(j)}$,

end

$\mathbf{G} = \sum_{\{j,k,\ell\} \in \mathcal{T}} \mathbf{Q}^{(\ell kj)\top} \mathbf{Q}^{(\ell kj)}$,

$\mathbf{w} = \sum_{\{j,k,\ell\} \in \mathcal{T}} \mathbf{Q}^{(\ell kj)\top} \text{vec}(\tilde{\mathcal{H}}^{(jkl)})$,

$\rho = \frac{1}{R} \text{tr}(\mathbf{G})$,

for at most $t_2 \in \llbracket 1, T_2 \rrbracket$ iterations **do**

$\boldsymbol{\mu}_{t_2} = (\mathbf{G} + \rho \mathbf{I}_R)^{-1} (\mathbf{w} + \rho(\boldsymbol{\lambda}_{t_2-1} + \mathbf{u}_{t_2-1}))$,

$\boldsymbol{\lambda}_{t_2} = \boldsymbol{\lambda}_{t_2-1} - \boldsymbol{\mu}_{t_2} + \mathbf{u}_{t_2-1}$,

 Projection of $\boldsymbol{\lambda}_{t_2}$ onto the simplex constraints [38],

$\mathbf{u}_{t_2} = \mathbf{u}_{t_2-1} + \boldsymbol{\lambda}_{t_2} - \boldsymbol{\mu}_{t_2}$.

end

vector $\boldsymbol{\lambda}$, PCTF3D solves the following coupled optimization problem:

$$\begin{aligned} \widehat{\boldsymbol{\lambda}}, \widehat{\mathbf{A}}^{(1)}, \widehat{\mathbf{A}}^{(2)}, \widehat{\mathbf{A}}^{(3)}, \widehat{\mathbf{A}}^{(4)} &= \underset{\boldsymbol{\lambda}, \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)}, \mathbf{A}^{(4)}}{\text{argmin}} \\ &\left\| \tilde{\mathcal{H}}^{(134)} - \llbracket \boldsymbol{\lambda}, \mathbf{A}^{(1)}, \mathbf{A}^{(3)}, \mathbf{A}^{(4)} \rrbracket \right\|_F^2 + \left\| \tilde{\mathcal{H}}^{(234)} - \llbracket \boldsymbol{\lambda}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)}, \mathbf{A}^{(4)} \rrbracket \right\|_F^2, \\ \text{s.t. } \boldsymbol{\lambda} &\geq 0, \quad \mathbf{A}^{(1)} \geq 0, \mathbf{A}^{(2)} \geq 0, \mathbf{A}^{(3)} \geq 0, \mathbf{A}^{(4)} \geq 0, \\ \mathbb{1}_R^\top \boldsymbol{\lambda} &= 1, \quad \mathbb{1}_I^\top \mathbf{A}^{(1)} = \mathbb{1}_R^\top, \mathbb{1}_I^\top \mathbf{A}^{(2)} = \mathbb{1}_R^\top, \mathbb{1}_I^\top \mathbf{A}^{(3)} = \mathbb{1}_R^\top, \mathbb{1}_I^\top \mathbf{A}^{(4)} = \mathbb{1}_R^\top. \end{aligned}$$

Because PCTF3D proceeds by alternating optimization, let us determine the $(t+1)$ -th update of the factor $\mathbf{A}^{(1)}$ denoted $\mathbf{A}_{t+1}^{(1)}$. This can be done by solving the following optimization algorithm:

$$\begin{aligned} \mathbf{A}_{t+1}^{(1)} &= \underset{\mathbf{A}^{(1)}}{\text{argmin}} \left\| \tilde{\mathcal{H}}_{(1)}^{(134)} - \mathbf{A}^{(1)} \text{diag}(\boldsymbol{\lambda}_t) \left(\mathbf{A}_t^{(4)} \odot \mathbf{A}_t^{(3)} \right)^\top \right\|_F^2 \\ \text{s.t. } \mathbf{A}^{(1)} &\geq 0, \quad \mathbb{1}_I^\top \mathbf{A}^{(1)} = \mathbb{1}_R^\top. \end{aligned}$$

This optimization problem does not feature the histogram $\tilde{\mathcal{H}}^{(234)}$. Indeed, this term is eliminated because it does not depend on the factor matrix $\mathbf{A}^{(1)}$. For the factor $\mathbf{A}^{(2)}$, the update is obtained similarly. Concerning the third factor, the optimization problem depends on both terms of the

sum:

$$\begin{aligned} \mathbf{A}_{t+1}^{(3)} = \operatorname{argmin}_{\mathbf{A}^{(3)}} & \left\| \tilde{\mathcal{H}}_{(1)}^{(314)} - \mathbf{A}^{(3)} \operatorname{diag}(\boldsymbol{\lambda}_t) \left(\mathbf{A}_t^{(4)} \odot \mathbf{A}_{t+1}^{(1)} \right)^\top \right\|_F^2 \\ & + \left\| \tilde{\mathcal{H}}_{(1)}^{(324)} - \mathbf{A}^{(3)} \operatorname{diag}(\boldsymbol{\lambda}_t) \left(\mathbf{A}_t^{(4)} \odot \mathbf{A}_{t+1}^{(2)} \right)^\top \right\|_F^2 \\ \text{s.t. } & \mathbf{A}^{(3)} \geq 0, \quad \mathbb{1}_I^\top \mathbf{A}^{(3)} = \mathbb{1}_R^\top. \end{aligned}$$

4.2 Coupling strategies

Partially coupled tensor factorization introduced the coupling \mathcal{T} . In this section, a coupling will be defined by using the formalism of hypergraphs. Some necessary conditions on \mathcal{T} will be proposed to ensure that PCTF3D allows to estimate the CPD factors. Then, some possible coupling strategies will be presented as well as a new *balanced* approach.

As a reminder, a coupling \mathcal{T} is a subset of triplets of variables. It represents the indices of the sum of the optimization problem (4.1). It is possible to represent couplings as shown in Figure 4.1. First, \mathcal{T} must feature all variables in $\llbracket 1, M \rrbracket$. Otherwise, if one missing variable is denoted \bar{m} , the optimization problem (4.1) is not a function of $\mathbf{A}^{(\bar{m})}$ (see Figure 4.1a). A coupling \mathcal{T} will be called *valid* if it satisfies the following assumptions.

Assumption 4.2.1. *All variables must be represented in \mathcal{T} :*

$$\bigcup_{\{j,k,\ell\} \in \mathcal{T}} \{j, k, \ell\} = \llbracket 1, M \rrbracket.$$

Assumption 4.2.2. *For all $(m_1, m_2) \in \llbracket 1, M \rrbracket^2$ there exists a finite set of triplets (τ_1, \dots, τ_T) such that $m_1 \in \tau_1$, $m_2 \in \tau_T$ and:*

$$\forall t \in \llbracket 1, T-1 \rrbracket, \tau_t \cap \tau_{t+1} \neq \emptyset.$$

Note that the Assumption 4.2.2 implies that all variables are present thus the Assumption 4.2.1.

For example, in Figure 4.1b, even if all variables are present in the coupling, the first 3 variables are decoupled from the last three. Uncoupling between sets of variables results in permutation ambiguities between factors. Figure 4.1c presents a valid coupling featuring $T = 3$ triplets for $M = 6$ variables.

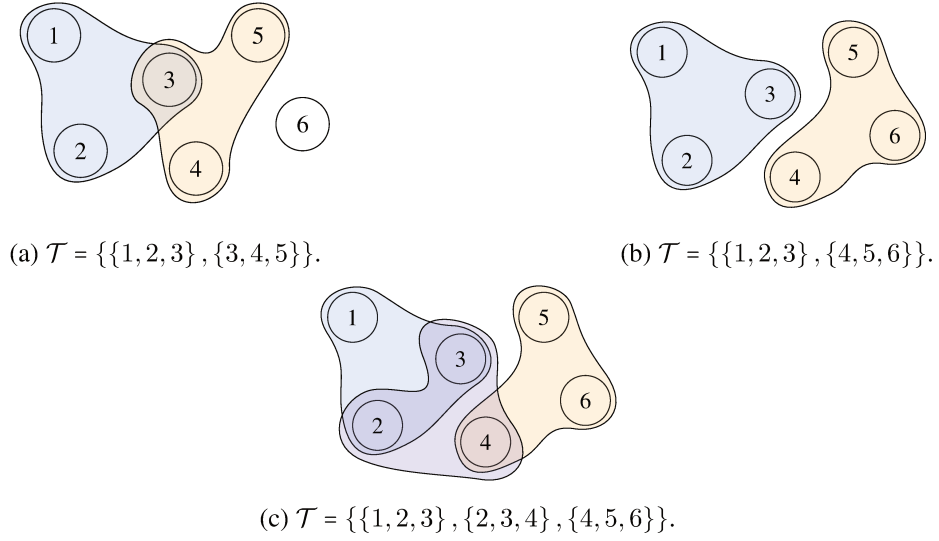


Figure 4.1: Representation of couplings for 6 variables.

4.2.1 Coupling strategies as hypergraphs

Coupling can be formalized with the theory of hypergraphs. A hypergraph is a graph whose edges are not limited to two nodes.

Definition 4.2.3. *r -uniform hypergraph [46]* – A hypergraph is a tuple $(\mathcal{V}, \mathcal{E})$ where \mathcal{V} is a finite set of vertices and \mathcal{E} is a set of hyper-edges (a set of non-empty subsets of \mathcal{V}). A hypergraph is called r -uniform if all elements in \mathcal{E} have cardinality r .

In our case of coupled factorization of 3D marginals, nodes are the M variables hence $\mathcal{V} = \llbracket 1, M \rrbracket$. Each triplet of a coupling \mathcal{T} can be seen as an edge of hypergraph. In particular, because all marginals considered are 3-dimensional, the associated hypergraph is then a 3-uniform hypergraph or shortly a 3-graph. In this thesis, only 3-uniform hypergraphs are considered and the terms hypergraph \mathcal{T} or 3-graph \mathcal{T} will be used to denote the coupling \mathcal{T} .

Definition 4.2.4. Incidence matrix – Let $\mathcal{T} = \{\{j_1, k_1, \ell_1\}, \dots, \{j_T, k_T, \ell_T\}\}$ be a hypergraph containing T triplets. The incidence matrix of \mathcal{T} is a $T \times M$ binary matrix denoted \mathbf{V} such that for all $(t, m) \in \llbracket 1, T \rrbracket \times \llbracket 1, M \rrbracket$:

$$v_{tm} := \begin{cases} 1 & \text{if } m \in \{j_t, k_t, \ell_t\} \\ 0 & \text{else} \end{cases} .$$

Note that for a 3-graph \mathcal{T} , it holds that $\mathbf{V}\mathbb{1}_M = 3\mathbb{1}_T$.

Definition 4.2.5. Connected hypergraph – For a hypergraph $\mathcal{T} = (\mathcal{V}, \mathcal{E})$, a path denotes a sequence of N distinct hyper-edges $(\varepsilon_1, \dots, \varepsilon_N) \in \mathcal{E}^N$ such that:

$$\forall n \in \llbracket 1, N-1 \rrbracket, \quad \text{Card}(\varepsilon_n \cap \varepsilon_{n+1}) > 0.$$

Two nodes (m_1, m_2) of a hypergraph \mathcal{T} are said to be connected if \mathcal{T} contains a path linking the two nodes – that is if $m_1 \in \varepsilon_1$ and $m_2 \in \varepsilon_N$. Therefore, a hypergraph is said connected if all pairs of nodes are connected.

Remark 4.2.6. Connectedness of a hypergraph \mathcal{T} is equivalent to the Assumption 4.2.2.

Definition 4.2.7. Sequence of degrees of a hypergraph – Let \mathcal{T} be a hypergraph of M variables. The sequence of degree of \mathcal{T} is defined as the row vector \mathbf{d} of size M that contains the number of occurrences of each node in \mathcal{T} :

$$d_m := \text{Card}\{\{j, k, \ell\} \in \mathcal{T} \mid m \in \{j, k, \ell\}\}.$$

In terms of incidence matrix, the degree sequence of a hypergraph is given by $\mathbb{1}_T^T \mathbf{V} = \mathbf{d}$.

Proposition 4.2.8. For a valid coupling \mathcal{T} of T triplets, the degree sequence is bounded by:

$$\forall m \in \llbracket 1, M \rrbracket, \quad d_m \in \left[\max\left(1, T - \binom{M-1}{3}\right), \min\left(T, \binom{M-1}{2}\right) \right].$$

Proof. Without the valid coupling constraint, d_m follows a hyper-geometric distribution. Indeed, picking triplets is a drawing of T triplets among $\binom{M}{3}$ without replacement. Moreover, if picking a triplet containing the variable m is considered a success, then there are $\binom{M-1}{2}$ samples that leads to this event. The support of a hyper-geometric distribution gives the following bounds for d_m :

$$d_m \in \left[\max\left(0, T + \binom{M-1}{2} - \binom{M}{3}\right), \min\left(T, \binom{M-1}{2}\right) \right].$$

To finish the proof, first it holds $T + \binom{M-1}{2} - \binom{M}{3} = T - \binom{M-1}{3}$. Finally, because only valid couplings are studied, the support of d_m is restricted to $d_m \geq 1$ (see the Assumption 4.2.1). \square

To give some examples, couplings of Figure 4.1 are featured in Table 4.1 with their associated incidence matrices and degree sequences. Note that the Assumption 4.2.1 is equivalent to ensuring that there are no zero-columns in \mathbf{V} . Concerning the Assumption 4.2.2, it can also be

obtained by ensuring that there does not exist a permutation of both rows and columns such that \mathbf{V} can be written as a diagonal block matrix.

Proposition 4.2.9 (Special case of [13, Theorem 3.9]). **Sufficient coupling condition** – Let \mathcal{T} a coupling 3-graph containing T different triplets. Then, if \mathcal{T} is connected, it means that:

$$\left\lfloor \frac{M}{2} \right\rfloor \leq T \leq \binom{M}{3}.$$

Definition 4.2.10. Step-1 degree sequence – For a coupling \mathcal{T} , a degree sequence is said step-1 if its maximal discrepancy is lower than 1:

$$\max \mathbf{d} - \min \mathbf{d} \leq 1.$$

Definition 4.2.11. Balanced coupling – A coupling \mathcal{T} is said balanced if its degree sequence is step-1. Couplings with strictly constant (or step-0) sequence of degrees will be called perfectly balanced couplings.

Definition 4.2.12. Number of pairs P – For a coupling \mathcal{T} , P denotes the number of different pairs contained inside the hypergraph \mathcal{T} :

$$P := \text{Card}\{\{m_1, m_2\} \subset \llbracket 1, M \rrbracket \mid m_1 \neq m_2 \text{ and } \exists \tau \in \mathcal{T}, \{m_1, m_2\} \subset \tau\}.$$

Table 4.1 gives examples of \mathbf{d} and P for different couplings. These two quantities are important especially for the study of identifiability and recoverability (see Section 5.1.5).

4.2.2 Examples of coupling strategies

A partial coupling comes to the choice of T different triplets from the set of $\binom{M}{3}$ triplets. This combinatorial problem cannot be exhaustively resolved as the number of possible couplings is $\binom{\binom{M}{3}}{T}$. In the following, a *coupling strategy* denotes a method to choose a coupling at a fixed T . Note that there is only one coupling which contains $\binom{M}{3}$ triplets and this fully coupled strategy results in CTF3D. It can also be noted that couplings presented in the following subsections are defined up to a permutation of variables. The different strategies presented in the following are all summarized in Table 4.2.

Table 4.1: Examples of couplings in the case of 6 variables (see Figure 4.1).

\mathcal{T}	\mathbf{V}	\mathbf{d}	P
$\{\{1, 2, 3\}, \{3, 4, 5\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 \end{bmatrix}$	$[1 \ 1 \ 2 \ 1 \ 1 \ 0]$	6
$\{\{1, 2, 3\}, \{4, 5, 6\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$	$[1 \ 1 \ 1 \ 1 \ 1 \ 1]$	6
$\{\{1, 2, 3\}, \{2, 3, 4\}, \{4, 5, 6\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$	$[1 \ 2 \ 2 \ 2 \ 1 \ 1]$	8
$\{\{j, k, \ell\} \subset \llbracket 1, 6 \rrbracket \mid j < k < \ell\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ & & : 20 \text{ rows} \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$	$[10 \ 10 \ 10 \ 10 \ 10 \ 10]$	15

Table 4.2: Coupled strategies properties and associated complexities.

Strategy	Number of triplets	Complexity of PCTF3D	Properties
Full	$\binom{M}{3}$	$\mathcal{O}(IR^2M^3)$	Corresponds to CTF3D [74]
'+2'	$\lfloor \frac{M}{2} \rfloor$	$\mathcal{O}(IR^2M)$	Least coupled strategy
'+1'	M	$\mathcal{O}(IR^2M)$	Balanced degree sequence
Random	T	$\mathcal{O}(IR^2T)$	Unbalanced degree sequence
Balanced	T	$\mathcal{O}(IR^2T)$	Balanced degree sequence

Least coupling strategy ('+2' strategy)

Proposition 4.2.9 shows that a coupling should contain at least $\lfloor \frac{M}{2} \rfloor$ triplets. Although there are several valid couplings of size $\lfloor \frac{M}{2} \rfloor$, we propose a coupling strategy that consists in adding '+2' between two triplets hence the name '+2'. For $M \geq 4$, the '+2' coupling is denoted $\mathcal{T}^{(+2)}$:

- If M is odd:

$$\mathcal{T}^{(+2)} := \{\{1, 2, 3\}, \{3, 4, 5\}, \dots, \{M-4, M-3, M-2\}, \{M-2, M-1, M\}\},$$

- If M is even:

$$\mathcal{T}^{(+2)} := \{\{1, 2, 3\}, \{3, 4, 5\}, \dots, \{M-3, M-2, M-1\}, \{M-2, M-1, M\}\}.$$

Examples of $\mathcal{T}^{(+2)}$ couplings are given in Table 4.3 for small values of M . The '+2' strategy gives couplings with the least possible amount of triplets. Indeed, if one removes a triplet, at least one variable will be missing (see Assumption 4.2.1). Moreover, if one decides to spread even more the triplets, the resulting couplings will no longer satisfy the Assumption 4.2.2.

The couplings obtained with the '+2' strategy have at most degree 2 and are by definition balanced. Note that there exist couplings with the least amount of triplets $T = \lfloor \frac{M}{2} \rfloor$ which have unbalanced degree sequences. For example, it is possible to choose each triplets with the first variable and the other two variables are not overlapping across triplets:

- If M is odd:

$$\mathcal{T} = \{\{1, 2, 3\}, \{1, 4, 5\}, \{1, 6, 7\}, \dots, \{1, M-1, M\}\},$$

$$\mathbf{d} = \begin{bmatrix} T & 1 & \dots & 1 \end{bmatrix},$$

- If M is even:

$$\mathcal{T} = \{\{1, 2, 3\}, \{1, 4, 5\}, \{1, 6, 7\}, \dots, \{1, M-1, M-2\}, \{1, M-1, M\}\},$$

$$\mathbf{d} = \begin{bmatrix} T & 1 & \dots & 1 & 2 & 1 \end{bmatrix}.$$

'+1' strategy

The '+1' strategy consists in choosing $T = M$ triplets such that consecutive triplets have a 2-variables overlap. Those couplings are denoted as $\mathcal{T}^{(+1)}$ and can be obtained by adding '+1' between two consecutive triplets:

$$\mathcal{T}^{(+1)} := \{\{1, 2, 3\}, \{2, 3, 4\}, \dots, \{M-2, M-1, M\}, \{M-1, M, 1\}, \{M, 1, 2\}\}.$$

Table 4.3: Examples of '+2' couplings associated with their incidence matrix and degree sequence.

	$\mathcal{T}^{(+2)}$	\mathbf{V}	\mathbf{d}	P
$M = 4$	$\{\{1, 2, 3\}, \{2, 3, 4\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix}$	$[1 \ 2 \ 2 \ 1]$	5
$M = 5$	$\{\{1, 2, 3\}, \{3, 4, 5\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix}$	$[1 \ 1 \ 2 \ 1 \ 1]$	5
$M = 6$	$\{\{1, 2, 3\}, \{3, 4, 5\}, \{4, 5, 6\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$	$[1 \ 1 \ 2 \ 2 \ 2 \ 1]$	8
$M = 7$	$\{\{1, 2, 3\}, \{3, 4, 5\}, \{5, 6, 7\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$	$[1 \ 1 \ 2 \ 1 \ 2 \ 1 \ 1]$	9
$M = 8$	$\{\{1, 2, 3\}, \{3, 4, 5\}, \{5, 6, 7\}, \{6, 7, 8\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$	$[1 \ 1 \ 2 \ 1 \ 2 \ 2 \ 2 \ 1]$	11

The examples featured in Table 4.4 shows that '+1' couplings are perfectly balanced with degrees 3. The complexity of PCTF3D with '+1' couplings is in $\mathcal{O}(IMR^2)$ which is lower by a M^2 factor compared to CTF3D.

Random strategies

As its name suggests, a random strategy consists in picking T triplets randomly while ensuring that the resulting hypergraph is connected. The advantage of random strategies is to be able to control the complexity of PCTF3D thanks to the choice of the number of triplets. However, those strategies introduces randomness which can alter reproducibility of the method. If the constraints of Assumptions 4.2.1 and 4.2.2 were relaxed, the degrees would follow a hypergeometric distribution. Therefore, in that case, the variance of d_m for a random coupling would be the following:

$$\text{Var}[d_m] = \frac{T \binom{M-1}{2} \binom{M}{3} - \binom{M-1}{2} \binom{M}{3} - T}{\binom{M}{3}} = \frac{3T(M-3) \binom{M}{3} - T}{M^2 \binom{M}{3} - 1}$$

In order to compare the actual distribution of degrees, 1000 random couplings were generated for $M = 7$ for each value of number of triplets. These couplings give us $7 \times 1000 = 7000$ realizations of the random variable d_m . Figure 4.2 shows the distribution of d_m averaged over the 7000 realizations. This distribution was obtained for each $T \in \left[\left\lceil \frac{M}{2}, \binom{M}{3} \right\rceil \right]$ and is compared to its associated hyper-geometric distribution. This figure shows that the distribution of d_m is very close to the hyper-geometric (without constraints). Due to the constraint of Assumption 4.2.1, the difference between the d_m distribution and the hypergeometric distribution is higher for low-triplet couplings. It is more difficult to define a distribution that approximates the distribution of P . However, the empirical distribution of P is shown in Figure 4.3 which was computed for a set of 1000 couplings. This figure shows that the maximal number of pairs is achieved for nearly all couplings if $T \geq 15$. Knowing the distribution of d_m like in Figure 4.2 gives insights on the possible recoverability bounds for a value of T . Indeed, for the defective cases presented in Section 5.1.6, $d_m = 1$ is a critical value that cannot be achieved for $T \geq 22$ according to the distribution. Analogously, the empirical distribution of P gives insights on recoverability because the higher the number of pairs, the higher is the recoverability bound (see Equation (5.13)).

Table 4.4: Examples of '+1' couplings with their associated incidence matrix and degree sequence.

	$\mathcal{T}^{(+1)}$	\mathbf{V}	\mathbf{d}	P
$M = 4$	$\{\{1, 2, 3\}, \{2, 3, 4\}, \{3, 4, 1\}, \{4, 1, 2\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \end{bmatrix}$	$[3 \ 3 \ 3 \ 3]$	6
$M = 5$	$\{\{1, 2, 3\}, \{2, 3, 4\}, \{3, 4, 5\}, \{4, 5, 1\}, \{5, 1, 2\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 \end{bmatrix}$	$[3 \ 3 \ 3 \ 3 \ 3]$	10
$M = 6$	$\{\{1, 2, 3\}, \{2, 3, 4\}, \{3, 4, 5\}, \{4, 5, 6\}, \{5, 6, 1\}, \{6, 1, 2\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ & & & \vdots & & \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}$	$[3 \ 3 \ 3 \ 3 \ 3 \ 3]$	12
$M = 7$	$\{\{1, 2, 3\}, \{2, 3, 4\}, \{3, 4, 5\}, \{4, 5, 6\}, \{5, 6, 7\}, \{6, 7, 1\}, \{7, 1, 2\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ & & & \vdots & & & \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$	$[3 \ 3 \ 3 \ 3 \ 3 \ 3 \ 3]$	14
$M = 8$	$\{\{1, 2, 3\}, \{2, 3, 4\}, \{3, 4, 5\}, \{4, 5, 6\}, \{5, 6, 7\}, \{6, 7, 8\}, \{7, 8, 1\}, \{8, 1, 2\}\}$	$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ & & & \vdots & & & & \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$	$[3 \ 3 \ 3 \ 3 \ 3 \ 3 \ 3 \ 3]$	16

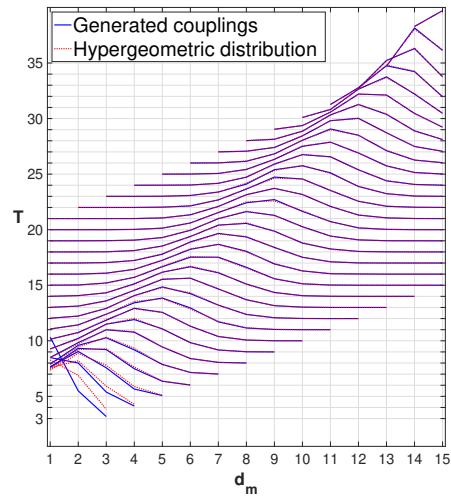


Figure 4.2: Degree distribution of 1000 random couplings for different number of triplets with $M = 7$. Empirical distributions are plotted with associated hyper-geometric distributions.

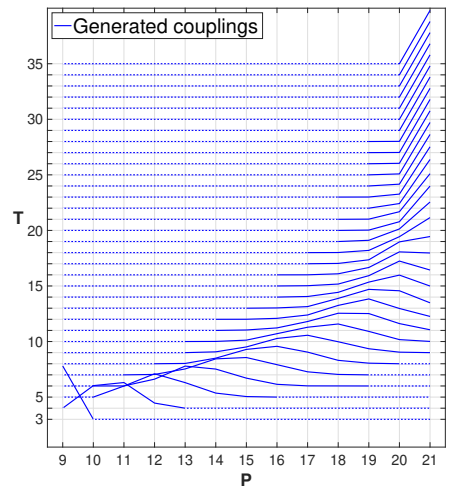


Figure 4.3: Empirical distribution of P for 1000 random couplings with different number of triplets and $M = 7$.

4.2.3 An algorithm for generating balanced couplings

The idea of balanced couplings is to ensure that all variables are coupled equally in the sense that the degree sequence is constant or nearly constant. Indeed, random couplings lead to random degree sequences which can lead to identifiability rank bounds much lower compared to strategies with constant degree sequence (see Section 5.1.6). It is not always possible to find a coupling with a constant degree sequence. Indeed, a necessary condition is that the value $\delta = \frac{3T}{M}$ is an integer. In this case, the degree sequence of a balanced coupling is

$$\mathbf{d} = \left[\delta \quad \dots \quad \delta \right]. \quad (4.4)$$

Note that this condition is easier to apply if M is a multiple of 3. If the condition $\delta = \frac{3T}{M}$ is satisfied, then there exist couplings with a constant degree sequence.

Otherwise, for all $T \in \left[\left[\frac{M}{2}, \binom{M}{3} \right] \right]$ such that $3T$ is not divided by M , there exist couplings with step-1 degree sequences. If integers (δ, ε) are defined such that $0 < \varepsilon < M$ and

$$3T = \delta M + \varepsilon,$$

then there exists \mathcal{T} with a step-1 degree sequence:

$$\mathbf{d} = \left[\delta \quad \dots \quad \delta \quad \overbrace{\delta + 1 \quad \dots \quad \delta + 1}^{\varepsilon} \right]. \quad (4.5)$$

A balanced hypergraph with T edges is then generated via an adaptation of a reconstruction algorithm of [46, Section 4]. The goal of Algorithm 4 is to generate incidence matrices $\mathbf{V} \in \{0, 1\}^{T \times M}$ such that variables are linked with a step-1 degree sequence (4.5). In the case where $3T = M\delta$, the algorithm will return the incidence matrix of a perfectly-balanced hypergraph.

To build such a hypergraph, the algorithm uses Lyndon words [68] in the alphabet of two letters $\{0, 1\}$ and their circular permutations.

Definition 4.2.13. Lyndon word (LW) – A Lyndon word (LW) is a word of length M denoted lw such that any circular permutation lw' of lw is strictly greater regarding the lexicographic order ($lw' > lw$). In our framework, a LW is a binary vector of length M containing exactly 3 times the letter 1.

For example, there are 3 LW for $M = 6$ which are:

$$\begin{aligned} lw_1 &= [0 \ 0 \ 0 \ 1 \ 1 \ 1], \\ lw_2 &= [0 \ 0 \ 1 \ 0 \ 1 \ 1], \\ lw_3 &= [0 \ 0 \ 1 \ 1 \ 0 \ 1]. \end{aligned}$$

Note that the LW definition does not allow for periodic patterns. This means that for $M = 6$, the word $[0 \ 1 \ 0 \ 1 \ 0 \ 1]$ is not considered a LW. Therefore, a LW has M different circular permutations (including itself) which permits to create a perfectly-balanced hypergraph. For example, for the LW $[0 \ 0 \ 0 \ 1 \ 1 \ 1]$ of length $M = 6$, the hypergraph \mathcal{T} below has a step-0 degree sequence equal to 3:

$$\mathbf{V} = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 \end{bmatrix}, \quad \mathcal{T} = \begin{aligned} &\{\{4, 5, 6\}, \\ &\{1, 5, 6\}, \\ &\{1, 2, 6\}, \\ &\{1, 2, 3\}, \\ &\{2, 3, 4\}, \\ &\{3, 4, 5\} \end{aligned} \quad (4.6)$$

$$\mathbf{d} = [3 \ 3 \ 3 \ 3 \ 3 \ 3]. \quad (4.7)$$

Lemma 4.2.14. *If two LW (lw_1, lw_2) share a circular permutation w , then necessarily $lw_1 = lw_2$.*

Proof. If w is a circular permutation of both LW, there exists p_1 and p_2 two circular permutations such that:

$$w = p_1(lw_1) \quad \text{and} \quad w = p_2(lw_2).$$

Hence, $p_1(lw_1) = p_2(lw_2)$ so there exists a circular permutation p such that $lw_1 = p(lw_2)$ and $lw_2 = p^{-1}(lw_1)$. If p is the identical permutation the proof is complete. If not, it means that lw_1 is a circular permutation of the LW lw_2 therefore $lw_1 > lw_2$ and conversely lw_2 is a circular permutation of the LW lw_1 so $lw_2 > lw_1$ which leads to a contradiction. Hence, p is the identical permutation and the proof is complete. \square

Note that the converse of this property states that for two different LW, the 2 sets of M triplets generated by circular permutations are disjoint.

Lemma 4.2.15. For a given $M \geq 3$, the number of LW is denoted $LW(M)$ and is given by:

$$LW(M) = \left\lfloor \frac{\binom{M}{3}}{M} \right\rfloor.$$

Proof. In [48], the number of LW was studied in a more general case for LW of densities k . Here, the density means the number of '1' contained in a LW and thus is equal to 3 for triplets. This leads to the following equation:

$$LW(M, k) = \frac{1}{M} \sum_{\substack{j=1 \\ j|k \text{ and } j|M}}^k \mu(j) \binom{\frac{M}{j}}{\frac{k}{j}},$$

where $\mu(j)$ is the Möbius function evaluated in j and where $j|k$ means " j is a divisor of k ". In our case $k = 3$ so the common divisors between M and 3 are either 1 if $3 \nmid M$ (M is not divided by 3) or $\{1, 3\}$ if M is a multiple of 3:

$$LW(M) = \begin{cases} \frac{1}{M} \binom{M}{3} - \frac{1}{M} \binom{\frac{M}{3}}{1} & \text{if } 3 \mid M \\ \frac{1}{M} \binom{M}{3} & \text{else} \end{cases}.$$

As $LW(M)$ is an integer for $3 \mid M$ and because $\frac{1}{M} \binom{\frac{M}{3}}{1} = \frac{1}{3} < 1$, the requested equation is obtained for both cases. \square

In the case of $3 \nmid M$, all $\binom{M}{3}$ triplets can be generated with LW and their circular permutations. Indeed, each of the $LW(M) = \frac{\binom{M}{3}}{M}$ LW counts for M different triplets which gives a total of $\binom{M}{3}$ triplets. However, if M is divided by 3, permuted LW only generate $\binom{M}{3} - \frac{M}{3}$ triplets. The remaining $\frac{M}{3}$ triplets can be generated with the circular permutation of the word:

$$\underbrace{0 \dots 0}_{\frac{M}{3}} 1 \underbrace{0 \dots 0}_{\frac{M}{3}} 1 \underbrace{0 \dots 0}_{\frac{M}{3}} 1. \tag{4.8}$$

Finally, we propose the Algorithm 4 based on [46, Section 4] which relies on circular permutations of $\lfloor \frac{T}{M} \rfloor$ LW (and the word (4.8) if necessary) and then add the missing triplets with permutation of the LW $'(0)^{M-3}(1)^3'$.

Algorithm 4: Balanced hypergraph generation
Input: M, \mathcal{T} .
Output: \mathbf{V}
Generate $\lfloor \frac{T}{M} \rfloor$ Lyndon words
Add the $M \lfloor \frac{T}{M} \rfloor$ circular permutations to \mathbf{V}
if M is divided by 3 and more than $\frac{M}{3}$ triplets are needed **then**
| Add the permutations of word (4.8) to \mathbf{V}
end
Add the missing \mathbf{V} rows with permutations of $(0)^{M-3}(1)^3$.

Remark 4.2.16. Note that the Algorithm 4 allows to generate random balanced couplings. The number of different couplings that can be generated depends on the number of different Lyndon word combinations, and thus (by Lemma 4.2.15) is equal to:

$$\binom{LW(M)}{\lfloor \frac{T}{M} \rfloor}$$

However, this algorithm may not explore the whole space of valid balanced couplings.

4.3 Numerical experiments

This section presents numerical experiments that were performed to study the performance of both PCTF3D and CTF3D. We intentionally do not present any results regarding flow cytometry experiments. Indeed, flow cytometry experiments will be performed and studied in Chapter 6.

To study the performance of the proposed methods, several criteria were proposed. First, it must be noted that the Frobenius norm between two tensors is often impossible to compute in practice, because of the prohibitive size of tensors:

$$D_M := \left\| \mathcal{H} - \left[\widehat{\boldsymbol{\lambda}}; \widehat{\mathbf{A}}^{(1)}, \dots, \widehat{\mathbf{A}}^{(M)} \right] \right\|_F^2.$$

To palliate this issue, the Factor Match Score was defined in [1] and consists in finding the best permutation such that for $r \in \llbracket 1, R \rrbracket$ estimated factors $\widehat{\mathbf{a}}_r^{(m)}$ can be compared to the theoretical $\mathbf{a}_r^{(m)}$:

$$D_{\text{FMS}}(\mathcal{H}, \widehat{\mathcal{H}}) := \sum_{r=1}^R \prod_{m=1}^M \frac{\mathbf{a}_r^{(m)\top} \widehat{\mathbf{a}}_r^{(m)}}{\|\mathbf{a}_r^{(m)}\|_2 \|\widehat{\mathbf{a}}_r^{(m)}\|_2}. \quad (4.9)$$

A lower-order metric can also be used which consists in summing the errors over lower-order

marginals:

$$\text{Err}_{1\text{D}} := \sum_{m=1}^M \left\| \mathbf{h}^{(m)} - \widehat{\mathbf{h}}^{(m)} \right\|_2^2,$$

$$\text{Err}_{3\text{D}} := \sum_{\{j,k,\ell\} \in \mathcal{T}_{\text{all}}} \left\| \mathcal{H}^{(j k \ell)} - \left[\widehat{\boldsymbol{\lambda}}; \widehat{\mathbf{A}}^{(j)}, \widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{A}}^{(\ell)} \right] \right\|_F^2,$$

where $\mathbf{h}^{(m)}$ represents the m -th 1D marginal and $\widehat{\mathbf{h}}^{(m)}$ its estimation. Note that $\text{Err}_{3\text{D}}$ is computed over all possible triplets. It is equivalent with the cost function (3.24) solved by CTF3D. For PCTF3D, the 3D marginals distance is computed even for marginals not considered in partial couplings.

4.3.1 Random strategies and performance regarding the number of triplets

To study the performance of coupling strategies, PCTF3D was applied for all strategies of Table 4.5 with $M = 10$ dimensions synthetic datasets. Multivariate Gaussian distributions of rank $R = 20$ were generated randomly and added with weights $\boldsymbol{\lambda}$ together to create a theoretical histogram \mathcal{H} . Different number of samples $N \in \{10^4, 3 \cdot 10^4, 10^5, 3 \cdot 10^5, 10^6\}$ were generated and the histograms were computed with $I = 30$ bins per dimension. We set $T_1 = 10^3$ outer iterations and $T_2 = 20$ inner iterations for the parameters of PCTF3D (Algorithms 1 to 3) and the target rank equal to R . The reconstruction error $\text{Err}_{1\text{D}}$ and $\text{Err}_{3\text{D}}$ were evaluated for each strategy and averaged over 10 experiments. Figure 4.4 shows that the coupling strategy '1/8' – where $T = \frac{1}{8} \binom{M}{3}$ are chosen randomly – yields similar performance to the fully coupled strategy. This shows that partial coupling strategies are beneficial in terms of computational cost as the computational complexity is linear with the number of triplets in consideration.

4.3.2 Comparison between random and balanced couplings.

To compare random and balanced strategies, we performed a similar experiment than in the previous section. Indeed, 10 multivariate distributions were generated with the following parameters:

- $M = 20$ variables,
- $I = 30$ bins per dimension,
- $R_{\text{th}} = 40$ rank-one mixture.

Table 4.5: Coupling strategies for $M = 10$.

Strategy	T	Triples
+2	5	$\{1, 2, 3\}, \{3, 4, 5\}, \{5, 6, 7\}, \{7, 8, 9\}, \{9, 10, 1\}$
+1	10	$\{1, 2, 3\}, \{2, 3, 4\}, \dots, \{9, 10, 1\}, \{10, 1, 2\}$
1/8	$15 = \frac{120}{8}$	random triples
1/4	$30 = \frac{120}{4}$	random triples
1/2	$60 = \frac{120}{2}$	random triples
1	$120 = \binom{10}{3}$	all triples (\sim CTF3D)

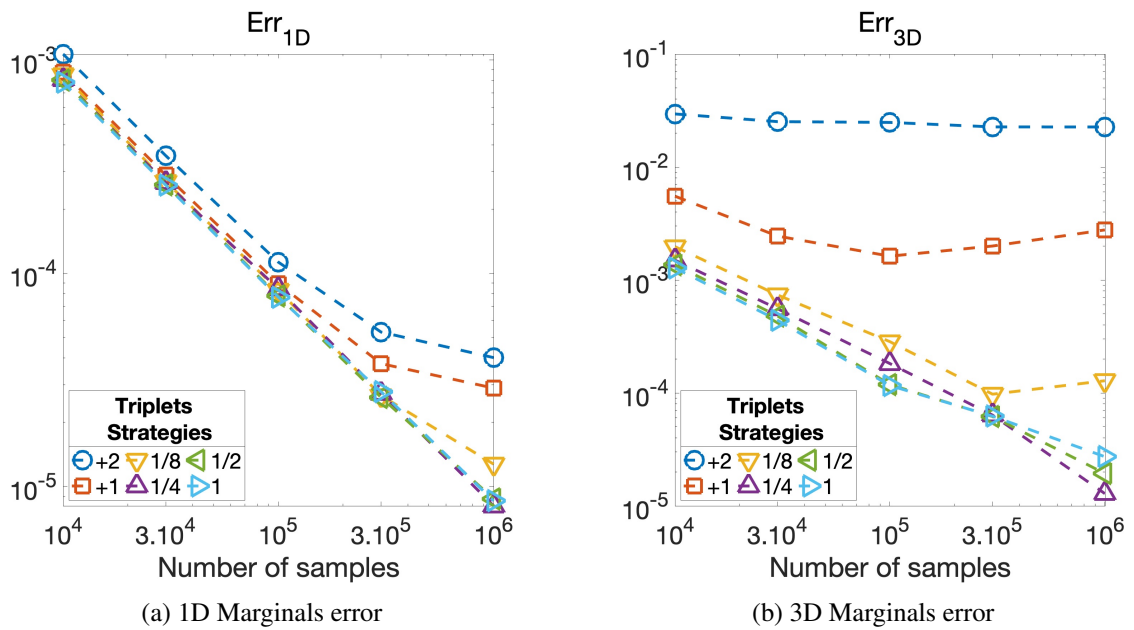


Figure 4.4: Error regarding the number of triplets for different coupling strategies.

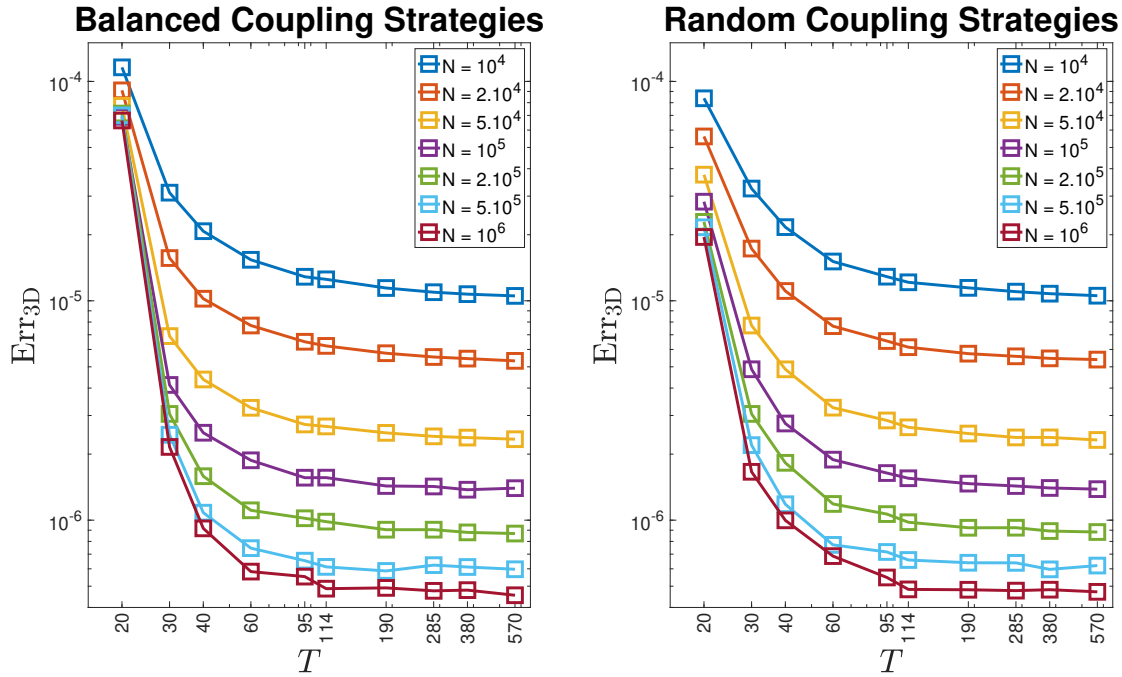


Figure 4.5: Error on 3D marginals for balanced and random coupling strategies.

Then, following these distributions, datasets with different sample sizes were generated:

$$N \in \{10^4, 2.10^4, 5.10^4, 10^5, 2.10^5, 5.10^5, 10^6\}.$$

For both strategies and for each distribution, a set of triplets was generated for different number of triplets $T \in \{20, 30, 40, 60, 95, 114, 190, 285, 380, 570\}$. PCTF3D is then applied with $R = R_{\text{th}}$, $T_1 = 1000$ outer iterations and $T_2 = 20$ inner iterations of ADMM.

The results of this experiment are shown in Figure 4.5 and show the same results as for the previous experiment. Indeed, error on marginals reaches a constant level above a certain number of triplets. This level is the same for both strategies. However, for low values of T , balanced strategies seems to have worse performance. As shown in Chapter 5, random strategies have better recoverability guarantees especially for low number of triplets. In that sense, there is a compromise between the choice of the coupling strategy. Note that because there are fewer balanced coupling strategies, it may be possible that the estimation of the error for balanced couplings is worse than for random couplings.

4.3.3 Comparison with a KLD-based method

Small cluster sensitivity experiment

In flow cytometry, end users search for small cell populations. Thus, an experiment to test the sensitivity of three density estimation methods is proposed. The first two methods are CTF3D and PCTF3D and will be compared with a KLD-based method: SQUAREM-PMF. This method based on a maximum likelihood estimator also estimates the factors of a discrete NBM. However, SQUAREM-PMF does not use histogram to estimate the CPD and rather use directly the data from the observation matrix. More details on SQUAREM-PMF are provided in [23]. In the following, PCTF3D was applied with half of all possible triplets (strategy '1/2'), chosen randomly.

To study the sensitivity of the three methods, $R_{\text{th}} = 3$ multivariate discrete Gaussian random variables were generated with $M = 7$ and $I = 20$. For each theoretical distribution, $N = 10^5$ samples were generated with 8 different proportions:

$$\lambda_1 \in \left\{ 0.002, 0.005, 0.01, 0.02, 0.05, 0.1, 0.2, \frac{1}{3} \right\}, \quad \lambda_2 = \frac{1}{3} \quad \text{and} \quad \lambda_3 = 1 - (\lambda_1 + \lambda_2).$$

With theoretical datasets generated, the three methods were applied with increasing rank until the smallest component was part of the estimate. Therefore, this procedure stopped if there existed an estimated CPD rank-one term such that:

$$\prod_{m=1}^M \left\| \hat{\mathbf{a}}_r^{(m)} - \mathbf{a}_r^{(m)} \right\|_2 < 10^{-4}.$$

Figure 4.6 shows the minimum rank required for each experiment averaged over 100 trials. It can be seen that, compared to CTF3D and PCTF3D, SQUAREM-PMF is more sensitive to small clusters as R must be higher if one wants to obtain a small component with (P)CTF3D. Unlike CTF3D and PCTF3D, for SQUAREM-PMF, the estimation rank can be chosen close to the expected number of clusters. In terms of computational load, SQUAREM-PMF is faster on harder problems because the smallest rank-one term is found with $R = 3$. As λ_1 increases, the minimum rank required is close to R_{th} for all methods and thus SQUAREM-PMF is slower than the other methods.

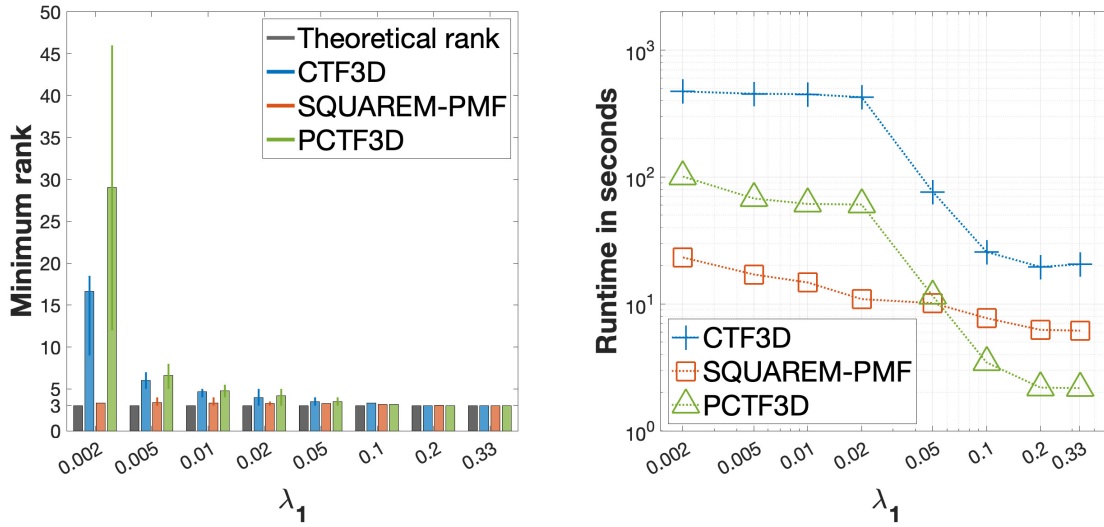


Figure 4.6: Results for the sensitivity experiment. 100 trials were run on datasets with $M = 7$ variables featuring $R_{th} = 3$ theoretical rank-one terms. **Left plot:** Mean value over 100 trials of the minimal rank that provides the desired rank-one term. **Right plot:** Median value over 100 trials of the runtime in seconds.

Runtime analysis

The three examined algorithms have different complexities. For example, the computational time of CTF3D and PCTF3D is less dependent on N in comparison with SQUAREM-PMF. In order to verify this empirically, three experiments were ran on each method and averaged the results over 100 trials. The performance criteria were the runtime and the FMS (4.9).

Evolution with respect to M – In the first experiment, $N = 10^5$ samples were generated from random factors with $M \in \llbracket 4, 10 \rrbracket$, $I \in \{5, 10, 15\}$ and $R = 5$. For each setting, the three proposed methods were run to obtain a CPD of rank $R = 5$. The left plots of Figure 4.7 show that CTF3D’s run times increase with M , and CTF3D becomes slower than SQUAREM-PMF for the highest values of M (note that the total number of elements in 3D histograms becomes comparable with the number of data points N). The run times for SQUAREM-PMF do not depend strongly on I .

Evolution with respect to R – Next, for 100 random observation matrices of size $(N = 10^4) \times (M = 5)$, the proposed methods were run to obtain a CPD with different ranks $R \in \llbracket 3, 20 \rrbracket$. The middle plots of Figure 4.7 show that the run times increase with R . For small values of

R , SQUAREM-PMF converges with less iterations and thus has a lower computation time for lower ranks; for higher ranks CTF3D and PCTF3D are faster. For CTF3D and PCTF3D, a drop of the runtime was observed for $I = 4$ and $R > 20$ (not shown), which can be explained by a loss of identifiability after those ranks.

Evolution with respect to N – Finally, to study the complexity regarding N , the proposed methods were applied on 100 datasets with $M = 6$ to obtain a rank $R = 5$ CPD. The right plots of Figure 4.7 show that the runtime of CTF3D does not depend on N , while achieving similar performance to that of SQUAREM-PMF. However, SQUAREM-PMF’s run times increase considerably with T (some trials took a few hours to run).

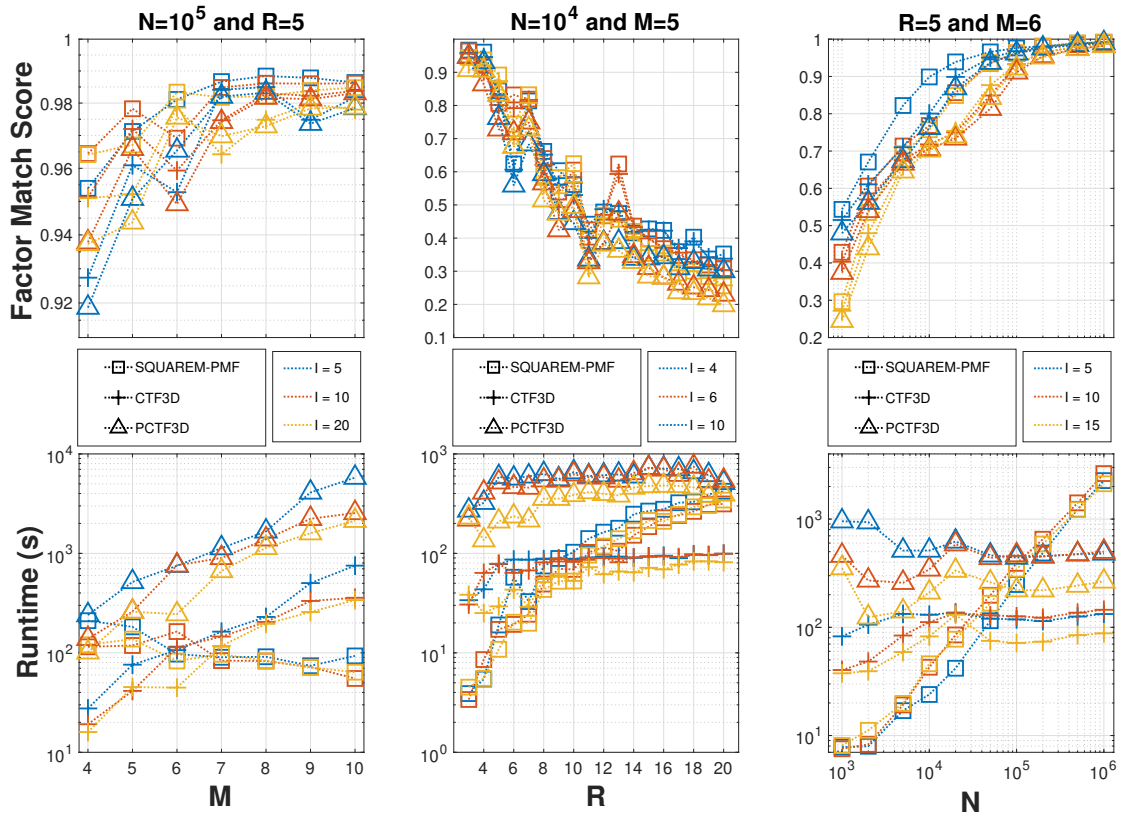


Figure 4.7: Comparison of CTF3D, PCTF3D and SQUAREM-PMF in terms of runtime and FMS. **Left plots:** M experiment with $M \in \llbracket 4, 10 \rrbracket$. **Middle plots:** R experiment with $R \in \llbracket 3, 20 \rrbracket$. **Right plots:** experiment regarding N .

4.4 Conclusion of Chapter 4

In this chapter, PCTF3D was proposed and presented as a new algorithm for probability mass tensor estimation that uses partially coupled tensor factorization. This method solves the problem of the third and last curse of dimensionality presented in our study. This algorithm permits to reduce and control the number of marginals considered in the coupling, hence the complexity of the algorithm.

This algorithm introduced the concept of coupling: a set of triplets of variables that defines the 3D marginals coupled in PCTF3D. Couplings can be formalized as hypergraphs which are graphs whose edges contains possibly more than two variables. In our case, a set of triplets permits to create 3-uniform hypergraphs or 3-graphs. With that formalism, several coupling strategies and their basic hypergraph properties were proposed. For example, an algorithm for balanced coupling generation was proposed. This algorithm permits to ensure that every variable is represented evenly in terms of occurrence in the hypergraph.

Finally, numerical experiments were performed to study the performance of PCTF3D. In a first experiment, partially coupled tensor factorization is shown to reduce the computational load without a big loss of estimation performance. In the second experiment, random couplings were compared to balanced couplings. Finally, PCTF3D was compared to a KLD-based method. Results show that a distance like the KLD may be more appropriate in the search of rare populations. However, the KLD method presented in this study may be hard to apply in practice to large datasets. In Appendix B, an application of PCTF3D to the MNIST dataset is reported. The key point here is to show that PCTF3D can estimate a PMF in dimension 256. The results of the decomposition are used to perform the digit classification using MAP-like classifiers that yields performances similar to those achieved by a Bayes classifier.

In the next chapter, the problem of uniqueness of the method PCTF3D will be examined. We will see that recoverability or identifiability of this method depend on the choice of triplets.

Chapter 5

Identifiability of coupled tensor models

Contents

5.1 Recoverability results	100
5.1.1 Reduced re-parametrization of factor matrices	100
5.1.2 Jacobian structure and algorithm for recoverability bound search	102
5.1.3 Recoverability necessary condition and number of degrees of freedom	104
5.1.4 Recoverability results for fully coupled tensor factorization	105
5.1.5 Recoverability results for randomly selected triplets	105
5.1.6 Analysis of defective cases	106
5.1.7 Recoverability results for balanced couplings	110
5.2 Identifiability of Cartesian product coupling	111
5.2.1 Re-parametrization of the coupled model and identifiability results	113
5.2.2 Identifiability of the equivalent tensor model	114
5.2.3 An algebraic proof of identifiability	115
5.3 Application of Theorem 5.2.3	119
5.3.1 Identifiability result in the general case of variable partition coupling	119
5.3.2 Identifiability result for an even partition of variables	120
5.3.3 Identifiability result for the fully coupled case	121
5.4 Conclusion of Chapter 5	121

In this chapter, we seek to find identifiability conditions on the partially coupled model introduced in Chapter 4. To do this, we address the problem of identifiability of the following coupled tensor decomposition model:

$$\begin{aligned}
 \mathcal{H}^{(jk\ell)} &= \left[\boldsymbol{\lambda}; \mathbf{A}^{(j)}, \mathbf{A}^{(k)}, \mathbf{A}^{(\ell)} \right] \quad \text{for } \{j, k, \ell\} \in \mathcal{T} \\
 \text{s.t. } \boldsymbol{\lambda} &\geq 0, \quad \mathbf{A}^{(1)} \geq 0, \quad \dots, \quad \mathbf{A}^{(M)} \geq 0 \\
 \mathbb{1}_R^\top \boldsymbol{\lambda} &= 1, \quad \mathbb{1}_I^\top \mathbf{A}^{(1)} = \mathbb{1}_R^\top, \quad \dots, \quad \mathbb{1}_I^\top \mathbf{A}^{(M)} = \mathbb{1}_R^\top,
 \end{aligned} \tag{5.1}$$

where \mathcal{T} denotes a valid coupling in the sense of PCTF3D (see Section 4.2).

First, the recoverability of coupled models is examined. To do this, the Jacobian of the parametrization of our problem is analyzed. An algorithm that computes the maximum recoverable rank is proposed. This algorithm gives a sufficient recoverability condition for specific cases. By applying this algorithm for different coupling strategies and sizes of tensors, we obtain recoverability results that depend on coupling strategies.

Secondly, we present a proof for a sufficient identifiability condition for a special class of couplings. The proof of this result relies on the re-parametrization of the model and allows the exploitation of the identifiability results for non-negative CPD provided in [83]. The sufficient identifiability condition will be applied to special cases and it will be shown that the identifiability guarantees of [74] can be improved.

5.1 Recoverability results

In this section, we show how to embed the coupled CP decomposition of 3D marginals (5.1) in the framework of additive decompositions described in Section 2.3. Because of simplex constraints and especially the sum-to-one constraints on the factors, the space of parameters needs to be reduced as shown in the next subsection.

5.1.1 Reduced re-parametrization of factor matrices

Relaxing constraint on the loading vector First, note that the sum-to-one equality constraint on $\boldsymbol{\lambda}$ can be relaxed in (5.1) for the context of recoverability and identifiability as shown below. Consider a coupled CP decomposition model of order-3 tensors $\{\mathcal{H}^{(jk\ell)}\}_{\{j,k,\ell\} \in \mathcal{T}}$ of size $I \times I \times I$

of the form:

$$\begin{aligned} \mathcal{H}^{(j k \ell)} &= \left[\left[\boldsymbol{\lambda}; \mathbf{A}^{(j)}, \mathbf{A}^{(k)}, \mathbf{A}^{(\ell)} \right] \right], \quad \text{for all } \{j, k, \ell\} \in \mathcal{T}, \\ \text{where } \mathbb{1}_I^\top \mathbf{A}^{(m)} &= \mathbb{1}_R^\top, \mathbf{A}^{(m)} \geq 0, \boldsymbol{\lambda} \geq 0. \end{aligned} \quad (5.2)$$

For each tensor in (5.2), the sum of elements is given by

$$\left(\left(\left(\mathcal{H}^{(j k \ell)} \bullet_3 \mathbb{1}_I^\top \right) \bullet_2 \mathbb{1}_I^\top \right) \bullet_1 \mathbb{1}_I^\top \right) = \sum_{r=1}^R \lambda_r.$$

This implies the following.

1. First, for probability tensors (see (5.1)), relaxing the sum-to-one constraint on $\boldsymbol{\lambda}$ does not introduce new decompositions (5.2). This is due to the fact that the sum-to-one non-negative factors guarantee scaling ambiguities.
2. Second, for any $\alpha > 0$, the simultaneously scaled tensors $\mathcal{Y}^{(j k \ell)} = \alpha \mathcal{H}^{(j k \ell)}$ admit decompositions

$$\mathcal{Y}^{(j k \ell)} = \left[\left[\alpha \boldsymbol{\lambda}; \mathbf{A}^{(j)}, \mathbf{A}^{(k)}, \mathbf{A}^{(\ell)} \right] \right],$$

and these are the only possible decompositions of $\mathcal{Y}^{(j k \ell)}$.

Therefore, studying recoverability or identifiability of (5.1) is equivalent to studying the coupled model (5.2).

Re-parametrization of constraints on factor matrices Second, simplex constraints on factor matrices can be removed by truncation. Let us consider the parametrization \mathcal{P} which maps a vector $\underline{\mathbf{a}} \in \mathbb{R}^{I-1}$ onto a sum-to-one vector $\mathbf{a} = \mathcal{P}(\underline{\mathbf{a}})$:

$$\mathcal{P}(\underline{\mathbf{a}}) = \left[a_1 \quad \cdots \quad a_{I-1} \quad 1 - \sum_{i=1}^{I-1} a_i \right]^\top. \quad (5.3)$$

With abuse of notation, $\mathcal{P}(\underline{\mathbf{A}})$ denotes the matrix of size $I \times R$ with \mathcal{P} applied column-wise to vectors of size $I - 1$. With truncated factors $\underline{\mathbf{a}}_r^{(m)} = \left[a_{r,1}^{(m)} \cdots a_{r,I-1}^{(m)} \right]^\top$, it is possible to define a vector of parameters $\boldsymbol{\theta}$ such that:

$$\boldsymbol{\theta} = \text{vec} \left(\lambda_1, \underline{\mathbf{a}}_1^{(1)}, \dots, \underline{\mathbf{a}}_1^{(M)}, \dots, \lambda_R, \underline{\mathbf{a}}_R^{(1)}, \dots, \underline{\mathbf{a}}_R^{(M)} \right), \quad (5.4)$$

and a set of parameters Θ (of positive Lebesgue measure) defined by

$$\Theta = \left\{ \boldsymbol{\theta} \in \mathbb{R}^{(M(I-1)+1)R} \mid \boldsymbol{\lambda} \geq 0, \mathbf{a}_r^{(m)} \geq 0, \mathbb{1}_I^\top \mathbf{a}_r^{(m)} \leq 1 \right\}. \quad (5.5)$$

It is possible to define the function μ that maps the parameters onto the set of 3D marginals of \mathcal{T} :

$$\mu : \begin{cases} \mathbb{R}^{(M(I-1)+1)R} & \longrightarrow & \mathbb{R}^{TI^3} \\ \boldsymbol{\theta} & \longmapsto & \mathbf{y} = \text{vec}(\mathcal{H}^{(\tau_1)}, \dots, \mathcal{H}^{(\tau_T)}) \end{cases}, \quad (5.6)$$

where

$$\mathcal{H}^{(\tau_i)} = \left[\boldsymbol{\lambda}; \mathcal{P}(\underline{\mathbf{A}}^{(j_\tau)}), \mathcal{P}(\underline{\mathbf{A}}^{(k_\tau)}), \mathcal{P}(\underline{\mathbf{A}}^{(\ell_\tau)}) \right], \quad (5.7)$$

and the reduced factor matrices $\underline{\mathbf{A}}^{(m)} \in \mathbb{R}^{(I-1) \times R}$ are defined as $\underline{\mathbf{A}}^{(m)} = \begin{bmatrix} \mathbf{a}_1^{(m)} & \dots & \mathbf{a}_R^{(m)} \end{bmatrix}$. Note that the model (μ, Θ) defined by (5.5)-(5.7) is a polynomial additive model as introduced in Section 2.3.2.

5.1.2 Jacobian structure and algorithm for recoverability bound search

In Section 2.3.2, we saw that the Jacobian matrix for an additive model can be separated column-wise in R blocks (2.6). But because the model parametrized by (5.6) is coupled, $\mathcal{J}_\mu(\boldsymbol{\theta})$ can also be separated in T blocks where $T = \text{Card}(\mathcal{T})$. Each block represents an observed marginal:

$$\mathcal{J}_\mu(\boldsymbol{\theta}) = \begin{bmatrix} \mathcal{J}_\mu^{(\tau_1)}(\boldsymbol{\theta}) \\ \vdots \\ \mathcal{J}_\mu^{(\tau_T)}(\boldsymbol{\theta}) \end{bmatrix} \quad (5.8)$$

where $\mathcal{J}_\mu^{(\tau_i)}(\boldsymbol{\theta})$ denotes the Jacobian of the parametrization of $\mathcal{H}^{(\tau_i)}$ defined in Equation (5.7). To give an example, if we consider $M = 4$ and the coupling strategy $\mathcal{T} = \{\{1, 3, 4\}, \{2, 3, 4\}\}$, the Jacobian matrix of the r -th rank one term is defined by

$$\mathcal{J}_\mu(\boldsymbol{\theta}_r) = \begin{bmatrix} \mathbf{a}_r^{(4)} \otimes \mathbf{a}_r^{(3)} \otimes \mathbf{a}_r^{(1)} & \\ \mathbf{a}_r^{(4)} \otimes \mathbf{a}_r^{(3)} \otimes \mathbf{a}_r^{(2)} & \lambda_r \mathbf{B}_r \end{bmatrix},$$

where

$$\mathbf{B}_r = \begin{bmatrix} \mathbf{a}_r^{(4)} \otimes \mathbf{a}_r^{(3)} \otimes \mathcal{J}_{\mathcal{P}} & \mathbf{0}_{I^3 \times I} & \mathbf{a}_r^{(4)} \otimes \mathcal{J}_{\mathcal{P}} \otimes \mathbf{a}_r^{(1)} & \mathcal{J}_{\mathcal{P}} \otimes \mathbf{a}_r^{(3)} \otimes \mathbf{a}_r^{(1)} \\ \mathbf{0}_{I^3 \times I} & \mathbf{a}_r^{(4)} \otimes \mathbf{a}_r^{(3)} \otimes \mathcal{J}_{\mathcal{P}} & \mathbf{a}_r^{(4)} \otimes \mathcal{J}_{\mathcal{P}} \otimes \mathbf{a}_r^{(2)} & \mathcal{J}_{\mathcal{P}} \otimes \mathbf{a}_r^{(3)} \otimes \mathbf{a}_r^{(2)} \end{bmatrix}$$

and where $\mathcal{J}_{\mathcal{P}}$ denotes the Jacobian matrix of the projection \mathcal{P} (5.3)

$$\mathcal{J}_{\mathcal{P}} = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ -1 & \dots & -1 & \end{bmatrix}. \quad (5.9)$$

Note that the structure of \mathbf{B} in (5.1.2) is directly linked to the incidence matrix \mathbf{V} of \mathcal{T} .

$$\mathbf{V} = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{bmatrix}$$

This means that, excluding λ , each row block of $\mathcal{J}_{\mu}^{(\tau_t)}(\theta_r)$ only contains 3 non-zeros blocks out of M .

Following the proposition on recoverability in Section 2.3.2 (and similarly to the approach in [31]), an algorithm that returns R_{\max} in the case of a coupled CP model is proposed. The principle of Algorithm 5 is to increase the rank of a randomly generated decomposition while (2.5) holds for the model defined by (5.5)–(5.6).

Algorithm 5: Maximum recoverable rank R_{\max} for coupled CP model (5.6)

Input: M, I, \mathcal{T}

Output: $R_{\max} = R$

Initialization: $R = 1, \theta = \text{vec}(\lambda_1, \underline{\mathbf{a}}_1^{(1)}, \dots, \underline{\mathbf{a}}_1^{(M)})$ random satisfying (5.5).

while $\text{rank}(\mathcal{J}_{\mu}(\theta)) = R((I-1)M+1)$ **do**

$\theta_{R+1} = \text{vec}(\lambda_{R+1}, \underline{\mathbf{a}}_{R+1}^{(1)}, \dots, \underline{\mathbf{a}}_{R+1}^{(M)})$ random, satisfying (5.5)

$\theta = \text{vec}(\theta_1, \dots, \theta_R, \theta_{R+1}), R = R + 1$

 Estimation of $\text{rank}(\mathcal{J}_{\mu}(\theta))$

end

Remark 5.1.1. In Algorithm 5, a parameter θ that satisfies (2.5) gives a certificate of generic recoverability for a given R . This follows from Remark 2.3.5. Such a point θ presents a computer

proof of recoverability, as long as the absence of numerical errors can be guaranteed (this can be done, for example, by choosing parameters with rational entries, as in [27]).

5.1.3 Recoverability necessary condition and number of degrees of freedom

Recall that by Proposition 2.3.4, if a CP model is recoverable, then the Jacobian matrix $\mathcal{J}_\mu(\boldsymbol{\theta})$ is full column rank. Therefore, a necessary condition for recoverability is that $\mathcal{J}_\mu(\boldsymbol{\theta})$ has more rows than columns, meaning that there are more observations than values to describe the model. This condition is described in [15] as the well-posedness of an algebraic model.

In the following proposition, we count the number of actual observations (degrees of freedom). While it may seem at first sight that the number of observations is TI^3 , it is actually smaller, due to redundancy of the information in different marginals.

Proposition 5.1.2. *For a coupling \mathcal{T} , the dimension of the image of the parametrization μ (from (5.6)) denoted $N_{obs}(\mathcal{T})$ and is given by:*

$$N_{obs}(\mathcal{T}) := \dim(\text{Im}(\mu)) = 1 + M(I - 1) + P(I - 1)^2 + T(I - 1)^3. \quad (5.10)$$

Proof. To count the number of possibly different observations of the set of T marginals, we separate the count for each order of lower-order marginals. The lower-order marginals can be obtained from the observations by contracting the 3D marginals with the vector of ones. In our case of 3D-marginals coupling, we must count up to order-3 marginals. There is one order-0 marginal which is the sum of all entries of \mathcal{Y} ($\mathbf{y} = \mu(\boldsymbol{\theta})$). Because, we relaxed the sum-to-one constraint on $\boldsymbol{\lambda}$, this sum is not fixed hence count as 1 in the sum (5.10).

For 1D-marginals, there are M different marginals. For the marginal $\mathbf{h}^{(m)}$, the sum over its entries is equal to the order-0 marginal. Therefore, the I values of $\mathbf{h}^{(m)}$ are constrained by one linear equation. Hence, each 1D-marginal contributes to $I - 1$ values.

With the same reasoning, an order-2 marginal contributes for $(I - 1)^2$ free entries as well as an order-3 marginal contributes for $(I - 1)^3$ new free entries. The proof is complete because there are T order-3 marginals and P different pairs of variables hence P order-2 marginals. \square

Proposition 5.1.3. *For a coupling strategy \mathcal{T} , if a model (μ, \mathbb{R}^n) with μ given in Equation (5.6) is recoverable, then*

$$R \leq \left\lfloor \frac{N_{obs}(\mathcal{T})}{M(I - 1) + 1} \right\rfloor. \quad (5.11)$$

Proof. Jacobian matrix $\mathcal{J}_\mu(\boldsymbol{\theta})$ cannot be full rank if the number of parameters of the model exceeds the dimension of the image of μ . Therefore and in addition to Proposition 5.1.2, we have

$$R(M(I-1) + 1) \leq N_{\text{obs}}(\mathcal{T}),$$

which leads to the proposed condition on R . \square

5.1.4 Recoverability results for fully coupled tensor factorization

When all triplets are considered in \mathcal{T} , PCTF3D becomes equivalent with full coupling hence CTF3D which was presented in [74]. In this case, the necessary condition for recoverability (5.11) states that the model is recoverable if:

$$R \leq \left\lfloor \frac{1 + M(I-1) + \binom{M}{2}(I-1)^2 + \binom{M}{3}(I-1)^3}{M(I-1) + 1} \right\rfloor. \quad (5.12)$$

By applying the Algorithm 5 to this case, it is possible to determine the recoverability bound for the fully coupled case. In Figure 5.1, the blue curve represents the most favorable identifiability sufficient condition proposed in [74]. The dotted orange curve represents Equation (5.12) which is a necessary condition for identifiability and recoverability. Finally, the orange full line represents the rank R_{max} obtained with Algorithm 5. Figure 5.1 shows that the necessary identifiability condition is way above the sufficient conditions provided by [74]. Moreover, recoverability results of Algorithm 5 are achieved for ranks up to the necessary condition of Equation (5.12).

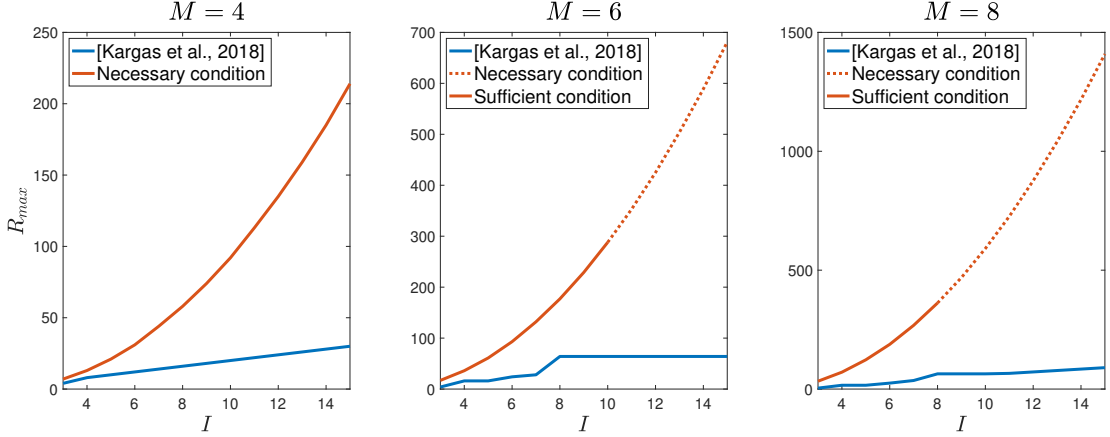
Our experiment suggests that the model recoverability is achieved for all ranks up to the necessary condition of Proposition 5.1.3, which permits us to formulate the following conjecture:

Conjecture 5.1.4. *For the full coupling \mathcal{T} , the coupled CP model is recoverable if:*

$$R \leq \left\lfloor \frac{N_{\text{obs}}(\mathcal{T})}{M(I-1) + 1} \right\rfloor.$$

5.1.5 Recoverability results for randomly selected triplets

When a random subset of triplets are considered, \mathcal{T} becomes random, hence the necessary condition for recoverability (5.11) also becomes possibly random (depending on the number of


 Figure 5.1: Recoverability bounds for different values of M and I in the case of the CTF3D.

triplets T):

$$R \leq \left\lfloor \frac{1 + M(I-1) + P(I-1)^2 + T(I-1)^3}{M(I-1) + 1} \right\rfloor. \quad (5.13)$$

For a coupling strategy of T triplets, the number of possible couplings is at most⁷ $\binom{M}{T}$. Therefore, even for small number of variables, it is difficult to study recoverability using Algorithm 5 exhaustively. To illustrate the behavior of the identifiability bound in the random case, we computed R_{\max} for 1000 different realizations of \mathcal{T} in the case of $M = 8$ variables and $I = 4$ bins per dimension. Figure 5.2 shows how the rank R_{\max} is distributed with respect to the values of T . Unlike the fully coupled case which is deterministic, the bound (5.13) is not always achieved as the left plot of Figure 5.2 suggests. As the number of triplets increases, the bound is increasing which is coherent. However, for a notable number of random couplings, R_{\max} is not exceeding the values 10 and 16. Those two cases are going to be properly studied in the following subsection. Consequently, the more triplets are considered, the less variability is expected in the recoverability bound.

5.1.6 Analysis of defective cases

In the following, let the sequence of degrees \mathbf{d} be sorted in ascending order, which can always be the case up to a permutation of variables. We present in this section two cases of a particular sequence of degrees that leads to identifiability and recoverability loss, which can explain the defective cases from the previous subsection.

⁷At most because some couplings are not valid (see Section 4.2)

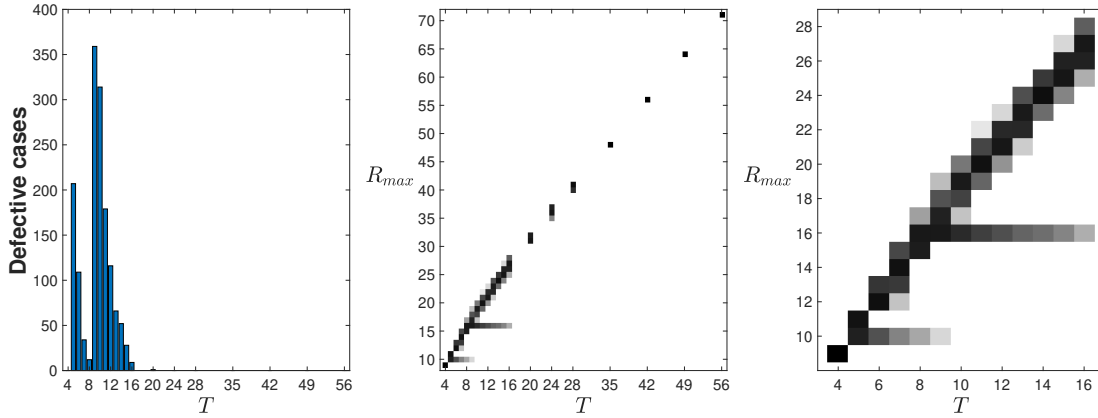


Figure 5.2: Recoverability bounds in the case of random partial coupling. **Left plot:** Number of defective cases (over 1000 realizations) for which R_{\max} is not equal to the maximal rank provided by the condition (5.11). **Middle plot:** Distribution of R_{\max} regarding T . **Right plot:** Zoom on the middle plot around the lower values of T . Defective cases appear distinctly at ranks $R_{\max} = 10$ and $R_{\max} = 16$.

First case: $d_1 = 1$ and $d_2 > 1$

This setup means that the first variable (and only this one) appears once in \mathcal{T} . This case can occur as soon as $T \leq \binom{M-1}{2}$ but is more and more likely as T is decreasing. After analyzing the couplings of Figure 5.2, those cases are leading to low recoverability bounds equal to $R_{\max} = 16$.

Proposition 5.1.5. *Let \mathcal{T} a coupling strategy such that $d_1 = 1$ and $d_2 > 1$. If the coupled CP model is identifiable, then we have that*

$$R_{\max} \leq I^2, \quad (5.14)$$

for this coupling strategy.

Proof. Let a coupled model with such coupling \mathcal{T} and suppose this model is identifiable. Then, it is necessary that $\mathcal{J}_\mu(\theta)$ is full rank. As seen in Section 5.1.2, the structure of $\mathcal{J}_\mu(\theta)$ is linked to the incidence matrix \mathbf{V} :

$$\mathbf{V} = \begin{bmatrix} 1 & v_{12} & \cdots & v_{1M} \\ 0 & v_{22} & \cdots & v_{2M} \\ \vdots & \vdots & & \vdots \\ 0 & v_{T2} & \cdots & v_{TM} \end{bmatrix}.$$

Table 5.1: Evolution over M of R_{\max} for most favorable coupling case with $d_1 = 1$.

M	4	5	6	7	...	15
$I = 3$	5	7	9	9	...	9
$I = 4$	8	13	16	16	...	16
$I = 6$	18	32	36	36	...	36

Therefore, it is possible to define $\mathbf{B}_r \in \mathbb{R}^{I^3 \times (I-1)}$ and $\mathbf{C}_r \in \mathbb{R}^{(T-1)I^3 \times (M-1)(I-1)+1}$ such that:

$$\mathcal{J}_\mu(\boldsymbol{\theta}_r) = \left[\begin{array}{c|c} \mathbf{B}_r & \dots \\ \hline \mathbf{0} & \mathbf{C}_r \end{array} \right]. \quad (5.15)$$

It follows that a necessary condition for $\mathcal{J}_\mu(\boldsymbol{\theta})$ is full rank is that the matrix $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 & \dots & \mathbf{B}_R \end{bmatrix}$ is full rank. Because of sum-to-one constraints and thus the structure of \mathcal{J}_μ , \mathbf{B} contains I^2 dependent rows. Therefore, necessarily the number of observations $I^3 - I^2$ should not be less than the number of parameters $R(I - 1)$ which completes the proof. \square

To verify this proposition, we proceed to the following experiment. For $M \in \llbracket 4, 15 \rrbracket$, Algorithm 5 was applied to the case where all triplets with variables $\{2, \dots, M\}$ are present and only the triplet $\{1, M - 1, M\}$ contains the first variable. This represents the most favorable case where one variable is present once in the coupling. The results of this experiment are plotted in Table 5.1. Even for this most favorable case, recoverability is guaranteed up to I^2 except for small number of M where the condition (5.13) is more restrictive.

Second case: $d_1 = d_2 = 1$ and $\tau_1 = \{1, 2, \ell_1\}$

In this case, the first two variables are present once in the first triplet and only in this triplet. Note that $d_1 = d_2 = 1$ ensures that $d_3 > 1$ because the hypergraph \mathcal{T} must be connected. Indeed, if $d_1 = d_2 = d_3 = 1$ and the first triplet is $\tau_1 = \{1, 2, 3\}$, then variables $\{1, 2, 3\}$ are not connected to the other variables. After analyzing the results of Figure 5.2, the couplings with $d_1 = d_2 = 1$ lead to a recoverability bound of $R_{\max} = 10$.

Table 5.2: Evolution over M of R_{\max} for most favorable coupling case with $d_1 = d_2 = 1$.

M	5	6	7	8	...	15
$I = 3$	4	6	6	6	...	6
$I = 4$	7	10	10	10	...	10
$I = 6$	16	21	21	21	...	21

Proposition 5.1.6. *Let \mathcal{T} a coupling strategy such that $d_1 = d_2 = 1$ and $\tau_1 = \{1, 2, \ell_1\}$. If the coupled CP model is identifiable, then we have*

$$R_{\max} \leq \frac{I(I+1)}{2}. \quad (5.16)$$

Proof. With the same method used in Proposition 5.1.5, let us consider an identifiable coupled model such that $d_1 = d_2 = 1$ and $\tau_1 = \{1, 2, \ell_1\}$. The incidence matrix \mathbf{V} is equal to

$$\mathbf{V} = \begin{bmatrix} 1 & 1 & v_{13} & \cdots & v_{1M} \\ 0 & 0 & v_{23} & \cdots & v_{2M} \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & v_{T3} & \cdots & v_{TM} \end{bmatrix} \quad (5.17)$$

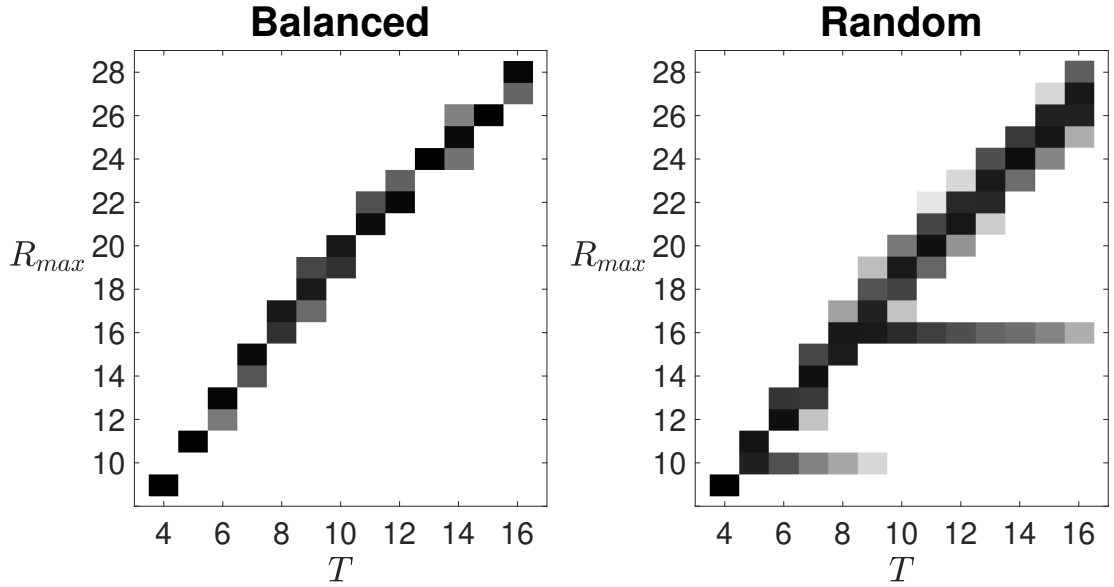
and therefore we can define $\mathbf{B}_r \in \mathbb{R}^{I^3 \times 2(I-1)}$ and $\mathbf{C}_r \in \mathbb{R}^{(T-1)I^3 \times (M-2)(I-1)+1}$ such that $\mathcal{J}_\mu(\boldsymbol{\theta})$ has the structure (5.15). Because now 2 factors are included, only I rows of $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 & \cdots & \mathbf{B}_R \end{bmatrix}$ are redundant. Hence, we obtain that $2R(I-1) \leq I^3 - I$ which leads to the condition (5.16). \square

By conducting a similar experiment, Table 5.2 shows that the bound verifies the condition (5.16). In this case, all triplets with variables $\{3, \dots, M\}$ were present and only the triplet $\{1, 2, M\}$ contains the first two variables. Maximal recoverability ranks are also bounded by $\frac{I(I+1)}{2}$ and (5.13) as expected.

To conclude this subsection, we would like to emphasize the fact that this list of defective cases may not be exhaustive. Indeed, for higher number of variables, it may be possible that the search of defective cases and the couplings that lead to those cases is a very difficult problem to tackle. However, for any number of variables, the two cases presented in this section can appear.

Table 5.3: Number of trials regarding the number of triplets considered in \mathcal{T} for the balanced couplings experiment.

T	4	5	6	7	8	9	10	11	12	13	14	15	16
Number of trials	36	10	60	30	187	30	13	247	13	7	195	21	313


 Figure 5.3: Recoverability bounds in the case of balanced partial couplings. **Left plot:** Distribution of R_{\max} regarding T for balanced couplings. **Right plot:** For comparison, the same distribution is shown for randomly chosen triplets.

5.1.7 Recoverability results for balanced couplings

As random strategies, balanced couplings are random hence the necessary condition for recoverability (5.11) is the same as for random strategies:

$$R \leq \left\lfloor \frac{1 + M(I-1) + P(I-1)^2 + T(I-1)^3}{M(I-1) + 1} \right\rfloor. \quad (5.18)$$

To study those couplings, we proceed to the same experiment shown in Section 5.1.5 where $M = 8$ and $I = 4$. However, because the number of Lyndon words is limited, it may be impossible to create 1000 different balanced couplings at a fixed value of T (see Table 5.3). Figure 5.3 presents the empirical distribution of R_{\max} versus T . For balanced strategies, the number of defect cases is not shown because the recoverability is guaranteed up to the necessary condition

(5.18) in every case of the experiment. Therefore, balanced strategies permit to reduce the variability of the recoverability bound compared to random couplings. Moreover, we can formulate and expand the results of this experiment in the following conjecture:

Conjecture 5.1.7. *Let \mathcal{T} a balanced coupling. If (5.18) is satisfied, then the coupled CP model is recoverable.*

Although this conjecture is very difficult to prove for every balanced coupling, tractable examples can be verified with Algorithm 5.

5.2 Identifiability of Cartesian product coupling

In this section, an identifiability condition for the model (5.1) is proved in the case of a particular type of coupling. This coupling is based on the partition of variables into 3 sets (and which serves as a base for proofs of Theorems 3.4.1 and 3.4.2). Let \mathcal{M}_1 , \mathcal{M}_2 and \mathcal{M}_3 be three non-empty sets such that:

$$\mathcal{M}_1 \cup \mathcal{M}_2 \cup \mathcal{M}_3 = \llbracket 1, M \rrbracket, \quad \text{and} \quad \mathcal{M}_1 \cap \mathcal{M}_2 = \mathcal{M}_1 \cap \mathcal{M}_3 = \mathcal{M}_2 \cap \mathcal{M}_3 = \emptyset.$$

In the following, M_1 (respectively M_2 and M_3) will denote the number of variables in \mathcal{M}_1 hence we have $\mathcal{M}_1 = \{j_1, \dots, j_{M_1}\}$ (respectively $\mathcal{M}_2 = \{k_1, \dots, k_{M_2}\}$ and $\mathcal{M}_3 = \{\ell_1, \dots, \ell_{M_3}\}$). Note that $M_1 + M_2 + M_3 = M$. As in [74], we consider the concatenation of the factor matrices:

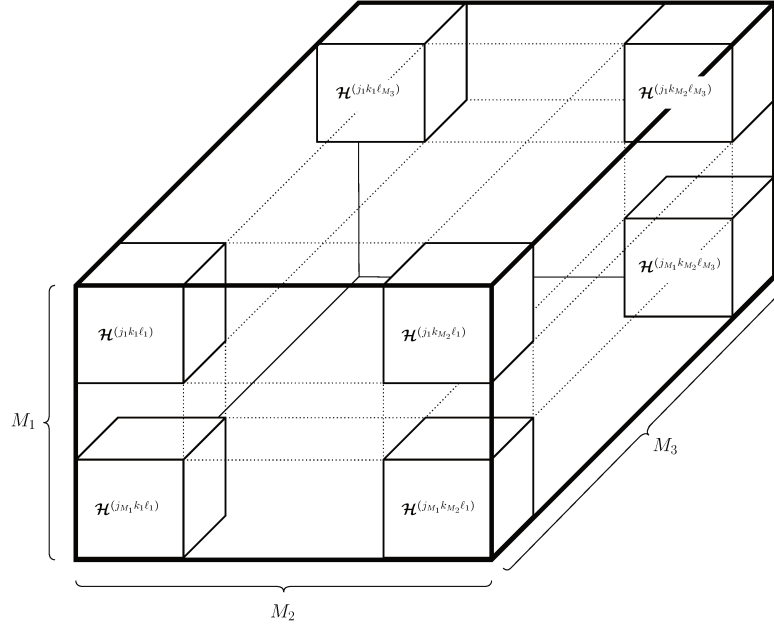
$$\mathbf{B}^{(1)} = \begin{bmatrix} \mathbf{A}^{(j_1)} \\ \vdots \\ \mathbf{A}^{(M_1)} \end{bmatrix} \in \mathbb{R}^{M_1 I \times R}, \quad \mathbf{B}^{(2)} = \begin{bmatrix} \mathbf{A}^{(k_1)} \\ \vdots \\ \mathbf{A}^{(M_2)} \end{bmatrix} \in \mathbb{R}^{M_2 I \times R}, \quad \mathbf{B}^{(3)} = \begin{bmatrix} \mathbf{A}^{(\ell_1)} \\ \vdots \\ \mathbf{A}^{(M_3)} \end{bmatrix} \in \mathbb{R}^{M_3 I \times R}. \quad (5.19)$$

In this setup, the order- M coupled model (5.1) is equivalent to the following order-3 model:

$$\mathcal{Y} = \llbracket \boldsymbol{\lambda}; \mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \mathbf{B}^{(3)} \rrbracket, \quad (5.20)$$

where \mathcal{Y} stacks all the marginals in \mathcal{T} (see Figure 5.4).

Example 5.2.1. *For $M = 5$ and a partition $\mathcal{M}_1 = \{1\}$, $\mathcal{M}_2 = \{2, 3\}$ and $\mathcal{M}_3 = \{4, 5\}$, the set of triplets becomes $\mathcal{T} = \{\{1, 2, 4\}, \{1, 2, 5\}, \{1, 3, 4\}, \{1, 3, 5\}\}$. Therefore, \mathcal{Y} contains $\{\mathcal{H}^{(124)}, \mathcal{H}^{(125)}, \mathcal{H}^{(134)}, \mathcal{H}^{(135)}\}$ as its blocks.*


 Figure 5.4: Structure of \mathcal{Y} as a concatenation of 3D marginals.

Because of the non-negativity constraint on factors $\mathbf{A}^{(m)}$, the factors of Equation (5.20) are also non-negative:

$$\lambda \geq 0, \quad \mathbf{B}^{(1)} \geq 0, \quad \mathbf{B}^{(2)} \geq 0, \quad \mathbf{B}^{(3)} \geq 0.$$

Concerning the sum-to-one constraints, the conditions $\mathbb{1}_I^\top \mathbf{A}^{(m)} = \mathbb{1}_R^\top$ can be summarized with the three following constraints:

$$\begin{aligned} (\mathbf{I}_{M_1} \otimes \mathbb{1}_I^\top) \mathbf{B}^{(1)} &= \mathbb{1}_{M_1} \otimes \mathbb{1}_R^\top, \\ (\mathbf{I}_{M_2} \otimes \mathbb{1}_I^\top) \mathbf{B}^{(2)} &= \mathbb{1}_{M_2} \otimes \mathbb{1}_R^\top, \\ \text{and } (\mathbf{I}_{M_3} \otimes \mathbb{1}_I^\top) \mathbf{B}^{(3)} &= \mathbb{1}_{M_3} \otimes \mathbb{1}_R^\top. \end{aligned} \tag{5.21}$$

The Cartesian product coupling was in fact used in the proofs of both Theorems 3.4.1 and 3.4.2.

Example 5.2.2. For Theorem 3.4.1, the partition is the following:

$$\mathcal{M}_1 = \{1\}, \quad \mathcal{M}_2 = \{2\} \quad \text{and} \quad \mathcal{M}_3 = \llbracket 3, M \rrbracket,$$

which leads to the following triplets $\mathcal{T} = \{\{1, 2, 3\}, \{1, 2, 4\}, \dots, \{1, 2, M\}\}$. Therefore, \mathcal{Y}

contains the set of $M - 2$ marginals

$$\{\mathcal{H}^{(123)}, \mathcal{H}^{(124)}, \dots, \mathcal{H}^{(12M)}\}.$$

For Theorem 3.4.2, a more balanced partition is used:

$$\mathcal{M}_1 = \left[\left[1, \left\lfloor \frac{M}{3} \right\rfloor \right] \right], \quad \mathcal{M}_2 = \left[\left[\left\lfloor \frac{M}{3} \right\rfloor + 1, \left\lfloor \frac{2M}{3} \right\rfloor \right] \right] \quad \text{and} \quad \mathcal{M}_3 = \left[\left[\left\lfloor \frac{2M}{3} \right\rfloor + 1, M \right] \right],$$

which means that \mathcal{Y} contains the set of $M_1 M_2 M_3$ marginals. An example of this partition of variable for $M = 5$ is given in Example 5.2.1.

The structure of \mathcal{Y} reveals three facts about the proofs of Theorems 3.4.1 and 3.4.2. First, the proofs presented in [74] are valid for a subset of marginals although CTF3D uses all 3D marginals in the CP model. Therefore, it seems possible to improve the Theorems 3.4.1 and 3.4.1. Secondly, as PCTF3D uses only a subset of triplets, the Theorems of [74] may be applied to PCTF3D. However, the sum-to-one constraints on the factors may have been overlooked in the article of CTF3D. In the following section, we will present a proof for generic identifiability of our coupled CP model. Finally, because the proofs presented show that an identifiability condition is valid for a set of triplets \mathcal{T} , it means that any coupling \mathcal{T}' such that $\mathcal{T} \subset \mathcal{T}'$ has the same identifiability results than the coupling \mathcal{T} .

5.2.1 Re-parametrization of the coupled model and identifiability results

In Section 5.1.1, the model (5.1) was relaxed to the model (μ, Θ) defined by (5.6) and (5.5). In the case of a Cartesian product coupling, the mapping μ in (5.6) is a vectorization⁸ of a tensor $\text{vec}(\mathcal{Y}) = \mu(\theta)$ (see Equation (5.20)) where the matrices $\mathbf{A}^{(m)}$ in (5.19) are defined by $\mathbf{A}^{(m)} = \mathcal{P}(\underline{\mathbf{A}}^{(m)})$. Therefore, the identifiability of the polynomial model (μ, Θ) can be studied via the identifiability of the tensor decomposition model of \mathcal{Y} . We can prove the following theorem.

Theorem 5.2.3. *Let $\mathcal{T} = \mathcal{M}_1 \times \mathcal{M}_2 \times \mathcal{M}_3$. If the non-negative CP model of sub-generic rank R for a tensor of size $((I - 1)M_1 + 1) \times ((I - 1)M_2 + 1) \times ((I - 1)M_3 + 1)$ is identifiable, then the model (5.1) of rank R is identifiable.*

We postpone the proof of this theorem to the following section. In fact, as it has been shown in Section 5.1.3, the tensor \mathcal{Y} of size $M_1 I \times M_2 I \times M_3 I$ contains marginals so its number of free

⁸where triplets (blocks of \mathcal{Y}) are ordered in the column-major order.

entries is reduced to $((I-1)M_1+1)((I-1)M_2+1)((I-1)M_3+1)$ (see Proposition 5.1.2). This can also be linked to the constraints of the CP model of \mathcal{Y} (5.21).

5.2.2 Identifiability of the equivalent tensor model

Let $\{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3\}$ a partition of M variables presented in Section 5.2. Let \mathcal{Z} be an order-3 tensor of size $(M_1(I-1)+1) \times (M_2(I-1)+1) \times (M_3(I-1)+1)$. From M matrices $\{\underline{\mathbf{D}}^{(m)}\}_{m=1}^M$ of size $I \times R$, we define three matrices denoted $\mathbf{D}^{(1)}$, $\mathbf{D}^{(2)}$ and $\mathbf{D}^{(3)}$:

$$\mathbf{D}^{(1)} = \begin{bmatrix} \alpha_1 \cdots \alpha_R \\ \underline{\mathbf{D}}^{(j_1)} \\ \vdots \\ \underline{\mathbf{D}}^{(j_{M_1})} \end{bmatrix}, \quad \mathbf{D}^{(2)} = \begin{bmatrix} \beta_1 \cdots \beta_R \\ \underline{\mathbf{D}}^{(k_1)} \\ \vdots \\ \underline{\mathbf{D}}^{(k_{M_2})} \end{bmatrix}, \quad \mathbf{D}^{(3)} = \begin{bmatrix} \gamma_1 \cdots \gamma_R \\ \underline{\mathbf{D}}^{(\ell_1)} \\ \vdots \\ \underline{\mathbf{D}}^{(\ell_{M_3})} \end{bmatrix}.$$

By defining $\alpha = [\alpha_1 \ \cdots \ \alpha_R] > 0$ (respectively β and γ), we consider the following model for \mathcal{Z} :

$$\begin{aligned} \mathcal{Z} &= \llbracket \mathbf{D}^{(1)}, \mathbf{D}^{(2)}, \mathbf{D}^{(3)} \rrbracket \\ \text{s.t. } \mathbf{D}^{(1)} &\geq 0, \quad \mathbf{D}^{(2)} \geq 0, \quad \mathbf{D}^{(3)} \geq 0, \\ \text{For all } j \in \mathcal{M}_1, \quad &\mathbb{1}_{I-1}^\top \underline{\mathbf{D}}^{(j)} \leq \alpha \\ \text{For all } k \in \mathcal{M}_2, \quad &\mathbb{1}_{I-1}^\top \underline{\mathbf{D}}^{(k)} \leq \beta \\ \text{For all } \ell \in \mathcal{M}_3, \quad &\mathbb{1}_{I-1}^\top \underline{\mathbf{D}}^{(\ell)} \leq \gamma \end{aligned} \tag{5.22}$$

For clarity, we define the three sets $\blacktriangle_{R,I}^{(1)}$, $\blacktriangle_{R,I}^{(2)}$ and $\blacktriangle_{R,I}^{(3)}$ such that:

$$\mathbf{D}^{(1)} \in \blacktriangle_{R,I}^{(1)}, \quad \mathbf{D}^{(2)} \in \blacktriangle_{R,I}^{(2)}, \quad \mathbf{D}^{(3)} \in \blacktriangle_{R,I}^{(3)},$$

The set $\blacktriangle_{R,I}^{(1)}$ represents the constraints on the factors in (5.22).

Lemma 5.2.4. *If the model (5.22) is identifiable for a given R , then the model (μ, Θ) (see (5.6) and (5.5)) for $\mathcal{T} = \mathcal{M}_1 \times \mathcal{M}_2 \times \mathcal{M}_3$ is also identifiable for the same R .*

Proof. Let us consider the three factor matrices:

$$\mathbf{C}^{(1)} = \underbrace{\begin{bmatrix} 1 & \dots & 1 \\ \underline{\mathbf{A}}^{(j_1)} \\ \vdots \\ \underline{\mathbf{A}}^{(j_{M_1})} \end{bmatrix}}_{(M_1(I-1)+1) \times R}, \quad \mathbf{C}^{(2)} = \underbrace{\begin{bmatrix} 1 & \dots & 1 \\ \underline{\mathbf{A}}^{(k_1)} \\ \vdots \\ \underline{\mathbf{A}}^{(k_{M_2})} \end{bmatrix}}_{(M_2(I-1)+1) \times R}, \quad \mathbf{C}^{(3)} = \underbrace{\begin{bmatrix} 1 & \dots & 1 \\ \underline{\mathbf{A}}^{(\ell_1)} \\ \vdots \\ \underline{\mathbf{A}}^{(\ell_{M_3})} \end{bmatrix}}_{(M_3(I-1)+1) \times R},$$

and the model (μ_2, Θ) defined on the same set of parameters (5.5) and such that

$$\mu_2(\theta) = \text{vec}(\mathcal{Z}) \quad \text{where } \mathcal{Z} = \llbracket \lambda; \mathbf{C}^{(1)}, \mathbf{C}^{(2)}, \mathbf{C}^{(3)} \rrbracket. \quad (5.23)$$

The model (μ_2, Θ) is a re-parametrization of (5.22) with scaling ambiguities removed by normalizing the first row of $\mathbf{D}^{(1)}$, $\mathbf{D}^{(2)}$ and $\mathbf{D}^{(3)}$. Therefore, (5.22) is essentially unique if and only if (5.23) is essentially unique.

Now, let \mathbf{Q}_1 (respectively \mathbf{Q}_2 and \mathbf{Q}_3) the matrix of size $M_1 I \times (M_1(I-1)+1)$ (respectively $M_2 I \times (M_2(I-1)+1)$ and $M_3 I \times (M_3(I-1)+1)$), defined as in Figure 5.5.

Then, the following factor matrices:

$$\mathbf{B}^{(1)} = \mathbf{Q}_1 \mathbf{C}^{(1)}, \quad \mathbf{B}^{(2)} = \mathbf{Q}_2 \mathbf{C}^{(2)} \quad \text{and} \quad \mathbf{B}^{(3)} = \mathbf{Q}_3 \mathbf{C}^{(3)}, \quad (5.24)$$

are non-negative and satisfy the constraints of (5.21). Vice-versa, for any non-negative factor matrix $\mathbf{B}^{(m)}$ satisfying (5.21), it is possible to retrieve $\underline{\mathbf{C}}^{(m)}$ by truncation of blocks of $\mathbf{B}^{(m)}$.

Finally, the mappings $\mu(\theta) = \text{vec}(\mathcal{Y})$ and $\mu_2(\theta) = \text{vec}(\mathcal{Z})$ are linked in the case of a Cartesian product coupling:

$$\mu(\theta) = \mathbf{Q}_3 \otimes \mathbf{Q}_2 \otimes \mathbf{Q}_1 \mu_2(\theta).$$

As each matrix \mathbf{Q} is full column rank, the operator $\mathbf{Q}_3 \otimes \mathbf{Q}_2 \otimes \mathbf{Q}_1$ is also full column rank. This proves the equivalence between identifiability of both models (μ, Θ) and (μ_2, Θ) . \square

5.2.3 An algebraic proof of identifiability

Recall that the identifiability of an order-3 CP model implies the identifiability of a non-negative order-3 CPD (see Section 2.2.2). In this section, we adapt the proof of that results to our re-parametrization in order to take properly into account the sum-to-one constraints. As we will

$$\mathbf{Q}_1 = \begin{bmatrix} 0 & & & & \\ \vdots & & & & \\ 0 & \mathcal{J}_{\mathcal{P}} & \mathbf{0}_{I \times I-1} & \dots & \mathbf{0}_{I \times I-1} \\ 1 & & & & \\ \hline 0 & & & & \\ \vdots & & & & \\ 0 & \mathbf{0}_{I \times I-1} & \mathcal{J}_{\mathcal{P}} & \dots & \vdots \\ 1 & & & & \\ \hline 0 & \vdots & \dots & \dots & \mathbf{0}_{I \times I-1} \\ \vdots & \vdots & \dots & \dots & \\ 0 & \vdots & \dots & \dots & \\ 1 & & & & \\ \hline 0 & & & & \\ \vdots & & & & \\ 0 & \mathbf{0}_{I \times I-1} & \dots & \mathbf{0}_{I \times I-1} & \mathcal{J}_{\mathcal{P}} \\ 1 & & & & \end{bmatrix},$$

Figure 5.5: The matrix \mathbf{Q}_1 . The matrix $\mathcal{J}_{\mathcal{P}} \in \mathbb{R}^{I \times (I-1)}$ is the Jacobian matrix of the mapping for truncated factors (see Equation (5.9)).

show, these techniques allow for an improvement of the identifiability bounds.

Proposition 5.2.5. *If the non-negative CP model of size*

$$(M_1(I-1)) \times (M_2(I-1)) \times (M_3(I-1))$$

is identifiable and rank R is strictly below the generic rank, then the model (5.22) is also identifiable.

To prove Proposition 5.2.5, we use a lemma about properties of the constraints in (5.22).

Lemma 5.2.6. *The set $\blacktriangle_{R,I}^{(k)}$ (for $k = 1, 2, 3$) contains an open Euclidean subset of $\mathbb{R}^{R(M_k(I-1)+1)}$.*

Proof. First, we prove that there exists a basis of $\mathbb{R}^{M_k(I-1)+1}$ spanned by elements of $\blacktriangle_{1,I}^{(k)}$. Let

$\mathbf{E} \in \mathbb{R}^{(M_1(I-1)+1) \times (M_1(I-1)+1)}$ the following square matrix:

$$\mathbf{E} = \begin{bmatrix} 1 & 1 & \dots & \dots & 1 \\ 0 & 0.5 & 0 & \dots & 0 \\ \vdots & 0 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 0 & 0 & \dots & \dots & 0.5 \end{bmatrix} = \begin{bmatrix} \mathbf{e}_1 & \dots & \mathbf{e}_{M_1(I-1)+1} \end{bmatrix}.$$

Each column \mathbf{e}_i follows the constraints of $\blacktriangle_{1,I}^{(k)}$. Moreover, this matrix is full rank hence represents a base for the set $\mathbb{R}^{M_k(I-1)+1}$.

Next, let us define the space \mathcal{W} the set of linear and positive combinations of columns of \mathbf{E} :

$$\mathcal{W} = \left\{ \epsilon_1 \mathbf{e}_1 + \dots + \epsilon_{M_1(I-1)+1} \mathbf{e}_{M_1(I-1)+1} \mid \epsilon_1 > 0, \dots, \epsilon_{M_1(I-1)+1} > 0 \right\}.$$

First, it is obvious that $\mathcal{W} \subset \blacktriangle_{1,I}^{(k)}$. For a $\mathbf{w} = \epsilon_1 \mathbf{e}_1 + \dots + \epsilon_{M_1(I-1)+1} \mathbf{e}_{M_1(I-1)+1} \in \mathcal{W}$ can be written:

$$w = \begin{bmatrix} \sum_{i=1}^{M_1(I-1)+1} \epsilon_i \\ 0.5\epsilon_2 \\ \vdots \\ 0.5\epsilon_{M_1(I-1)+1} \end{bmatrix}.$$

Therefore, for any $j \in \mathcal{M}_1$, because the combination is strictly positive, we have:

$$\sum_{i=1}^{M_1(I-1)+1} \epsilon_i > \sum_{i=1}^{I-1} \epsilon_{(j-1)(I-1)+1+i}.$$

This means that $\mathcal{W} \subset \blacktriangle_{1,I}^{(k)}$, and \mathcal{W} which is open by construction. Finally, we note that $\blacktriangle_{R,I}^{(k)} = (\blacktriangle_{1,I}^{(k)})^R \subset \mathbb{R}^{R(M_k(I-1)+1)}$ is also an open subset, which completes the proof. \square

Before proving Proposition 5.2.5, let us present a simple example where $M_1 = 1$ and $I = 3$ and $R = 1$. A factor matrix $\mathbf{D}^{(1)} \in \mathcal{D}_{\blacktriangle}^{(1)}$ is a $(M_1(I-1)+1) \times R = 3 \times R$ matrix such that a factor can be written $\mathbf{d}_r^{(1)} = \begin{bmatrix} \alpha_r & d_{r,1}^{(1)} & d_{r,2}^{(1)} \end{bmatrix}^\top$. Each factor is constrained by the following constraints:

$$\alpha_r \geq 0, \quad d_{r,1}^{(1)} \geq 0, \quad d_{r,2}^{(1)} \geq 0 \quad \text{and} \quad d_{r,1}^{(1)} + d_{r,2}^{(1)} \leq \alpha_r.$$

Those constraints are plotted in Figure 5.6 and shows that they form a polyhedron. This is due to the fact that the constraints are linear which creates a set of hyper-planes that defines the polyhedron.

Proof of Proposition 5.2.5. The proof of Proposition 5.2.5 follows from the proof of [83, Theorem 24]. Recall that [83, Theorem 24] considers identifiability of the additive polynomial map Σ_R that maps the factor matrices $\mathbf{D}^{(1)}$, $\mathbf{D}^{(2)}$ and $\mathbf{D}^{(3)}$ to the CPD with these factors $[\mathbf{D}^{(1)}, \mathbf{D}^{(2)}, \mathbf{D}^{(3)}]$. [83, Theorem 24] states that the following implication holds true (for sub-generic R):

$$\begin{aligned} & (\Sigma_R, \mathbb{R}^{M_1(I-1)+1} \times \mathbb{R}^{M_2(I-1)+1} \times \mathbb{R}^{M_3(I-1)+1}) \text{ identifiable} \\ \Rightarrow & (\Sigma_R, \mathbb{R}_+^{M_1(I-1)+1} \times \mathbb{R}_+^{M_2(I-1)+1} \times \mathbb{R}_+^{M_3(I-1)+1}) \text{ identifiable;} \end{aligned} \quad (5.25)$$

Note that the set of constraints in (5.25) by $\blacktriangle_{R,I}^{(1)} \times \blacktriangle_{R,I}^{(2)} \times \blacktriangle_{R,I}^{(3)}$ is a (semi-algebraic) subset of $\mathbb{R}_+^{M_1(I-1)+1} \times \mathbb{R}_+^{M_2(I-1)+1} \times \mathbb{R}_+^{M_3(I-1)+1}$, which contains an open subset by Lemma 5.2.6 (and thus has positive measure). Therefore, the conditions of Lemma 2.3.10 are satisfied, and therefore the following model

$$(\Sigma_R, \blacktriangle_{R,I}^{(1)} \times \blacktriangle_{R,I}^{(2)} \times \blacktriangle_{R,I}^{(3)}), \quad (5.26)$$

which is exactly the model (5.22), is identifiable as well. □

Finally, we show that Theorem 5.2.3 follows from the previously proved results.

Proof of Theorem 5.2.3. Theorem 5.2.3 follows from the equivalence of models in Lemma 5.2.4

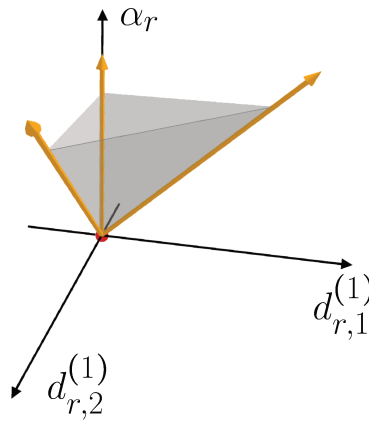


Figure 5.6: Visualization of the constraints which form a polyhedron in a 3D space.

and by invoking Proposition 5.2.5. □

5.3 Application of Theorem 5.2.3

5.3.1 Identifiability result in the general case of variable partition coupling

Thanks to Theorem 5.2.3, for the coupling $\mathcal{T} = \mathcal{M}_1 \times \mathcal{M}_2 \times \mathcal{M}_3$ it is possible to apply the generic uniqueness result from Proposition 2.2.12.

Corollary 5.3.1. *Let $M > 3$ and a partition such that $M_1 \geq M_2 \geq M_3$. The model (5.1) is identifiable from the marginals $\{\mathcal{H}^{(j k \ell)}\}_{\{j, k, \ell\} \in \mathcal{T}}$ if $R < R_g$, where*

$$R_g = \left\lceil \frac{1 + M(I - 1) + (M_1 M_2 + M_1 M_3 + M_2 M_3)(I - 1)^2 + M_1 M_2 M_3 (I - 1)^3}{M(I - 1) + 1} \right\rceil$$

. This result is applicable for every set of parameters I , M_1 , M_2 and M_3 , except for the following cases:

1. if $(M_1(I - 1) + 1)(M_2(I - 1) + 1)(M_3(I - 1) + 1) > 15000$,
2. if $M_1 = M_2 = M_3 = 3$ and $I = 2$,
3. if $M_1 = M_2 = 3$, $M_3 = 2$ and $I = 2$,
4. if $M_1 = M_2 = 5$, $M_3 = 2$ and $I = 2$,
5. if $M_1 > M_2 M_3 (I - 1)$.

The first exception is justified because the proof of this theorem [27] was performed symbolically for a limited number of parameters. The exception 2 is obtained with the case (presented in [83, Corollary 5.13]) of an order-3 tensor of size $4 \times 4 \times 4$ which corresponds to the partition mentioned in the Corollary. Same argument for exception 3 but for order-3 tensors of size $4 \times 4 \times 3$. Same argument for exception 4 but for order-3 tensors of size $6 \times 6 \times 3$. The last exception relies on the fact that the theorem does not hold if the studied tensor has a size too big compared to the other two sizes. This is defined by the condition:

$$\begin{aligned} M_1(I - 1) + 1 &> (M_2(I - 1) + 1)(M_3(I - 1) + 1) - (M_2(I - 1) + 1) - (M_3(I - 1) + 1) \\ \implies M_1 &> M_2 M_3 (I - 1). \end{aligned}$$

5.3.2 Identifiability result for an even partition of variables

In this section, we consider the partition of variables such that $M = 3K + \varepsilon$ where $\varepsilon < 3$:

- If $M = 3K$, $M_1 = M_2 = M_3 = K$,
- If $M = 3K + 1$, $M_1 = K + 1$ and $M_2 = M_3 = K$,
- If $M = 3K + 2$, $M_1 = M_2 = K + 1$ and $M_3 = K$.

Note that this partition groups variables in the most balanced way possible. With the same reasoning of Corollary 5.3.1, we obtain the following result:

Corollary 5.3.2. *Let $M > 3$ and a partition previously presented and the coupling $\mathcal{T} = \mathcal{M}_1 \times \mathcal{M}_2 \times \mathcal{M}_3$. The model (5.1) is identifiable from the marginals $\{\mathcal{H}^{(j k \ell)}\}_{\{j, k, \ell\} \in \mathcal{T}}$ if $R < R_g$. This result is applicable for every set of parameters I and M , except for the following cases:*

1. if $(M_1(I - 1) + 1)(M_2(I - 1) + 1)(M_3(I - 1) + 1) > 15000$,
2. if $M = 9$ and $I = 2$,
3. if $M = 4$ and $I = 2$,
4. if $M = 8$ and $I = 2$.

This condition is the best identifiability result possible. This is due to the fact that partitioning variables evenly permits to maximize the value of the generic rank R_g . Note that because this case evens the number of variables in the partition, the case of an unbalanced tensor only occurs in the case $M = 4$ and $I = 2$. Indeed, let us consider all possible cases:

If $M = 3K$

For this case, $M > 3$ so $K > 1$ and $M_1 = M_2 = M_3 = K$. The first possible exceptional case is the case of $4 \times 4 \times 4$ tensors. This case occurs if $K = 3$ ($M = 9$) and $I = 2$:

$$M_1(I - 1) + 1 = M_2(I - 1) + 1 = M_3(I - 1) + 1 = K + 1 = 4.$$

The unbalanced case is not possible in this configuration. Indeed, $K > K^2(I - 1)$ would imply that $K < 1$ which is not possible.

If $M = 3K + 1$

For this case, $M > 3$ so $K \geq 1$ and $M_1 = K + 1$ and $M_2 = M_3 = K$. The only exceptional case is the case of unbalanced tensors:

$$K + 1 > K^2(I - 1) \xRightarrow{I-1 \geq 1} K + 1 > K^2 \implies K = 1.$$

This is the case that gives the exception $M = 4$ and $I = 2$.

If $M = 3K + 2$

For this case, $M > 3$ so $K \geq 1$ and $M_1 = M_2 = K + 1$ and $M_3 = K$. The first possible exception case is the case of $4 \times 4 \times 3$ tensors. This case occurs if $K = 2$ ($M = 8$) and $I = 2$.

The unbalanced case is not possible in this configuration. Indeed, $K + 1 > (K + 1)K(I - 1)$ would imply that $K < 1$ which is not possible.

5.3.3 Identifiability result for the fully coupled case

As mentioned in Section 5.2, the identifiability of the fully coupled case is guaranteed by the identifiability of any partial model. Therefore, by choosing the model of Section 5.3.2, we obtain the most favorable identifiability result:

Corollary 5.3.3. *The model (5.1) is identifiable from its marginals $\{\mathcal{H}^{(j k \ell)}\}_{\{j, k, \ell\} \in \mathcal{T}_{all}}$ if $R < R_g$ (with an even partition). This result is applicable for any M and I except for the cases presented in Corollary 5.3.2*

Because this corollary does not take all triplets into account, we believe that this bound can be improved (see Conjecture 5.1.4). With those results, Figure 5.7 plots the same Figure 5.1 but adds our identifiability results. This figure shows that the identifiability bound of this corollary is better than the bound of the original CTF3D paper [74].

5.4 Conclusion of Chapter 5

In this chapter, the recoverability and the identifiability of our coupled model were studied. Concerning the recoverability, an algorithm was proposed to find the maximal rank that ensures

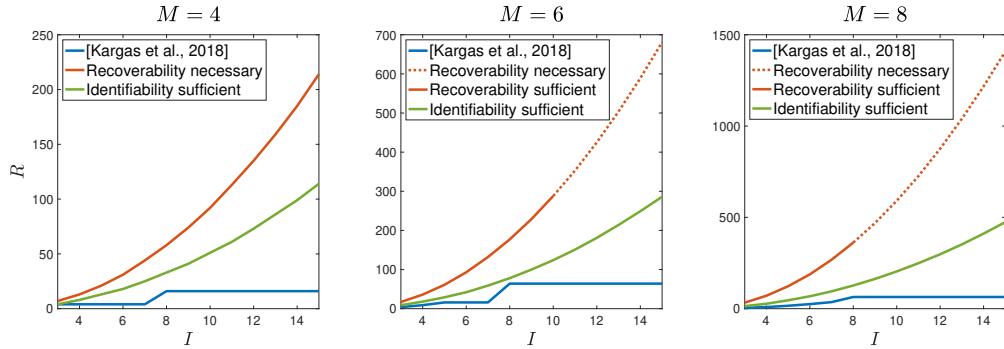


Figure 5.7: Identifiability bounds in the fully coupled case along with recoverability results.

recoverability hence giving us a sufficient recoverability condition in specific cases. This algorithm is based on the computation of the rank of the Jacobian matrix of our parametrization by verifying that it is full column rank.

This algorithm was applied to the case of both fully coupled tensor factorization and partially coupled tensor factorization. In both cases, necessary conditions were proposed. In the case of full coupling, our computations suggest that the recoverability holds till the necessary condition, which we state as a conjecture. For partial coupling strategies, a study showed that a non-negligible part of random couplings have low recoverability bound. Some of those so-called “defective cases” were studied and the reason for this drop in recoverability was explained by the structure of Jacobian which depends a lot on the coupling. On the other hand, balanced strategies provided in practice a good alternative to random couplings as balanced couplings did not feature defective cases.

Concerning the problem of identifiability, an analysis of the proofs of theorems of [74] showed that the sum-to-one constraints on factors may have been overlooked. In the second part of this chapter, a proof of identifiability based on identifiability of polynomial models is proposed to handle the simplex constraints properly. First, marginals can be stored in an order-3 tensors whose CP model identifiability can be studied. Second, the set of constraints can be reparametrized to use the identifiability for a non-negative CPD of a smaller tensor. The bounds obtained with our results improve over existing results, but are lower than the necessary bound on recoverability. Hence it may be possible to prove stronger identifiability results.

Chapter 6

CTFlowHD : a tensor-based method for flow cytometry data analysis

Contents

6.1	Presentation of the method <i>CTFlowHD</i>	124
6.1.1	Naive Bayes model and flow cytometry data	124
6.1.2	<i>CTFlowHD</i> workflow of analysis	125
6.1.3	Implementation	127
6.1.4	Application to a controlled flow cytometry dataset	128
6.2	Clustering and visualizations	131
6.2.1	Visualizations based on hierarchical clustering	132
6.2.2	Marginal visualizations	134
6.2.3	<i>CTFlowHD</i> : a versatile visualization tool	136
6.3	Application to real flow cytometry data	140
6.3.1	Application to 8-marker datasets	141
6.3.2	Application to 24-marker datasets	146
6.4	Conclusion of Chapter 6	148

This chapter aims at presenting the computational tools developed for flow cytometry (FCM) data analysis. The tools are gathered in a package called *CTFlowHD*. The targeted readers are end-users which are not much interested in the theoretical development of the algorithm PCTF3D (see Chapter 4) but rather by its practical use for analysing FCM data. Thus, it is written as a guide aiming at illustrating the possible uses of the package.

The core method of the package *CTFlowHD* is the algorithm PCTF3D (see Chapter 4) which estimates the factors of a naive Bayes model (NBM) and which is applied to pre-processed data. It outputs the factors of the NBM (fluorescence properties and associated proportions).

The NBM factors need to be post-processed to provide a clear biological/medical interpretation. Thus, the package also includes clustering and visualization tools acting directly on NBM factors. By isolating PCTF3D and visualization, *CTFlowHD* then becomes a versatile tool for flow cytometry data analysis. In particular, it is showed how to apply several tools with low additional computational burden, including computational methods widely used in the flow cytometry community.

The use of the *CTFlowHD* will be illustrated on a "simple to analyze" 4-marker controlled data set for which the ground truth is available. More challenging situations will also be considered: a 8-marker dataset with partial biological/medical information available and a dataset featuring 24 markers for which no interpretations are available.

6.1 Presentation of the method *CTFlowHD*

6.1.1 Naive Bayes model and flow cytometry data

For M fluorescence biomarkers, our method takes as an input an observation matrix \mathbf{X} . This matrix is of size $N \times M$ and contains for each of the N cells the M preprocessed values of fluorescence. Note that the geometric variables (FSC,SSC) are typically excluded from the analysis of *CTFlowHD*.

If fluorescence values are considered random variables $\{X^{(m)}\}_{m=1}^M$, it is possible to apply a naive Bayes model (NBM) to model the multidimensional density:

$$p(X^{(1)}, \dots, X^{(M)}) = \sum_{r=1}^R \Pr(L = r) \prod_{m=1}^M p(X^{(m)} | L = r).$$

When this model is applied to a FCM dataset, the latent variable L can be interpreted as a set of R cell populations. Each cell population is represented by its loading vector value $\Pr(L =$

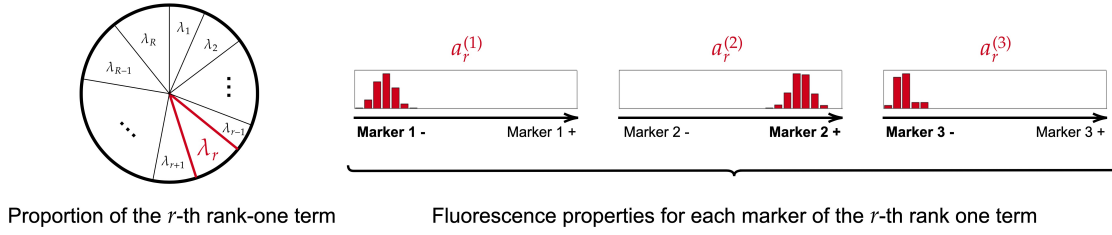


Figure 6.1: Naive Bayes model for flow cytometry data analysis.

r). This value can be interpreted as the percentage of the r -th cell population in the pool of cells. Concerning the fluorescence values, a cell population is characterized by M 1-dimensional profiles of fluorescence. Therefore, applying the NBM to a FCM dataset permits to separate cells into populations and characterize those populations in terms of fluorescence.

If each variable is discretized, the NBM can be interpreted as a low-rank tensor decomposition (see Chapter 3). In this case, the density is a histogram in M dimensions which can be modeled with a sum of product of 1D histograms as shown in Figure 6.1.

6.1.2 CTFlowHD workflow of analysis

CTFlowHD stands for Coupled Tensor for Flow cytometry in High Dimensions. This method is separated in different modules which are presented in the workflow of Figure 6.2. First, data is preprocessed with the three steps presented in Section 1.2. It results in an observation matrix \mathbf{X} which contains N rows (one for each cell) and M columns (one for each biomarker).

After this step, the algorithm PCTF3D can be performed. This algorithm has three parameters that need to be chosen by end users. The first parameter is the number of bins per dimension denoted I . Each fluorescence space will be uniformly discretized into I bins. This value can be chosen according to the error of marginal estimation (see Equation (3.10)) or the complexity of PCTF3D (see Section 4.1).

The second argument of PCTF3D is the choice of triplets of variables. Indeed, PCTF3D couples 3D histograms to estimate the NBM of a FCM dataset. For a high number of biomarkers, not all 3D histograms can be computed and coupled (see Section 3.4.3). To choose the set of triplets \mathcal{T} , coupling strategies have been presented in Section 4.2. Triplets can be either chosen randomly or fixed manually by end users. Most coupling strategies permit to choose the number of triplets T considered in the method thus end users may adapt the level of complexity to their application.

As soon as \mathcal{T} and I are set, order-3 histograms can be computed from the observation matrix

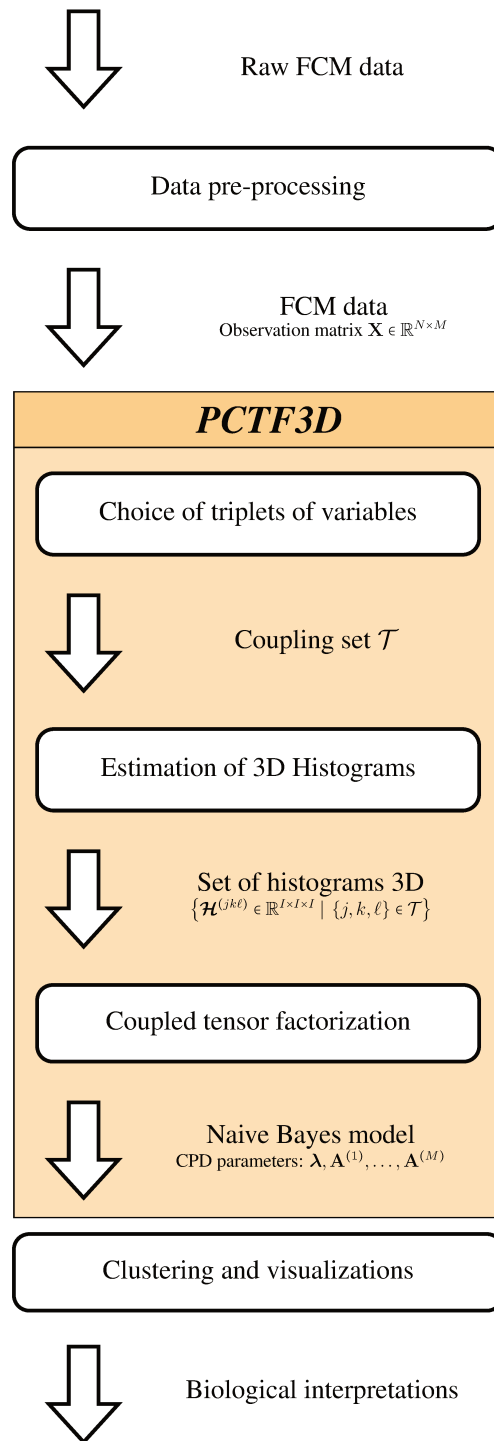


Figure 6.2: CTFLOWHD workflow that contains the algorithm PCTF3D.

X. The last parameter needed to perform PCTF3D is the rank of decomposition. This parameter represents the number of cell populations featured in the NBM. This number must remain below both identifiability and recoverability bounds presented in Chapter 5.

With the NBM of a FCM dataset, it may be difficult to extract biological interpretations from individual rank-1 components. Typically, the rank-1 components should be grouped to obtain biologically interpretable populations. To enhance the interpretability of *CTFlowHD*, clustering and visualization methods can be applied to the NBM components. This step will be thoroughly presented in a following section.

6.1.3 Implementation

CTFlowHD is available on MATLAB[®] (apart from preprocessing steps) at https://github.com/philippflores/fcm_ctflowhd. The package is divided into three categories that are presented in Figure 6.2: data pre-processing, PCTF3D and clustering/visualization.

First, preprocessing step scripts are provided in the repository [fcm_ctflowhd/tree/main/preprocessingR](https://github.com/philippflores/fcm_ctflowhd/tree/main/preprocessingR). One script for each preprocessing step is available as well as a script to apply preprocessing steps all-at-once. These scripts can be run on R and use the packages `flowCore` [41], `flowAI` [73], `lattice` [93], `flowViz` [40], `flowClust` [66] and `flowTrans` [42].

Second, the algorithm PCTF3D is provided in the context of FCM. This algorithm takes three parameters as input: the number of bins per dimension I , the coupling \mathcal{T} and the number of NBM components R . To choose those parameters and apply PCTF3D, Appendix C presents how the algorithm PCTF3D can be applied in the context of FCM. In addition to MATLAB[®] packages, PCTF3D uses the tensor package provided by [5].

Finally, *CTFlowHD* provides a lot of clustering and visualization tools to interpret the results of PCTF3D (see Section 6.1.4). Each visualization method is presented in a separated file whose name begins with “CTFlowHD_plot_” associated with a text file (README) that presents the principles and the different argument of each tool. One script featuring all methods is also available. When a visualization is based on an already existing clustering/visualization method, we provide a script to apply directly the said method for comparison in the repository [fcm_ctflowhd/tree/main/other_methods](https://github.com/philippflores/fcm_ctflowhd/tree/main/other_methods).

Note that the package does not feature any FCM datasets. If one wants to use *CTFlowHD* on the 4D controlled experiment datasets, it may be provided by asking the author.

Table 6.1: Properties of the 3 populations used in the controlled experiment. + is high marker expression and - low expression.

Population	Marker expression			
	CFSE	TCR	CTV	MHC-II
Macrophage	-	-	+	+
Lymphocyte B	+	-	-	++
Lymphocyte T	-	+-	-	-

6.1.4 Application to a controlled flow cytometry dataset

To illustrate the principles of *CTFlowHD*, this section features an example in 4 dimensions. For this first example, a synthetic dataset with no real data difficulties was not an ideal choice. On the other hand, it would have been hard if not impossible to use a FCM dataset used in an applicative environment, as no ground truth would have been available. The data for this experiment was provided by Guillaume Harlé and Stéphanie Grandemange from CRAN.

The controlled dataset features 3 cell populations: Lymphocytes T or T cells (LT), Lymphocytes B or B cells (LB) and Macrophages (MP). T cells were marked via TCR with the fluorescent molecule PE (Phycoerythrin). B cells were marked with CFSE (Carboxyfluorescein succinimidyl ester) combined with the fluorophore FITC (Fluorescein isothiocyanate). Macrophages were marked with the antibody CTV attached to the fluorophore BV-421. The fourth marker is MHC class 2 or MHC-II. This marker is not expressed by T cells while B cells highly express this marker. Macrophages also express MHC-II but not as much as B cells. Table 6.1 summarizes the properties of cell populations mixed in the controlled experiment. LT are noted +- in this table because the TCR marking was not very efficient in this experiment, as shown in the following results.

Ground truth for this study was obtained via manual gating. This is not an issue in this case because the 3 cell populations are well separated in terms of fluorescence properties. This was verified in practice in the bivariate plots of Figure 6.3.

After preprocessing steps, the observation matrix \mathbf{X} with $N = 53881$ cells and $M = 4$ markers was obtained. As shown in Section 3.2.2, the number of bins I is limited by the curse of dimensionality and I can be chosen in accordance with Equation (3.10)

$$I_m^* = \left\lceil \frac{S_m \sqrt[5]{N}}{3.5\sigma_m} \right\rceil.$$

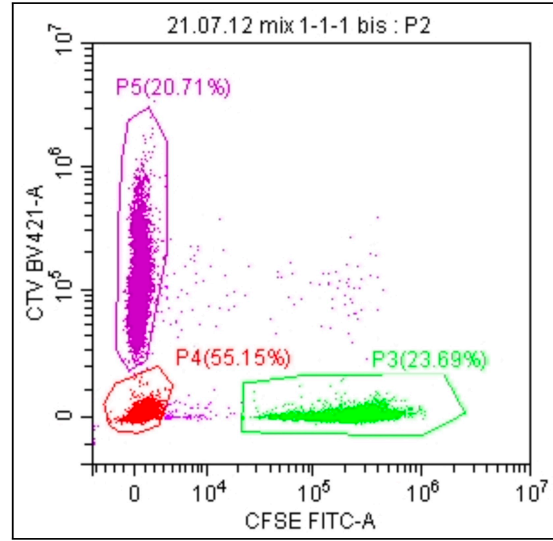


Figure 6.3: Manual gating providing ground truth for the controlled experiment.

The support for each variable is chosen as $S_m = \max_{n \in \llbracket 1, N \rrbracket} (X_{nm}) - \min_{n \in \llbracket 1, N \rrbracket} (X_{nm})$ and the standard deviation σ_m is estimated numerically for each variable. For this particular controlled dataset, the optimal number of bins for each variable is then $\{10, 23, 13, 13\}$. Because the optimal number of bin rule only provides an indication for the choice of I , it has been decided to set I equal to 15 bins in practice.

In this case with $M = 4$ markers, there could be at most 4 triplets but only $T = 2$ triplets were picked which resulted the following coupling:

$$\mathcal{T} = \{\{1, 3, 4\}, \{2, 3, 4\}\}.$$

Therefore, the lower-order histograms $\mathcal{H}^{(134)}$ and $\mathcal{H}^{(234)}$ are estimated from \mathbf{X} (see Equation (3.23)).

The NBM of this FCM dataset is obtained by applying PCTF3D. For this case, the necessary condition on recoverability is:

$$R_{\max} = \left\lceil \frac{1 + 4(20 - 1) + 5(20 - 1)^2 + 2(20 - 1)^3}{4(20 - 1) + 1} \right\rceil = 202.$$

Moreover, it is possible to apply the Algorithm 5 to verify if this bound is achieved in the generic case. This algorithm applied for this controlled experiment shows that for $R = 50$, the model is generically recoverable.

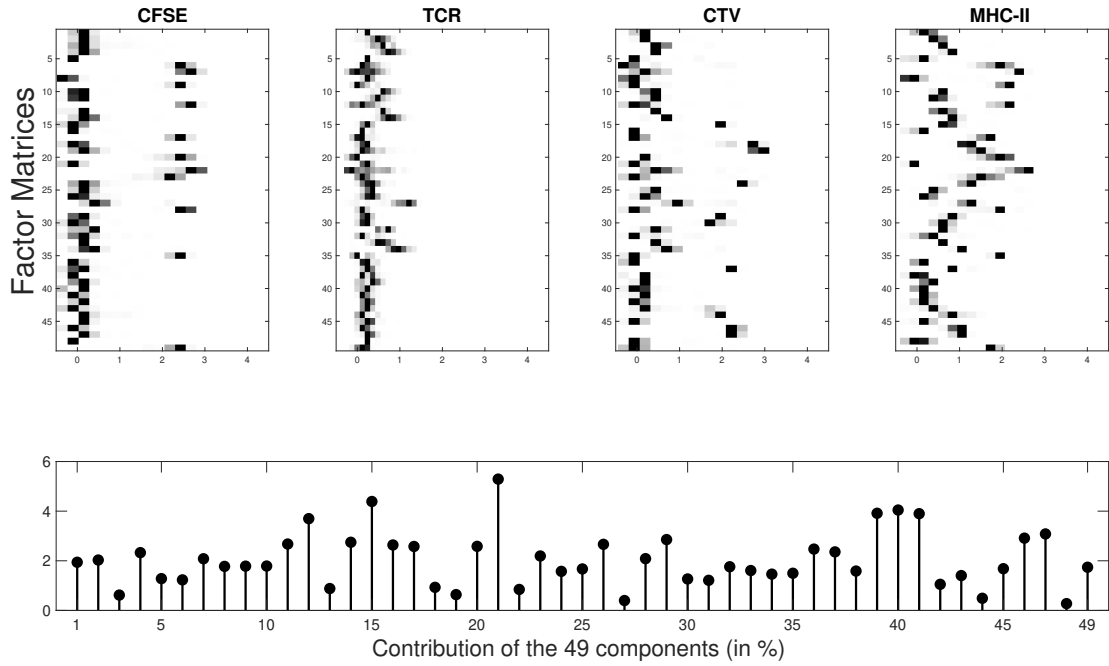


Figure 6.4: *CTFlowHD* output for a controlled experiment before the visualization step.

When PCTF3D is applied, the rank $R = 50$ decomposition featured a component with λ_r equal to 0. Therefore, after the deletion of this rank-one term, it defined $R = 49$ as the new rank of the decomposition. Because those components have a loading value of 0, they do not affect the sum of rank-one terms.

The estimated NBM of this dataset is plotted in Figure 6.4. On this figure, the bottom plot is a stem representation of the loading vector λ . Therefore, each value of this plot represents the proportion of cells in \mathbf{X} that has the fluorescence properties presented in the upper plots. Indeed, the $M = 4$ factor matrices are plotted such that each row corresponds to a transposed factor $\mathbf{a}_r^{(m)\top}$. The more a factor is located to the right part of a marker plot, the more this component expresses this marker. To summarize, an end user can characterize each cell population obtained by the NBM of *CTFlowHD* by examining its proportion and its M fluorescence properties.

However, this raw visualization alone is limited. Indeed, a lot of rank-one components obtained in Figure 6.4 are similar but their order is random. This is due to the fact that a low-rank tensor decomposition is defined up to a permutation of rank-one terms. Figure 6.5 shows a visualization that permits a better overall interpretation. This clustering/visualization step will be studied in details in the following section.

With this visualization step, it is now possible to extract the three cell populations that were

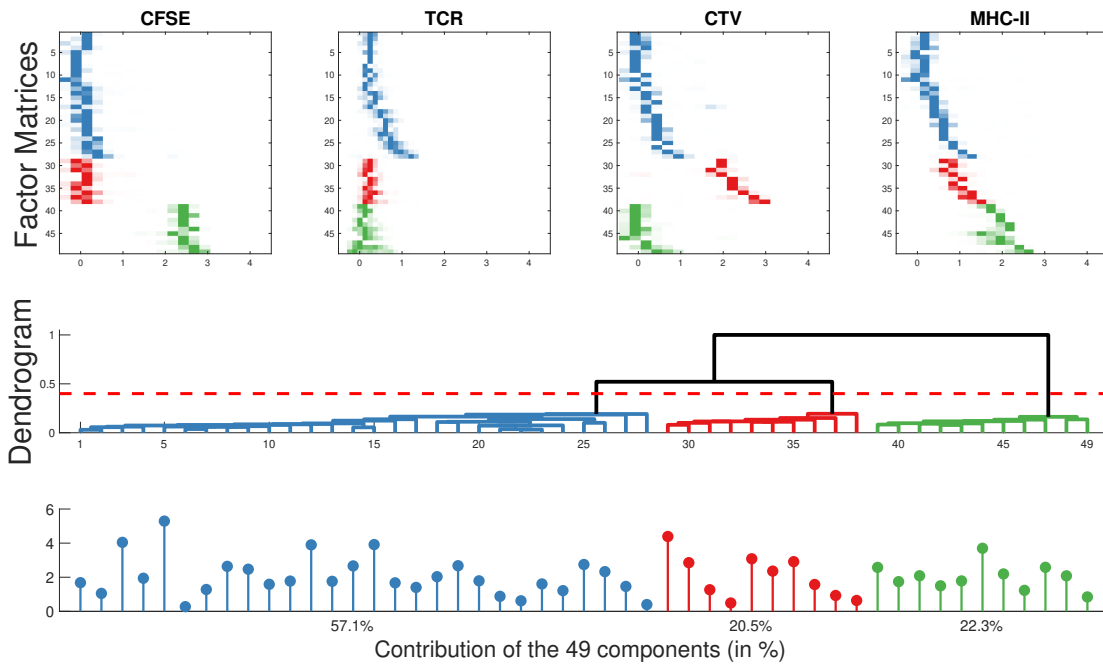


Figure 6.5: *CTFlowHD* output for a controlled experiment after a hierarchical clustering applied to rank-one terms (single linkage).

mixed at the start of the experiment. Indeed, the **blue** population expresses the most the TCR marker hence corresponds to T cells. As mentioned before, a better marking for this cell population was expected. Nonetheless, the proportion of T cells obtained in Figure 6.3 was 55.15% which is similar to 57.3% obtained with *CTFlowHD*. This value was obtained by summing the loading vector values for the rank-one terms of the **blue** cluster.

Concerning the other two cell populations, **red** cells are CTV+ hence corresponds to Macrophages while **green** cells are CFSE+ and corresponds to B cells (see Table 6.1).

6.2 Clustering and visualizations

The controlled experiment of Section 6.1.4 highlighted the fact that the NBM components obtained with PCTF3D is not sufficient to interpret a FCM dataset. In this section, clustering and visualization tools that can be applied to tensor decomposition factors are presented.

6.2.1 Visualizations based on hierarchical clustering

Hierarchical clustering is a clustering method based on a dendrogram. To build a dendrogram tree diagram from a base of R points, the agglomerative method (in opposition to the divisive method) consists in merging branches until there remains only one ending branch containing all R points. Merges are done by comparing branches with a particular distance. An example is given in Figure 6.5.

To apply hierarchical clustering to rank-one components, a distance between rank-one terms must be defined. First, a vector of fluorescence centroids $\mathbf{f}_r \in \mathbb{R}^M$ is computed for each rank-one term:

$$\mathbf{f}_r := \begin{bmatrix} \sum_{i_1=1}^I \overline{\Delta}_{i_1}^{(1)} \widehat{a}_{i_1 r}^{(1)} \\ \vdots \\ \sum_{i_M=1}^I \overline{\Delta}_{i_M}^{(M)} \widehat{a}_{i_M r}^{(M)} \end{bmatrix},$$

where $\overline{\Delta}_{i_m}^{(m)}$ represents the centroid of the i_m -th bin $\Delta_{i_m}^{(m)}$. Then, it is possible to compare two rank-one components by evaluating the Euclidean distance between the fluorescence centroids:

$$d(r, s) := \|\mathbf{f}_r - \mathbf{f}_s\|_2^2. \quad (6.1)$$

The principle of hierarchical clustering is recursive. The objective of each recursive step is to merge the two closest branches. At first, there are R rank-one branches, one for each component. The two closest regarding Equation (6.1) are grouped into a new branch. This process is repeated until there is a unique branch containing all rank-one terms.

To compare two branches that may contain several rank-one terms (necessary different components), a distance must be defined to find the pair of branches that are the closest. This process of decision is called the linkage. There exist *single linkage* and *complete linkage* which permit to define a distance between branches $\mathcal{B}_1, \mathcal{B}_2 \subset \llbracket 1, R \rrbracket$ such that $\mathcal{B}_1 \cap \mathcal{B}_2 = \emptyset$:

- Single linkage:

$$D(\mathcal{B}_1, \mathcal{B}_2) := \min \{d(r, s) \mid r \in \mathcal{B}_1, s \in \mathcal{B}_2\},$$

- Complete linkage:

$$D(\mathcal{B}_1, \mathcal{B}_2) := \max \{d(r, s) \mid r \in \mathcal{B}_1, s \in \mathcal{B}_2\}.$$

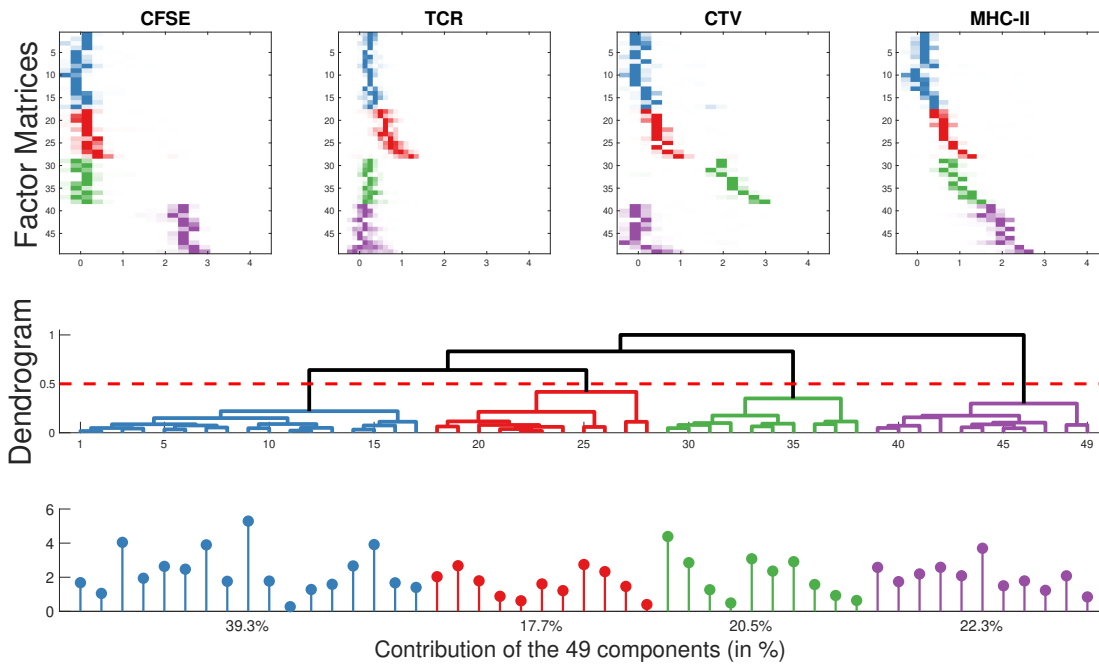


Figure 6.6: *CTFlowHD* output for a controlled experiment after a hierarchical clustering applied to rank-one terms (complete linkage).

Finally, the two branches that have the closest distance are merged to create a new class. A new step of hierarchical clustering will be applied to this new class with the other branches.

Single linkage was applied in Figure 6.5 while complete linkage was applied to the same NBM factors in Figure 6.6. This figure shows that complete linkage creates more classes compared to single linkage. This is due to the fact that with a max distance, classes contain packs of rank-one terms while single linkage creates chains of rank-one terms because of the min distance. Nevertheless, it is still possible to obtain the 3 cell populations in both cases and with the same proportion values.

To create clusters of rank-one components, a linkage tree can be cut at a threshold. The number of branches above the threshold denotes the number of clusters. In Figure 6.6, there are 4 clusters. Even though the controlled experiment featured three types of cells, choosing a threshold such that the 4 clusters can differentiate between T cells that did not express TCR and T cells that express this marker. This means that the threshold must be chosen properly by end users, as this choice can change the biological interpretation.

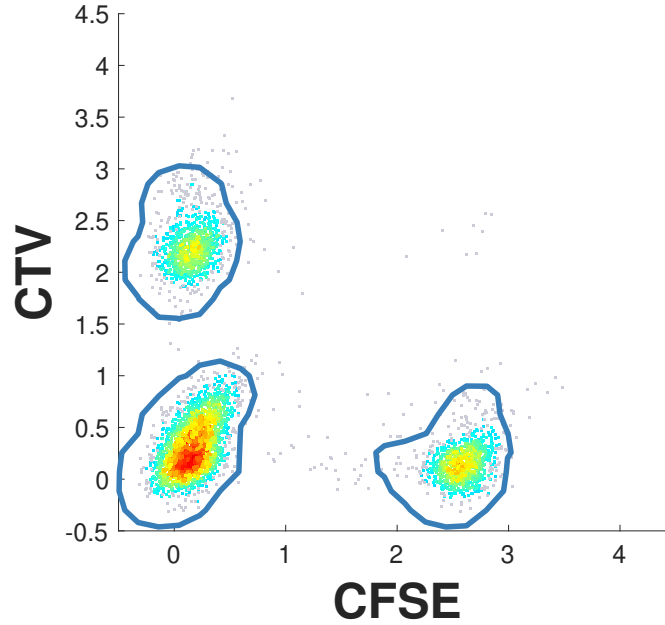


Figure 6.7: Bivariate plot CFSE/CTV for the controlled experiment combined with *CTFlowHD* 2D marginals.

6.2.2 Marginal visualizations

The principle of the naive Bayes model is to represent the order- M distribution by a set of 1-dimensional histograms. This model permits to circumvent the curse of dimensionality as the NBM can be obtained with a reasonable amount of FCM data. However, this same curse of dimensionality prohibits end users to compute the full histogram in M dimensions from the NBM factors. Indeed, for large number of biomarkers, I^M values are difficult to store and visualize.

However, it has been shown in the previous section that 1D-marginals can be plotted altogether to visualize fluorescence properties all-at-once. It is also possible to plot higher-order marginals, for example 2D histograms:

$$\hat{\mathbf{H}}^{(jk)} = \sum_{r=1}^R \lambda_r \hat{\mathbf{a}}_r^{(j)} \circ \hat{\mathbf{a}}_r^{(k)}. \quad (6.2)$$

Figure 6.7 plots an estimated 2D marginals in the case of the controlled experiment of Section 6.1.4. This bivariate plot is similar to the ground truth bivariate plot of Figure 6.3. Indeed, both figures shows the repartition of cells for markers CFSE and CTV. For the gating figure,

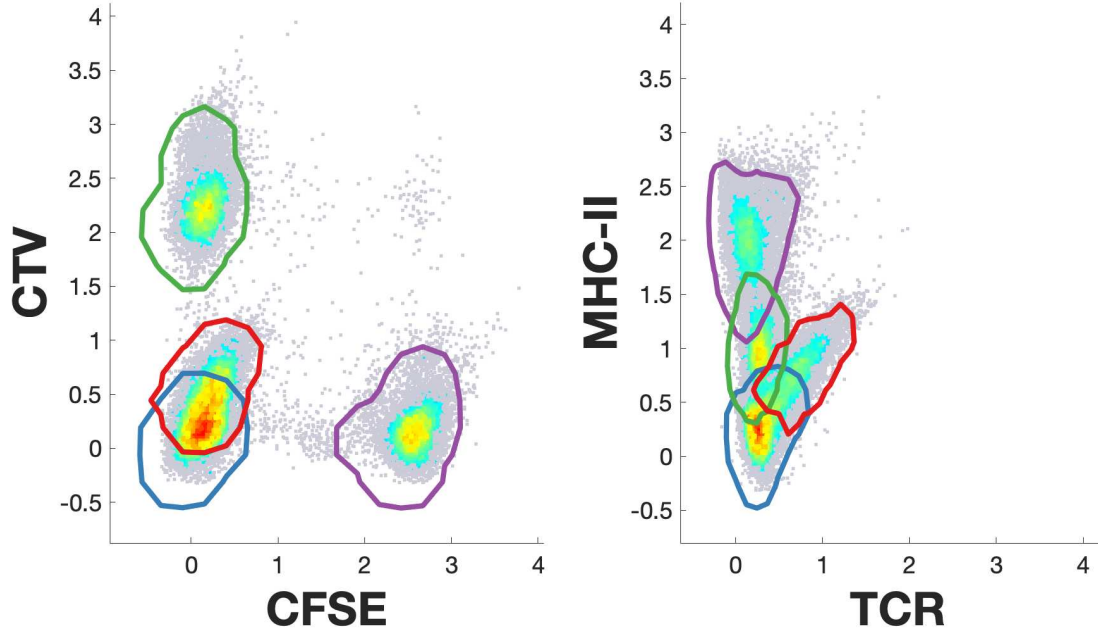


Figure 6.8: *CTFlowHD* 2D marginals combined with hierarchical clustering (complete linkage).

the gates have been selected manually by an end user. Whereas for *CTFlowHD*, the contoured region represents the estimated marginal obtained with Equation (6.2).

In addition to bivariate plots, the previously presented hierarchical clustering is combined with this visualization. It is then possible to plot lower order histograms for a subset of rank-one terms. If \mathcal{B} denotes a subset of rank-1 terms obtained by hierarchical clustering, the corresponding histograms can be defined as:

$$\mathbf{h}_{\mathcal{B}}^{(m)} := \sum_{r \in \mathcal{B}} \lambda_r \widehat{\mathbf{a}}_r^{(m)},$$

$$\mathbf{H}_{\mathcal{B}}^{(jk)} := \sum_{r \in \mathcal{B}} \lambda_r \widehat{\mathbf{a}}_r^{(j)} \circ \widehat{\mathbf{a}}_r^{(k)}.$$

For the controlled experiment, this tool can enhance biological interpretations. First, Figure 6.8 uses the complete linkage clustering of Figure 6.6 with the same colors for each cell population. These plots permit to visualize how cell populations are distributed among others. To visualize a cell population with a gating analysis, other cells must be deleted after each gate. *CTFlowHD* can plot a cell population without eluding other cells.

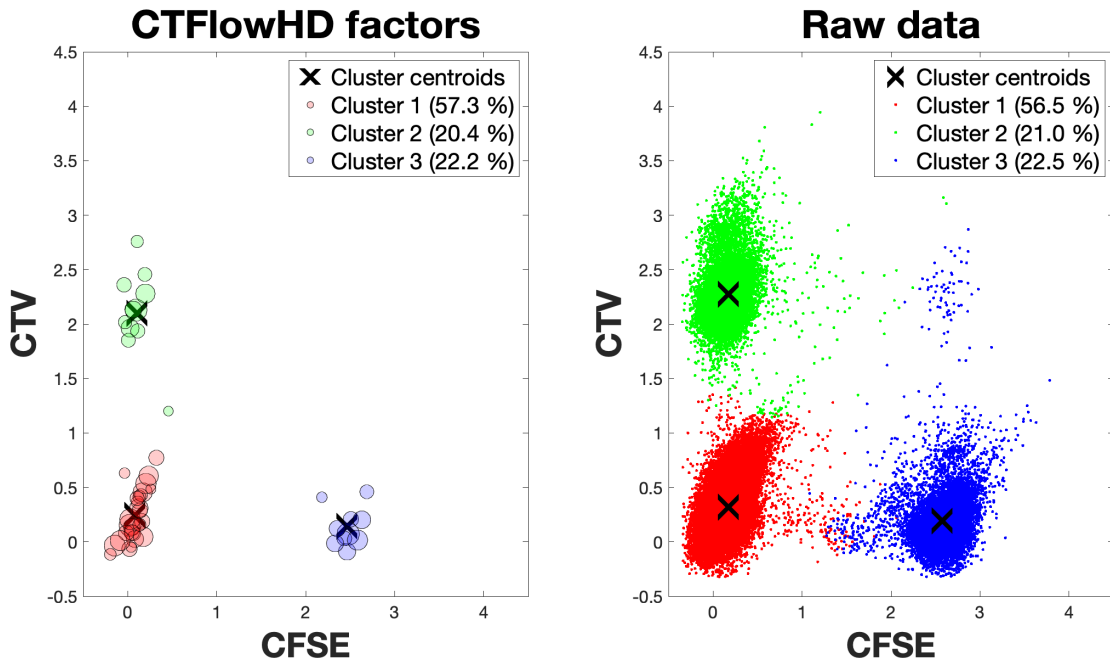


Figure 6.9: Controlled experiment visualized with a K-means clustering method. **Left plot:** K-means is applied to *CTFlowHD* factors. **Right plot:** K-means is applied directly to pre-processed data.

6.2.3 CTFlowHD: a versatile visualization tool

Clustering and visualizing NBM components is computationally very cheap compared to the cost of obtaining the NBM factors. Moreover, as soon as factors are estimated, the clustering step is independent of steps that permit to estimate the NBM factors. Therefore, it is possible to apply various visualization tools at little cost. This includes already existing clustering methods, even methods applied to flow cytometry. To give some examples, a few representative clustering methods were applied to the controlled experiment of Section 6.1.4.

First, K-means [70] was applied with a $K = 3$ classes. Figure 6.9 plots the NBM factors clustered with K-means. The size of each point in this plot represents the proportion of a rank-1 term. The same interpretation than Figures 6.3 and 6.5 was obtained. When the same K-means is applied to FCM data directly, the interpretation remains similar, but it take much more time to cluster raw data (see Table 6.2).

Then, *CTFlowHD* was applied to two of the most used flow cytometry data analysis tools [25]. viSNE [4] is an unsupervised visualization tool that is based on the t-SNE [109]. Figures 6.10 and 6.11 show the results of t-SNE on both raw data and the factors obtained with

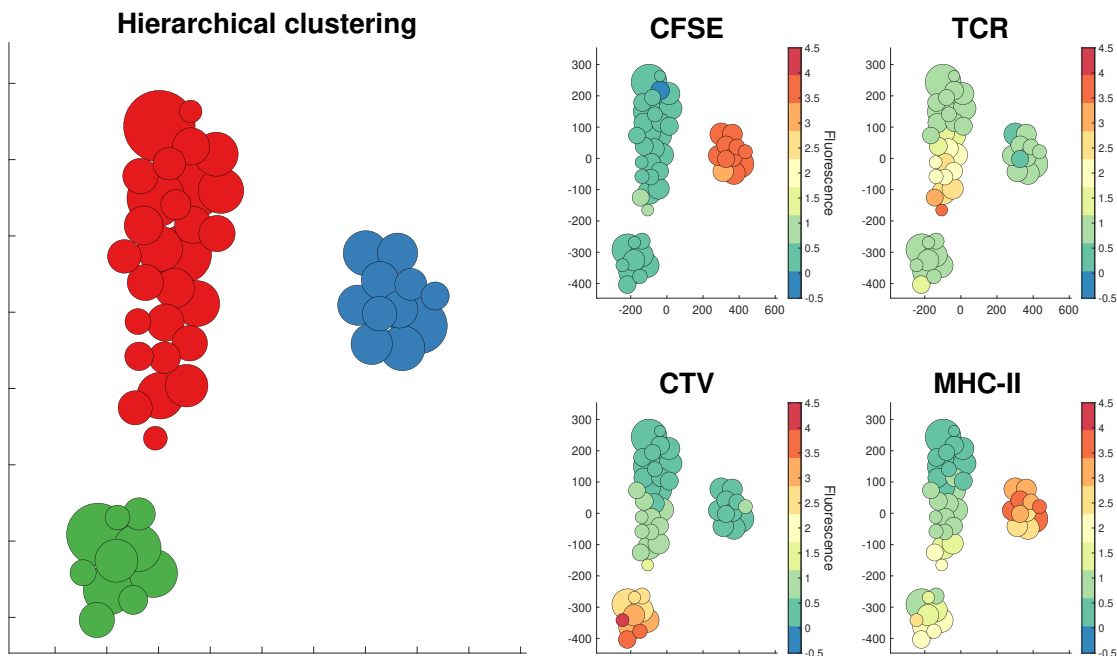


Figure 6.10: *CTFlowHD* applied with a t-SNE visualization (runtime: 2 seconds). **Left plot:** t-SNE visualization combined with the single linkage (see Figure 6.5). **Right plots:** t-SNE maps colored with mean marker expression for each rank-1 term.

CTFlowHD. Similarly to the tool presented in Section 6.2.2, t-SNE can be combined with hierarchical clustering. Figure 6.10 shows that the three populations clustered with single linkage are well separated with t-SNE. As for t-SNE applied directly to FCM data, Figure 6.11 also separates the 3 populations. But the visualizations are less clear than the *CTFlowHD*'s. For example, the TCR+ population (blue cluster in Figure 6.10) is not easily visualized in Figure 6.11.

For our last example, *CTFlowHD* was applied with a visualization step similar to SPADE [84]. The principle of SPADE consists in plotting a minimum spanning tree visualization [49] whose nodes are the results of a K-means clustering. The results for the controlled experiment are shown in Figure 6.12. To apply the same visualization to *CTFlowHD*, the clustering has already been done because NBM factors are estimated. Therefore, the same SPADE minimum spanning tree algorithm permitted to obtain a similar visualization on Figure 6.13.

Table 6.2 shows data analysis runtimes for the controlled experiment. It shows that *CTFlowHD* can be used in combination with many visualization tools. Even if there exists clustering methods that are computationally faster than *CTFlowHD*, our method offers the possibility to obtain visualizations rapidly. Indeed, in most cases, the dimension reduction obtained with

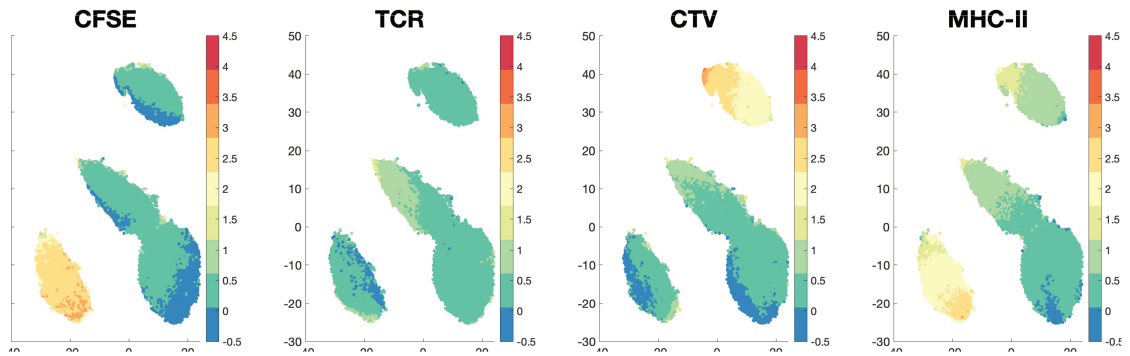


Figure 6.11: t-SNE applied to raw data from the controlled experiment (runtime: 550 seconds). For each marker, a t-SNE map colored with mean marker expression for each cell is plotted.

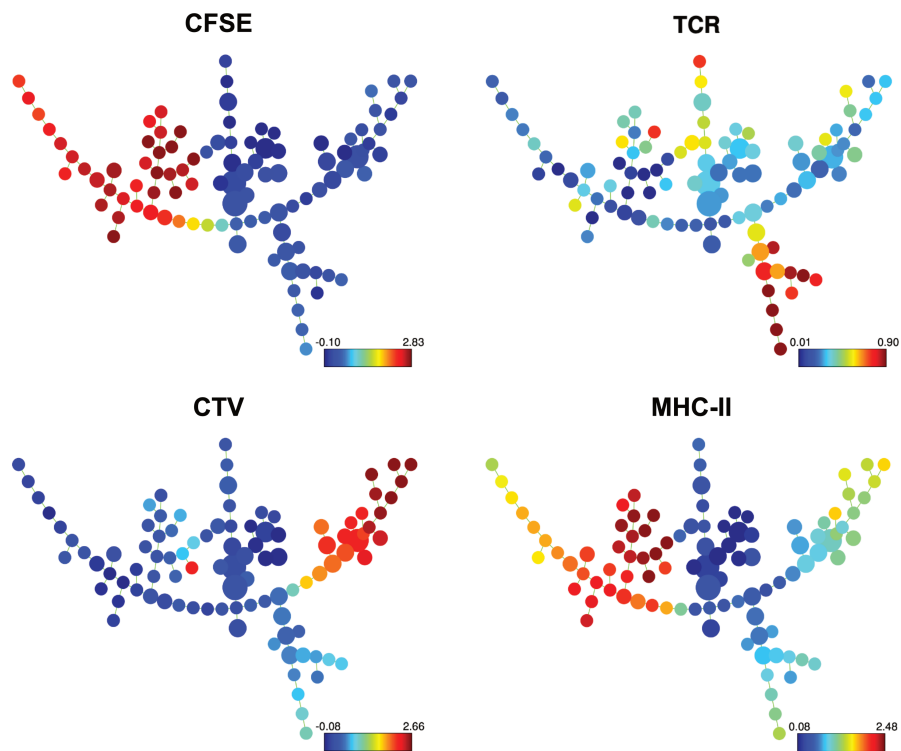


Figure 6.12: SPADE applied to raw data from the controlled experiment (runtime: 12.5 seconds). For each marker, the minimum spanning tree is colored with mean marker expression for each K-means cluster.

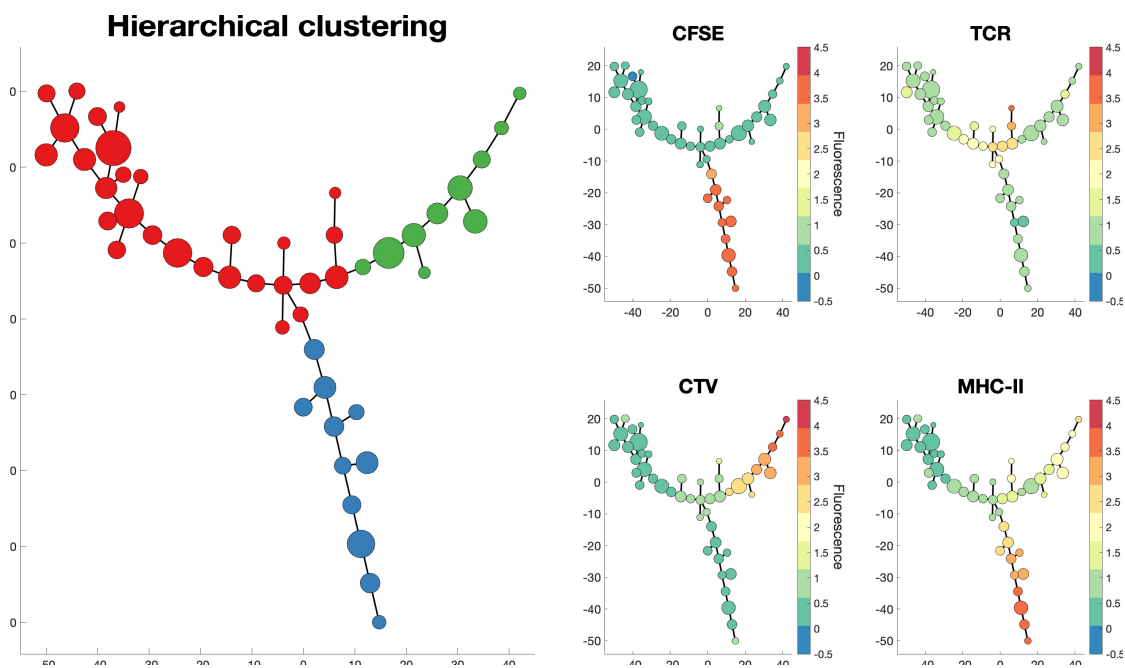


Figure 6.13: *CTFlowHD* applied to the controlled experiment with a SPADE visualization (run-time: 3.5 seconds). **Left plots:** Tree visualization colored with single linkage (see Figure 6.5). **Right plots:** For each marker, the minimum spanning tree is colored with mean marker expression for each K-means cluster.

Table 6.2: Analysis run times for the controlled experiment.

	Runtime	
	FCM pre-processed data	CTFlowHD
PCTF3D	–	13.6s
Linkage	–	2.3s
K-means	0.05s	0.002s
t-SNE	550s	2s
SPADE	12.5s	3.5s
Total	570s	8s

PCTF3D makes the clustering/visualization step of the R rank-one terms computationally free.

6.3 Application to real flow cytometry data

Biological context

As mentioned before, flow cytometry is widely used in immunology. Our study is focused on leukemia which is the most common type of blood cancers. People suffering from leukemia do not develop blood cells normally. This can result in an abnormal proportion of a certain type of cells, usually a cell population in abundance called blast cells [91]. The two main types of leukemia are *Acute Lymphoblastic Leukemia* and *Acute Myeloid Leukemia*.

To differentiate each type of leukemia, the notion of hematopoiesis has to be introduced first. Blood cells are created in the bone marrow. Each blood cell is different but was at first a stem cell from the bone marrow. The hematopoiesis is the succession of transformations that leads a stem cell to a particular blood cell. For example, a stem cell can transform into two cells: a myeloid stem cell or a lymphoid stem cell. From other possible states, a myeloid stem cell can change into a macrophage while a lymphoid stem cell into a T cell or a B cell. Figure 6.14 shows a visualization of blood cell hematopoiesis with intermediate cell populations called progenitors.

Acute lymphoblastic leukemia is a leukemia where lymphoid stem cells develop into lymphoid blasts or lymphoblasts. Those cells are produced in high numbers and do not develop into normal lymphocytes. As a consequence, other blood cells have difficulties to develop which unbalances the hematopoiesis [65]. Similarly, myeloid leukemia blasts are myeloid stem cell

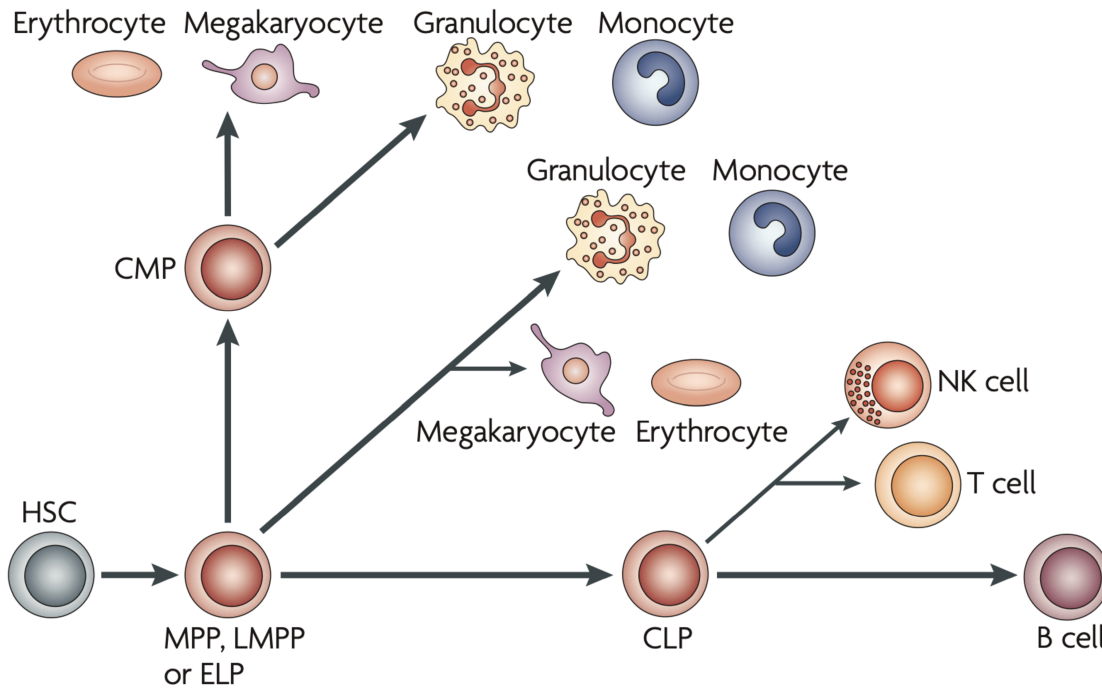


Figure 6.14: Hematopoiesis tree [19]. **HSC**: Hematopoiesis stem cell. **CMP-MMP-CLP-LMPP-ELP**: Progenitor cells which are intermediate cell populations.

that develop into abnormal cells called myeloblasts. Those cells do not develop into functional myeloid cells which can result in a risk of infections [91].

The main treatment against leukemia is a hematopoietic stem cell transplantation usually combined with chemotherapy. In the case of leukemia, chemotherapy stops cell development which amounts to with immunosuppression [20]. After this step, a hematopoietic stem cell graft sample is given to the patient. The aim of this treatment is to replace the previously malfunctioning immune system with a healthy sample. Because the older cell populations are outnumbered by the new graft cells, the new cells supersede the older immune system [79].

6.3.1 Application to 8-marker datasets

In this section, the method *CTFlowHD* was applied to a graft sample used in the treatment of leukemia. The goal of this study was to examine if it is possible with the method to interpret biologically this real dataset. The FCM dataset acquired contained $N = 1509790$ cells after pre-processing steps. The values of fluorescence for each cell was obtained for $M = 8$ markers: CD11b, CD3, PD-L1, CD33, CD14, CD34, CD15, HLA-DR.

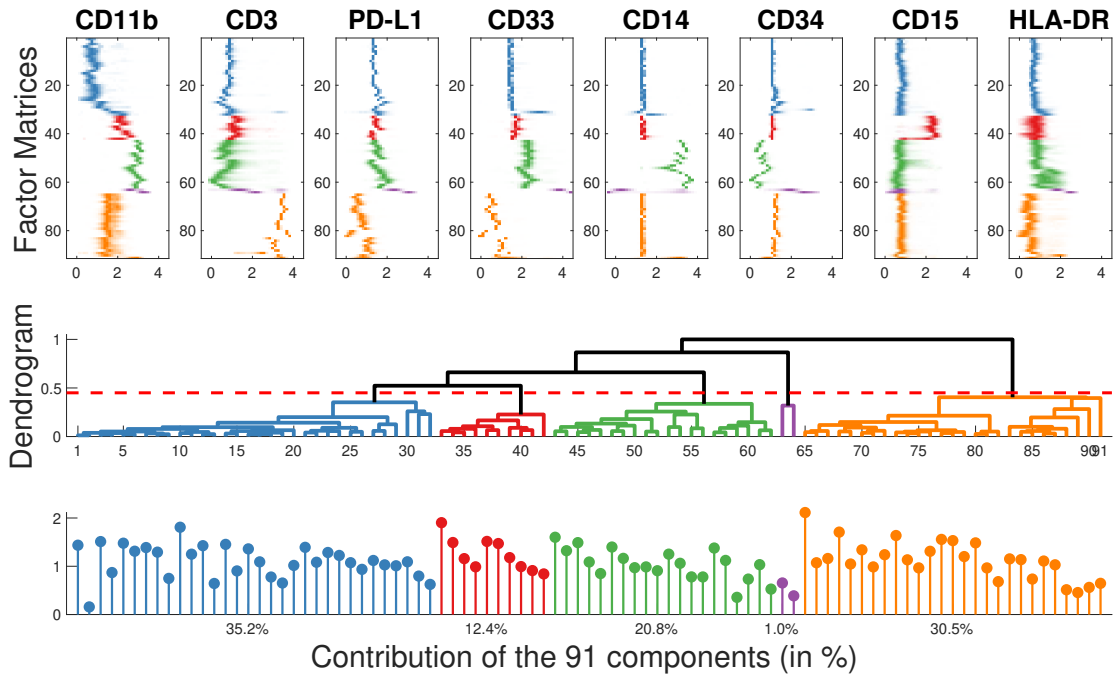


Figure 6.15: *CTFlowHD* applied to a graft sample dataset with 8 markers (full coupling).

First, the following parameters were used to apply *CTFlowHD*:

- $I = 30$ bins per dimension,
- all triplets were considered $\mathcal{T} = \mathcal{T}_{\text{all}}$ ($T = 56$ triplets),
- $R = 100$ rank-one terms,
- $T_1 = 2,000$ outer iterations and $T_2 = 20$ inner iterations.

Complete linkage hierarchical clustering was applied to the results of PCTF3D and plotted in Figure 6.15. Note that the figure only shows $R = 91$ components, because the remaining 9 factors had a λ_r value numerically equal to 0. With an expertise help to define each cell population clustered in this figure, each cluster was identified as a cell population⁹. Table 6.3 permits to compare the results obtained with *CTFlowHD* and a manual data analysis. In this experiment, even if no ground truth is available, the results obtained with those two methods are similar. Moreover, it is possible to locate each cluster of Table 6.3 on the hematopoiesis tree of Figure 6.14. Indeed, the dendrogram clustering separated granulocytes, MDSC which are particu-

⁹Here, we want to thank Anne-Béatrice Notarantonio for providing key expertise for this dataset but also more generally for her whole help regarding the flow cytometry context of this work.

Table 6.3: Identification of cell populations obtained with *CTFlowHD* with hierarchical clustering (see Figure 6.15). MDSC stands for Myeloid-derived suppressor cells.

<i>CTFlowHD</i>		Manual Gating		
Hierarchical clustering	Proportion	Cell population	Marker expression	Proportion
Red cluster	12.4%	Granulocytes	CD15+	10.7%
Green cluster	20.8%	MDSC	CD14+/HLA-DR-	16%
Purple cluster	1.0%	Stem cells	CD34+	0.61%
Orange cluster	30.5%	T cells	CD3+	30.1%
Blue cluster	35.2%	Negative cells	-	28.3%

lar cells generated from CMP (Common Myeloid Progenitors), stem cells which are referenced as HSC (Hematopoietic Stem Cells) in Figure 6.14 and T cells. Because this FCM experiment does not include discriminant markers for other cell populations – i.e. B cells, negative cells corresponds to cells that are neither of cell populations mentioned in the Table 6.3.

Finally, another experiment for this dataset was conducted with the same parameters except \mathcal{T} . Indeed, PCTF3D was applied for a random strategy with $T = 14$ triplets ($1/4$ of all triplets). Figure 6.16 shows the complete linkage visualization for partially coupled tensor factorization. Even if the components are not exactly the same compared to full coupling, the resulting proportions and fluorescence properties are very similar to the one obtained in Figure 6.15. Note that for a balanced coupling strategy with the same number of triplets, the same results have been observed. In order to compare both PCTF3D results, it is also possible to plot the signatures for each cluster and for each triplet setup like in Figures 6.17 and 6.18. Each cluster signature features the M fluorescence properties. For example, the orange cluster is CD3+. Both figures permit to differentiate cell populations in terms of fluorescence and also permit to compare the results from one strategy to another. In addition to displaying characteristics for cell populations, dendrogram visualizations permit to evaluate how a cell population can be separated into smaller clusters.

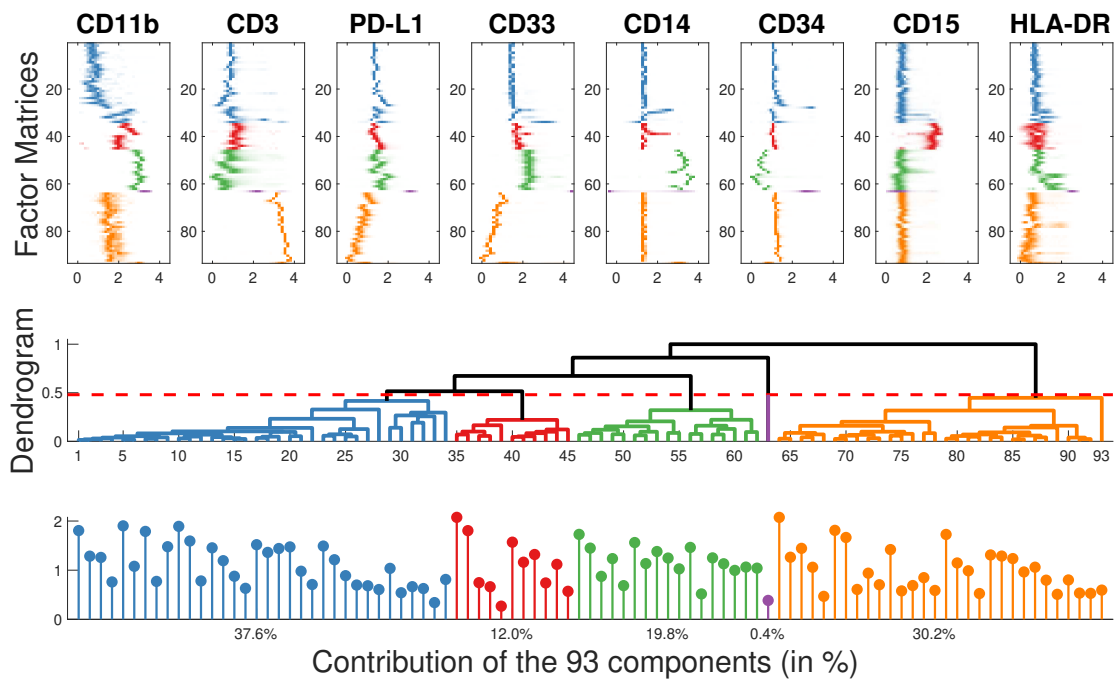


Figure 6.16: *CTFlowHD* applied to a graft sample dataset with 8 markers (partial coupling).

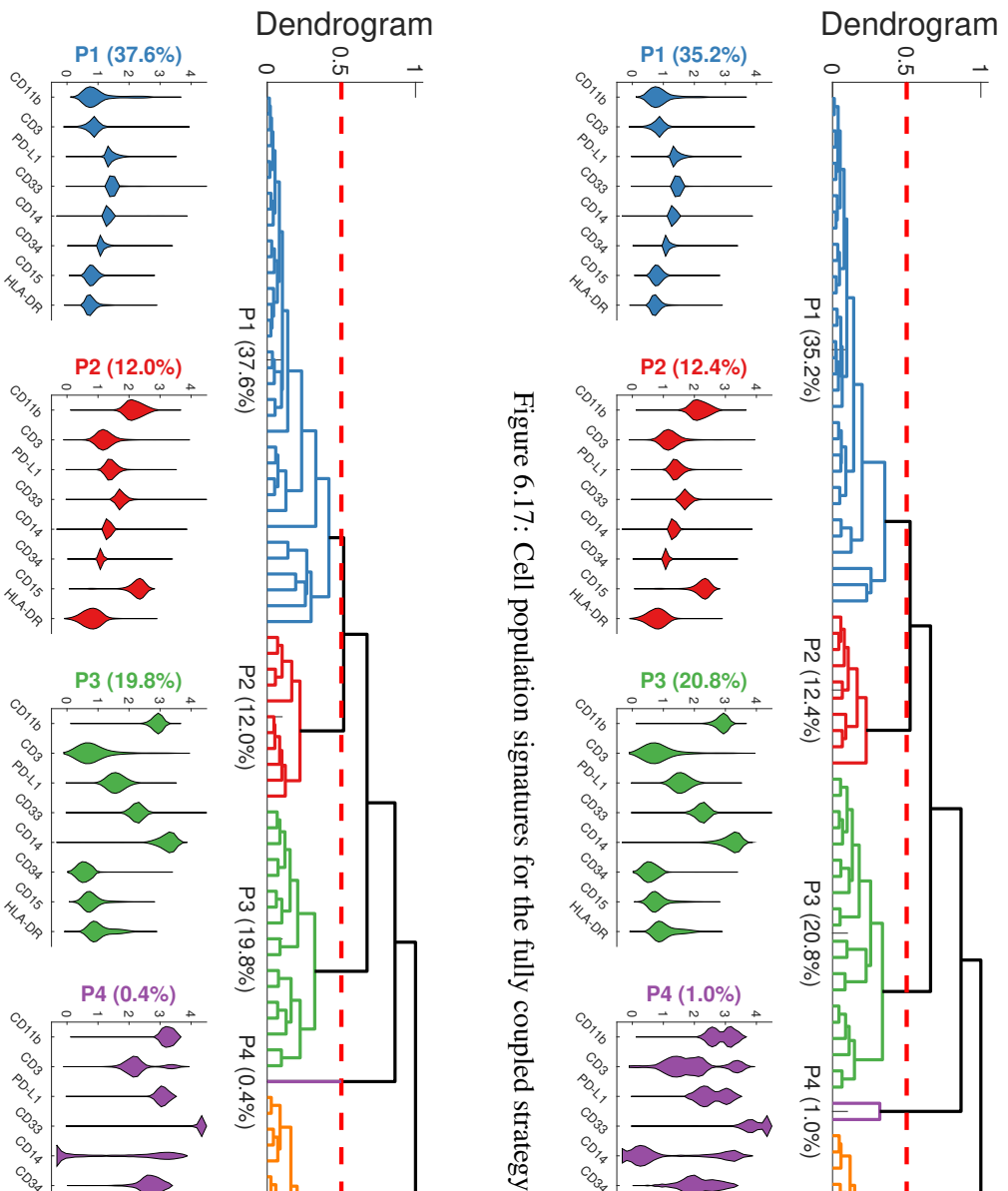


Figure 6.17: Cell population signatures for the fully coupled strategy

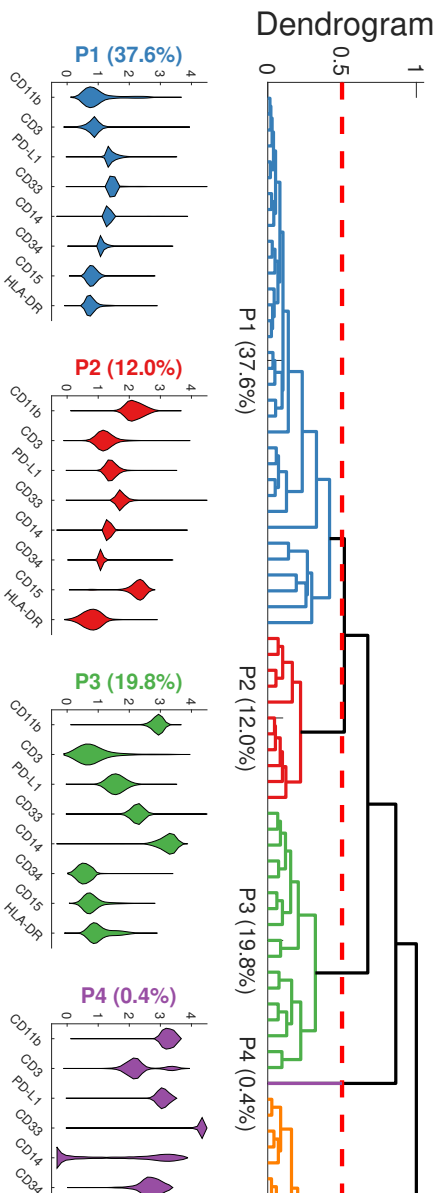


Figure 6.18: Cell population signatures for the partially coupled strategy

6.3.2 Application to 24-marker datasets

In this section, we are interested in the study of a graft relapse sample. It is possible, after a treatment of combined graft and chemotherapy, that a patient have a leukemia relapse sometimes months after the treatment. This phenomenon can occur if chemotherapy did not eliminate all problematic cells. Even if they have been outnumbered by the graft cells, cancer cells can be produced again after the treatment. Moreover, relapse risk is enhanced by the problem of Graft versus Host Disease (GvHD). Indeed, grafted T cells can proliferate and activate on organs of the patient [34]. There exists counter-treatments to control GvHD but these are interfering with the first treatment called Graft versus Tumor (GvT). There are several objectives to this study. First, it would be a big advantage clinically if relapses were predicted instead of diagnosed. Moreover, if relapses are predicted with a specific cell population, it would give valuable information on why relapses occur.

To conduct this study, samples were taken from patients at different moments of the treatment: the day of the graft (D0), 30 days after the draft (D30) and 90 days after the graft (D90). In some cases, relapse had already occurred at these milestones. To analyze these datasets, samples were analyzed in a flow cytometer with a pool of $M = 24$ markers. *CTFlowHD* was applied on a sample at D30 that contains $N = 1.8$ million cells with the following parameters:

- $I = 20$ bins per dimension,
- $T_1 = 2000$ outer iterations and $T_2 = 20$ inner iterations,
- a balanced coupling \mathcal{T} with 200 triplets (approximately 10% of all triplets),
- $R = 100$ components.

The computation time for this dataset was about 45 minutes. Note that *CTF3D* was also applied and its computation time for this experiment was nearly 12 hours. The biological interpretation of this result (see Figure 6.19) is still in progress. However, *CTFlowHD* is able to reconstruct a 24-variables NBM with all triplets. Because this dataset features a high number of cells and a high number of variables, SPADE and t-SNE cannot be applied to this dataset.

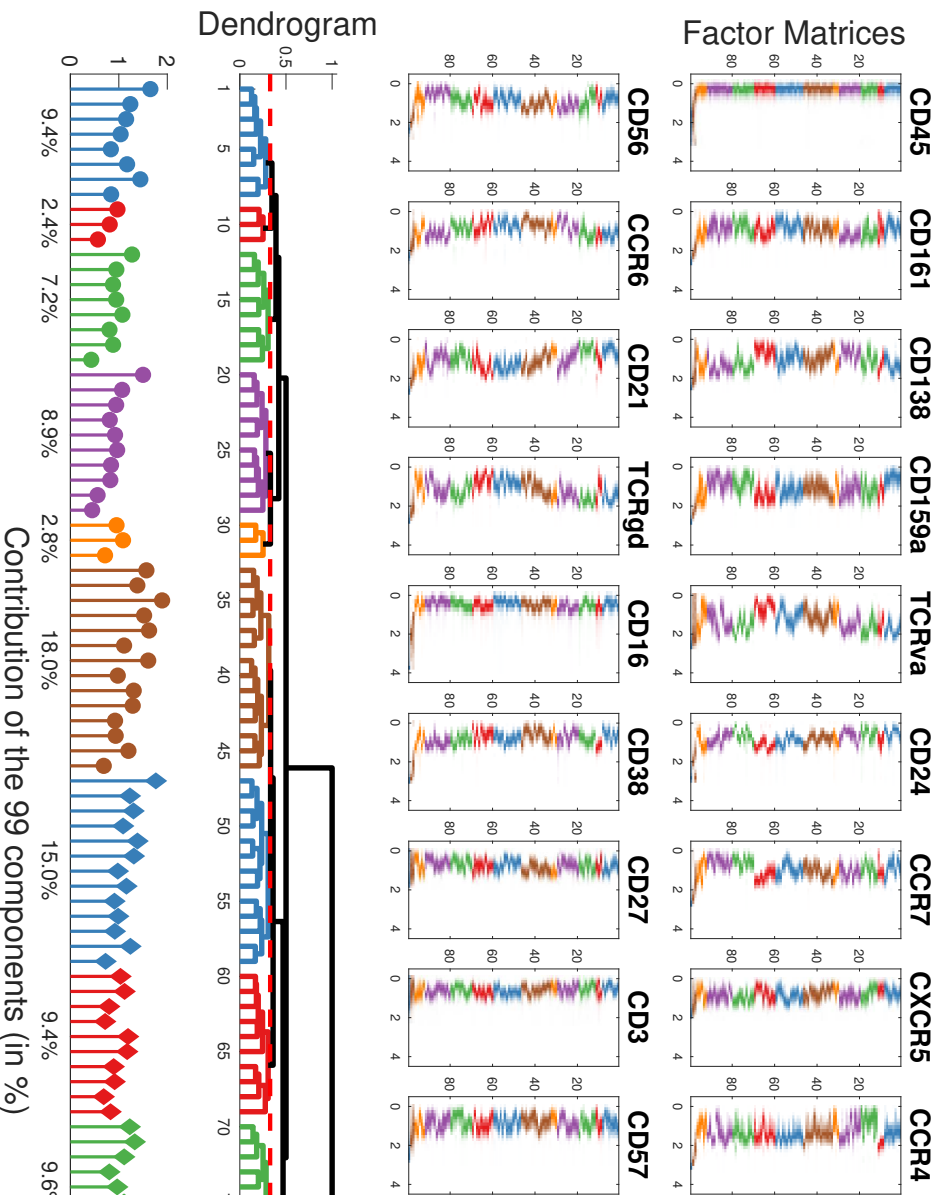


Figure 6.19: CTFloWHD applied to a 24-marker FCM dataset.

6.4 Conclusion of Chapter 6

This chapter presented a new method for flow cytometry data analysis. This method is based on the coupled tensor algorithm PCTF3D presented in Chapter 4. The workflow of this method is well segmented, as all steps of the method are independent. Indeed, pre-processing steps are performed then PCTF3D is applied to obtain the NBM factors. After NBM factors are estimated, a visualization step is needed and several tools for visual interpretation are proposed in this chapter. For a better understanding of the method by end users, tools that are currently used in the community of flow cytometry were presented. The main advantage of *CTFlowHD* is to have an overview of a dataset along with a panel of visualizations at a low computational load.

We believe that the results provided in this chapter are promising but the method still has a room for improvements. For example, our data pre-processing routine does not handle the pooling of multiple datasets by normalizing FCM datasets. This step would allow for analyzing data from two groups instead of analyzing individual samples one by one. Concerning the biological interpretation, we hope that our method will help provide results in the research on cancer relapse for example. Because of *CTFlowHD*'s versatility, other clustering/visualization methods may be applied to NBM component factors. The development of new tools for *CTFlowHD* must be driven by end users, as their needs may permit to enlarge the abilities of *CTFlowHD*.

General Conclusion

Conclusion of the manuscript

The objective of this manuscript was to address the problem of flow cytometry data analysis with a probabilistic approach using low-rank tensor decompositions. By considering cell fluorescence values as random variables, the problem of flow cytometry data analysis can be viewed as a joint PDF estimation problem.

In a first chapter, the problem of flow cytometry data analysis was presented. Flow cytometry is a widely used method that permit to analyze biological properties of high-numbers of cells at low cost. While the amount of data increased in the past decades, flow cytometry data was mainly analyzed manually with gating: an iterative, subjective and time-consuming method. Therefore, automated flow cytometry data analysis methods were developed in the so-called field of Computational Flow Cytometry. As it is shown in the first chapter of this thesis, challenges remain such as reducing runtimes, scale to more variables and more cells, or detect rare cell populations. The idea of this work was then to tackle some of the remaining challenges of flow cytometry data analysis.

In the third chapter, we recall that the problem of PDF estimation is governed by the curse of dimensionality at several levels. Indeed, when the number of variables increases, the complexity of PDF estimation increases exponentially. In this work, we proposed to use graphical models and more particularly the Naive Bayes Model (NBM). This model reduces the complexity of the problem since its complexity remains linear with the number of dimensions. Moreover, the NBM can be interpreted as a Canonical Polyadic Decomposition (CPD). Therefore, the problem of PDF estimation reduces to the estimation of CP tensor decomposition factors. At this point, another level of curse of dimensionality is observed. Indeed, obtaining the CPD requires a quantity of data increasing exponentially with the number of dimensions. This issue was tackled by the CTF3D of [74] which proposed a coupled tensor factorization algorithm that uses order-3

marginals to estimate the CPD. Nevertheless, the curse of dimensionality still remains since the number of 3D marginals which increases cubically with the number of dimensions.

In the fourth chapter of this manuscript, we proposed a PDF estimation algorithm called Partially Coupled Tensor Factorization of 3D marginals or PCTF3D that estimates a PDF with only a subset of 3D marginals. Coupling a subset of marginals permits to control the complexity of the method, as PCTF3D's complexity which directly depends on the number of 3D marginals. Different coupling strategies, *i.e.* the choice of the marginals to be considered, were examined in the framework of 3-uniform hypergraphs. We proposed an algorithm that generates balanced couplings which ensures that all variables are evenly represented. Finally, numerical experiments have shown that PCTF3D reduces computational load while keeping similar performances compared to CTF3D. A comparison with another method showed that the Kullback-Leibler Divergence (KLD) [23] may be a more adequate distance, especially for the detection of rare cell populations.

Then, uniqueness guarantees of our coupled model were examined in the fifth chapter. Unlike previous works, the simplex constraints were handled. First, the recoverability of the coupled models was studied. To do so, we proposed an algorithm that searches for the maximum recoverable rank. This algorithm is computing the rank of the Jacobian matrix for our constrained parametrization and works for any specific case of coupling but becomes computationally expensive for big-sized tensors. By applying this algorithm for several coupling strategies, various necessary conditions were proposed as well as a conjecture on the recoverability bound in the general case. Moreover, random strategies have been shown to induce defective cases for which recoverability is necessarily not guaranteed for low values of rank. Balanced couplings do not lead to those defective cases and reach the maximum recoverability bound. On the other hand, we proposed a sufficient identifiability condition by using the framework of algebraic geometry and polynomial additive model identifiability. The theorem was proven using the semi-algebraic structure of the set of constraints of the coupled tensor decompositions. Compared to the sufficient condition of [74], our theorem gives improved identifiability guarantees.

Finally, the new algorithm PCTF3D was applied to flow cytometry data analysis. Indeed, assuming a NBM on flow cytometry data can be interpreted as a separation of cells into a set of cell populations which can be fully described with the factors of the model. In that context, we proposed a workflow of analysis named *CTFlowHD* for Coupled Tensor factorization for Flow cytometry in High Dimensions. The main advantage of our method is its versatility. Indeed, the dimension reduction provided by the NBM enables us to apply any clustering/visualization

steps afterwards at little cost.

Perspectives

The development of a flow cytometry data analysis tool has been done during the course of this thesis. First, we hope to popularize our approach in the flow cytometry community. As mentioned earlier, visualization tools may be further developed to enhance the versatility of our workflow. We do believe that the future improvements of *CTFlowHD* are going to answer flow cytometry end-users problems.

Concerning the coupled tensor factorization method, PCTF3D may be used in other application setups. Moreover, the comparison with a KLD-based method highlighted the fact that the KLD may be a better distance for the optimization problem. Nevertheless, the complexity of the KLD-based method presented in Chapter 4 depends too much on the number of samples. A promising perspective may be to develop an algorithm that couples lower-order marginals while minimizing a KL divergence.

In terms of coupled tensor model identifiability, recoverability bounds obtained for random couplings showed defective cases. We think that it is possible that other defective cases can be studied, especially for higher number of dimensions. Moreover, the different conjectures proposed in Chapter 5 shows how much it can be done on that matter. Indeed, because of the inherent combinatorial structure of couplings, it may be challenging to prove results in the general case.

Appendix A

Brochure of BD for FACSymphony flow cytometers



BD FACSymphony™ Flow Cytometer

Special Order Research Product

Customized solutions for high-parameter cell analysis



Driving deeper scientific insights

High-parameter flow cytometry is a powerful analytical tool that enables scientists to identify and analyze distinctive phenotypes in heterogeneous populations. The BD FACSymphony™ flow cytometer is a novel cell analyzer that leverages the inherent benefits of flow cytometry and enables the simultaneous measurement of up to 50 different characteristics of a single cell.

This advanced instrument features an ultra-quiet VPX electronics system that supports up to 50 high-performance photomultiplier tubes (PMTs) and improves detection sensitivity to enable you to identify and analyze rare cell types and events. The capabilities of this platform technology uniquely allow you to conduct deep and broad phenotyping and gain richer scientific insights by fully leveraging the broad portfolio of BD Horizon Brilliant™ reagents.

With early access to newly developed BD Horizon Brilliant dyes, this platform helps you to overcome research challenges such as collecting maximal information from a precious sample and increases lab throughput with broad phenotyping panels that combine multiple cell line specific panels.

This highly customizable platform can be configured so you can select from multiple laser wavelengths and power ratings and choose the positions of decagon detection arrays to address the requirements of your specific research application.



Customizable models provide flexibility for your research lab

BD FACSymphony™ A5

- Configure to your needs today with room for growth tomorrow
- Up to 50 detection parameters (including FSC and SSC) featuring decagon arrays for up to 10 parameters on a single laser line
- Select and configure up to a maximum of 10 lasers* from various wavelengths with multiple power ratings



*Dependent on laser choice

Custom optics for your application

BD SORP 2016 – 25 Wavelength Laser Portfolio



355 nm	505 nm	637 nm
375 nm	514 nm	640 nm
405 nm	532 nm	647 nm
420 nm	552 nm	660 nm
445 nm	561 nm	685 nm
458 nm	568 nm	730 nm
460 nm	588 nm	785 nm
473 nm	592 nm	980 nm
488 nm	628 nm	

In the spirit of Special Order Research Products (SORP), there are 25 laser wavelengths to choose from to optimally configure your BD FACSymphony instrument for your specific research application. Additionally, there are multiple power ratings for most lasers that can be adjusted, stored and recalled using digital laser command and control functionality.

Innovation in detection array technology has allowed for a decagon formation to detect 10 parameters on a single laser line. The arrays can be configured on the laser of your choice.



Fluorochrome availability and excitation characteristics across various wavelengths should be discussed during the configuration process to identify the best use of reagents for your research. Optimal laser power settings for certain fluorochromes may be available.

Highlighted wavelengths are common laser choices

Broad portfolio of high-quality dyes and conjugates expand options for experimental design

BD's broad portfolio of fluorochromes featuring the BD Horizon Brilliant™ dyes offers flexibility for experimental design. Leverage the principles of antigen density and relative fluorochrome brightness to optimally design your panel.

Note: For specificities not yet available in the catalog or through BD OptiBuild™ custom reagents, high-parameter users have access to a small scale custom conjugation program for the BD Horizon Brilliant™ dyes, including early access to the high-parameter dye menu described below.

BD OptiBuild™ custom reagents offer on-demand access to hundreds of specificities associated with a range of BD Horizon Brilliant™ dyes, available in small sizes with quick turnaround times. This new portfolio of over 1,000 recently released conjugates complements the existing catalog reagents with a wide selection of cell surface antibodies that previously had few color options to choose from. Revisit this portfolio often, as we continue to expand the BD OptiBuild offering so that you can simplify the addition of markers to your experiments without the limitations of reagent availability.

Excitation Laser Line	Channel	Recommended Filter	Fluorochrome	Ex-Max (nm)	Em-Max (nm)	Relative Brightness
UV	1	379/28	BD Horizon™ BUV395	348	395	■ ■ ■ ■
	2	515/30	BD Horizon™ BUV496	348	496	■ ■ ■ ■
	3	585/15	BD Horizon™ BUV563	348	563	■ ■ ■ ■
	4	•	BD Horizon™ BUV615-P	349	616	■ ■ ■ ■
	5	670/25	BD Horizon™ BUV661	348	661	■ ■ ■ ■
	6	740/35	BD Horizon™ BUV737	348	737	■ ■ ■ ■
	7	820/60	BD Horizon™ BUV805	348	805	■ ■ ■ ■
Violet	8	450/40	BD Horizon™ BV421	407	421	■ ■ ■ ■
		450/40	BD Horizon™ V450	404	448	■ ■ ■ ■
		450/40	Pacific Blue™	401	452	■ ■ ■ ■
	9	525/40	BD Horizon™ BV480	436	478	■ ■ ■ ■
		525/50	BD Horizon™ V500	415	500	■ ■ ■ ■
		525/40	BD Horizon™ BV510	405	510	■ ■ ■ ■
		•	BD Horizon™ BV570	407	574	■ ■ ■ ■
	10	•	BD Horizon™ BV570	407	574	■ ■ ■ ■
	11	610/20	BD Horizon™ BV605	407	602	■ ■ ■ ■
	12	660/20	BD Horizon™ BV650	407	650	■ ■ ■ ■
13	710/50	BD Horizon™ BV711	407	711	■ ■ ■ ■	
14	•	BD Horizon™ BV750-P	407	748	■ ■ ■ ■	
15	780/60	BD Horizon™ BV786	407	786	■ ■ ■ ■	
Blue	16	530/30	BD Horizon™ BB515	490	515	■ ■ ■ ■
		530/30	Alexa Fluor® 488	495	519	■ ■ ■ ■
		530/30	FITC	494	519	■ ■ ■ ■
	17	•	BD Horizon™ BB630-P	484	631	■ ■ ■ ■
	18	•	BD Horizon™ BB660-P	484	667	■ ■ ■ ■
	19	695/40	PerCP**	482	678	■ ■ ■ ■
		•	BD Horizon™ BB700-P	484	695	■ ■ ■ ■
20	695/40	PerCP-Cy™ 5.5**	482	695	■ ■ ■ ■	
Yellow-Green	21	•	BD Horizon™ BYG584-P	563	584	■ ■ ■ ■
		575/26	PE*	496	578	■ ■ ■ ■
	22	610/20	BD Horizon™ PE-CF594*	564	612	■ ■ ■ ■
	23	670/14	PE-Cy™ 5*	564	667	■ ■ ■ ■
24	780/60	PE-Cy™ 7*	564	785	■ ■ ■ ■	
Red	25	660/20	APC	650	660	■ ■ ■ ■
		660/20	Alexa Fluor® 647	650	668	■ ■ ■ ■
	26	730/45	BD Horizon™ APC-R700	652	704	■ ■ ■ ■
		730/45	Alexa Fluor® 700	696	719	■ ■ ■ ■
	27	780/60	APC-Cy7	650	785	■ ■ ■ ■
		780/60	BD™ APC-H7	650	785	■ ■ ■ ■

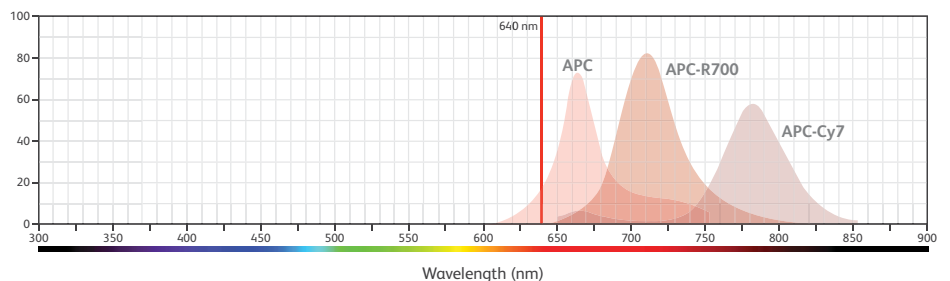
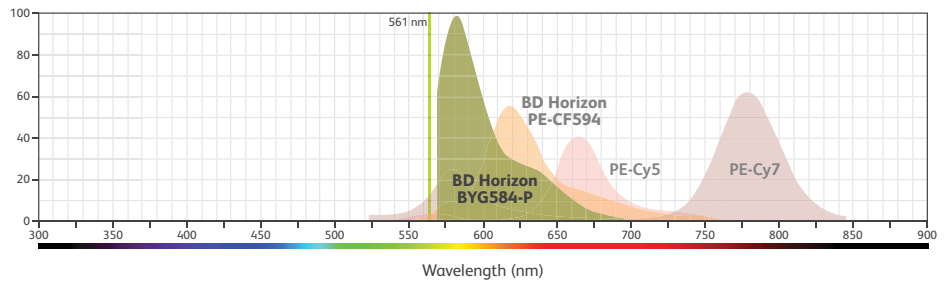
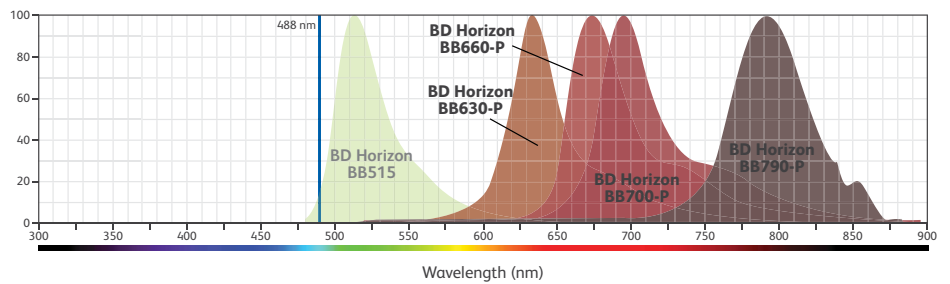
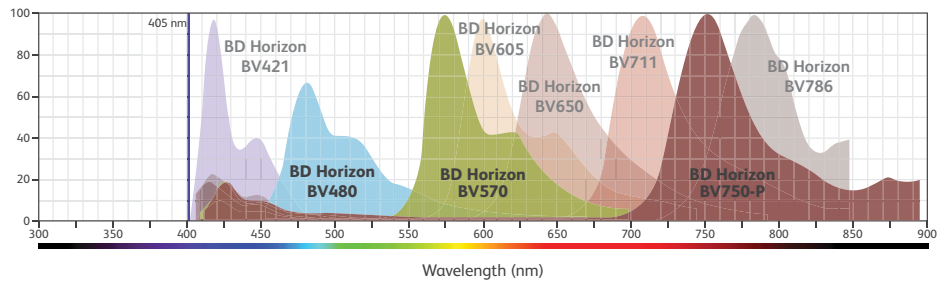
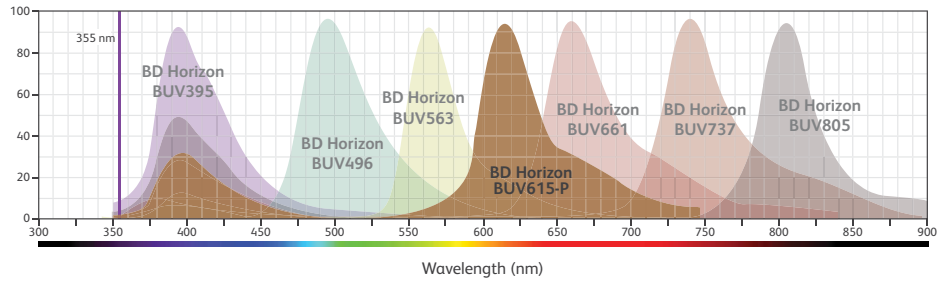
•Filter recommendations will be provided based on instrument configuration
 *Excited by 488 nm, 532 nm, and 561 nm
 **Excited by 488 nm and 532 nm



Prototypes of BD Horizon Brilliant™ dyes (-P)

BD Life Sciences is committed to continuing to develop new BD Horizon Brilliant™ dyes across various laser lines to improve spectral properties of dyes and minimize the need for compensation in higher order panels.

BD FACSymphony owners receive early access to a suite of prototype dyes for use in high-parameter panel design. Although these dyes are near completion and have received initial quality specifications, they may undergo additional development that could result in minor performance changes. The -P nomenclature indicates the prototypic nature of the dyes, and any significant changes to the structure of the dye to optimize performance will be appropriately communicated to customers.



Exclusive high-parameter reagent access and specialized support

Reagent availability is critical for high-parameter panel design. The high-parameter custom reagent program is specifically designed to cater to the needs of researchers looking to achieve >20 parameter flow cytometry analysis.

While many of the BD Horizon and BD Horizon Brilliant dyes are featured in the BD catalog, the prototype dyes are exclusively available through the high-parameter custom reagent program. This program allows you to acquire small-scale custom reagents on the prototype dyes and any other dyes in the BD Horizon Brilliant family to optimally design your complex multicolor panels.

Additionally, all systems will come with access to reagents and specialized support to get you up and running as quickly as possible. This includes onsite and offsite support from our team of dedicated high-parameter application specialists to discuss your research goals and consult with you about reagent choice to simplify your panel design activities.

The reagents provided with purchase of a BD FACSymphony instrument will assist you in setting up your instrument, identifying spectral characteristics when running various fluorochromes simultaneously, and beginning design of your initial panels. The reagents will include a fluorochrome evaluation kit, a suite of human CD4 SK3 reagents in nearly every color option to evaluate detection capabilities of your custom BD FACSymphony configuration. The kit also includes samples of specificities on the color of your choice for your research needs. Where available, reagent access includes high concentration, mass size human reagents to avoid dilution effects in high-parameter cocktails.

As part of the program, you will be given a dedicated point of contact for ordering reagents, contacting the specialized applications team, and answering your high-parameter questions.

Contact BDB_HPS_VIP@bd.com for all your high-parameter needs.



BD Life Sciences – Biosciences Regional Offices

bdbiosciences.com/contact

Australia

Toll Free 1800.656.100
Tel 61.2.8875.7000
Fax 61.2.8875.7200

Canada

Tel 866.979.9408
Fax 888.229.9918

China

Tel 86.21.3210.4610
Fax 86.21.5292.5191

Europe

Tel 32.2.400.98.95
Fax 32.2.401.70.94

India

Tel 91.124.2383566
Fax 91.124.2383224/25/26

Japan

Nippon Becton Dickinson
Toll Free 0120.8555.90
Fax 81.24.593.3281

Latin America/Caribbean

Toll Free 0800.771.71.57
Tel 55.11.5185.9688

New Zealand

Toll Free 0800.572.468
Tel 64.9.574.2468
Fax 64.9.574.2469

Singapore

Tel 65.6690.8691
Fax 65.6860.1593

United States

US Orders 855.236.2772
Technical Service
877.232.8995
Fax 800.325.9637

Office locations are available on our websites.

Class 1 Laser Product.

For Research Use Only. Not for use in diagnostic or therapeutic procedures.

Alexa Fluor® is a registered trademark of Life Technologies Corporation.

Cy™ is a trademark of GE Healthcare. Cy™ dyes are subject to proprietary rights of GE Healthcare and Carnegie Mellon University, and are made and sold under license from GE Healthcare only for research and in vitro diagnostic use. Any other use requires a commercial sublicense from GE Healthcare, 800 Centennial Avenue, Piscataway, NJ 08855-1327, USA.

Trademarks are the property of their respective owners.

23-18657-00

BD Life Sciences, San Jose, CA, 95131, USA

bdbiosciences.com



Appendix B

PCTF3D applied to supervised classification on the MNIST dataset

In this appendix, we show how PCTF3D can be used to in a high dimensional ($M = 256$) supervised classification setup. We perform a proof-of-concept experiment of the MNIST dataset (note that the classification on MNIST dataset can be achieved with near-100% accuracy using convolutional neural networks).

MNIST dataset and data preprocessing

The MNIST dataset contains 70,000 black and white (BW) images separated in a training set ($N_{\text{tr}} = 60,000$ images) and a test set ($N_{\text{te}} = 10,000$ images). Each image is representing a hand-written digit as shown in Figure B.1. Each image is a 28×28 matrix whose pixels are integers between 0 and 255. In this experiment, labelled images are used to perform supervised classification. The number of samples associated with each digit is given in Table B.1. First, we proceed to a cropping of images where pixels are equal to 0. This permit to center images on the hand-written digit. The result of this cropping step can be found in Figure B.2a. Note

Table B.1: Number of samples for each digit in the labelled MNIST dataset.

Dataset	Digit									
	0	1	2	3	4	5	6	7	8	9
Training set	5923	6742	5958	6131	5842	5421	5918	6265	5851	5949
Test set	980	1135	1032	1010	982	892	958	1028	974	1009

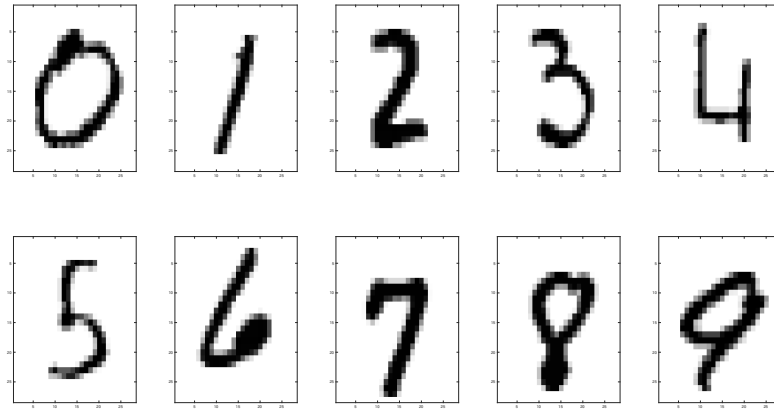


Figure B.1: Examples of digit images from the MNIST dataset.

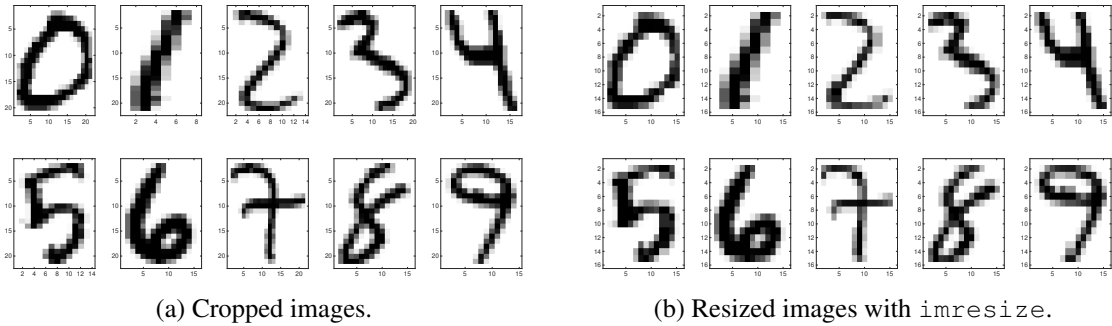


Figure B.2: Cropping and resizing of MNIST images.

that after cropping all images have not the same sizes. Therefore, the next preprocessing step is to resize cropped images to obtain 16×16 matrices (see Figure B.2b). Finally, we proceed to a last preprocessing step which consists in the quantization of the dataset. We use the MATLAB functions `imbinarize` for Figure B.3a and `imquantize` for Figure B.3b.

Supervised classification problem

The classification problem of MNIST consists in finding which digit class is represented on an image. We follow the Bayesian approach. Each datapoint is a pair (θ, \mathbf{x}) , where $\theta \in \llbracket 0, 9 \rrbracket$ is the label, considered as a random variable, and \mathbf{x} is the random vector that collects the pixels of a given image. Therefore, as images are 16×16 matrices $\mathbf{x} = (X^{(1)}, \dots, X^{(256)})$ is a vector containing $M = 256$ variables. Each random variable is taking I possibly different values, I depends on the levels of quantization chosen in the preprocessing (see Figure B.3).

We use the Bayesian framework. In this case, using the maximum a posteriori principle, the

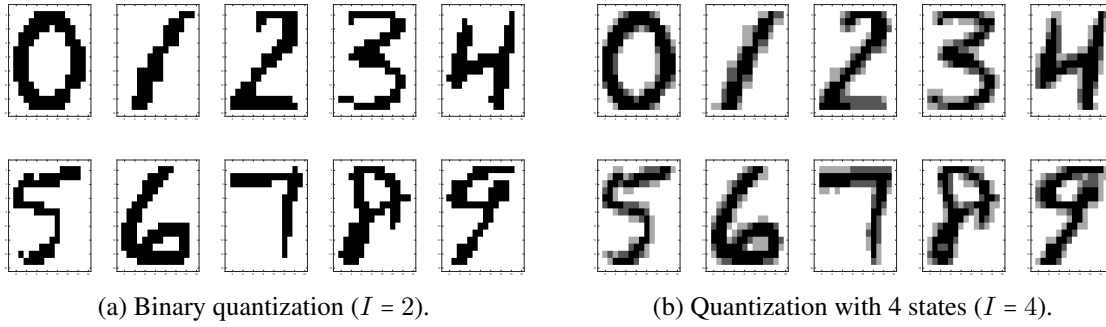


Figure B.3: Quantization of MNIST images.

classification problem can be posed as:

$$\hat{\theta} = \operatorname{argmax}_{d \in \llbracket 0,9 \rrbracket} \Pr(\theta = d \mid \mathbf{x}). \quad (\text{B.1})$$

With the Bayes theorem, the problem (B.1) is formalized with the prior and the likelihood:

$$\hat{\theta} = \operatorname{argmax}_{d \in \llbracket 0,9 \rrbracket} \Pr(\theta = d) p(\mathbf{x} \mid \theta = d).$$

To solve this problem, we have to estimate both the prior $\Pr(\theta = d)$ and the likelihood $p(\mathbf{x} \mid \theta = d)$ with a set of N_{tr} training images. The matrix $\mathbf{X}^{(\text{tr})}$ will then denote the training set as a $N_{\text{tr}} \times M$ matrix whose rows are vectorized images noted $\{\mathbf{x}_{n \cdot}^{(\text{tr})}\}_{n=1}^{N_{\text{tr}}}$. On the other hand, the labels for the training set will be noted $\Theta_n^{(\text{tr})} = \{\theta_n^{(\text{tr})}\}_{n=1}^{N_{\text{tr}}}$. Then, the prior can be estimated by counting the occurrences of each digit in the training set:

$$\Pr(\theta = d) \approx \hat{p}_d = \frac{1}{N_{\text{tr}}} \operatorname{Card}\{n \in \llbracket 1, N_{\text{tr}} \rrbracket \mid \Theta_n^{(\text{tr})} = d\}.$$

Concerning the likelihood, it can be interpreted as a set of 10 order- M PMFs of size $\overbrace{I \times \dots \times I}^{256 \text{ times}}$. To apply PCTF3D, we model each conditional distribution as a NBM:

$$p(\mathbf{x} \mid \theta = d) = \sum_{r=1}^R \Pr(L = r \mid \theta = d) \prod_{m=1}^M p(X^{(m)} \mid L = r, \theta = d).$$

We then apply PCTF3D to estimate those 10 high-dimensional distributions $\{\mathcal{H}_d\}_{d=0}^9$ which

represent the likelihoods of the problem:

$$\begin{aligned} \widehat{\mathcal{H}}_d &= \sum_{r=1}^R \widehat{\lambda}_{dr} \widehat{\mathbf{a}}_{dr}^{(1)} \circ \dots \circ \widehat{\mathbf{a}}_{dr}^{(M)} \\ \text{s.t. } \widehat{\lambda}_d &\geq 0, \quad \mathbb{1}_R^\top \widehat{\lambda}_d = 1, \\ \widehat{\mathbf{a}}_{dr}^{(M)} &\geq 0, \quad \mathbb{1}_I^\top \widehat{\mathbf{a}}_{dr}^{(M)} = \mathbb{1}_R. \end{aligned} \quad (\text{B.2})$$

In Equation (B.2), the loading vectors $\{\widehat{\lambda}_d\}_{d=0}^9$ represent the estimation of the conditional probability $\Pr(L = r \mid \theta = d)$. Each factor $\widehat{\mathbf{a}}_{dr}^{(M)}$ is the estimation of the distribution of the m -th pixel $p(X^{(m)} \mid L = r, \theta = d)$. The rank of the decomposition is chosen in accordance with identifiability results presented in Chapter 5. We therefore chose values of R below 6 for 16×16 images.

PCTF3D triplet strategy

To apply PCTF3D, the coupling hypergraph \mathcal{T} must be chosen. We use in this experiment 2 different strategies which will be compared in the following. Choosing a triplet in this experiment corresponds to the choice of triplet of pixels. The first triplet strategy that we consider is completely random. In the second strategy, we choose triplets with pixels close to each others. In Figure B.4, the two different sorts of triplets inside \mathcal{T}_{pix} are presented. For $L \times l$ images, the number of triplets of \mathcal{T}_{pix} is noted T_{pix} :

$$T_{\text{pix}} = (16 - 2) \times 16 + (16 - 2) \times 16,$$

which gives $T_{\text{pix}} = 448$ triplets. Each 3D marginal $\mathcal{H}_d^{(j k \ell)}$ is a $I \times I \times I$ tensor which represents either $I^3 = 8$ entries for binary images and $I^3 = 64$ entries for images quantized with $I = 4$. Hence, empirical histograms $\widetilde{\mathcal{H}}_d^{(j k \ell)}$ can be accurately estimated with the amount of data available.

Classifier and application of the Maximum A Posteriori (MAP) principle

For a test image whose vectorization is noted $\mathbf{x}^{(te)}$, the Maximum A Posteriori (MAP) classifier is the function $\widehat{\theta}$ that solves the following problem:

$$\widehat{\theta}(\mathbf{x}^{(te)}) = \underset{d \in \llbracket 0,9 \rrbracket}{\operatorname{argmax}} \Pr(\theta = d) \Pr(\mathbf{x} = \mathbf{x}^{(te)} \mid \theta = d).$$

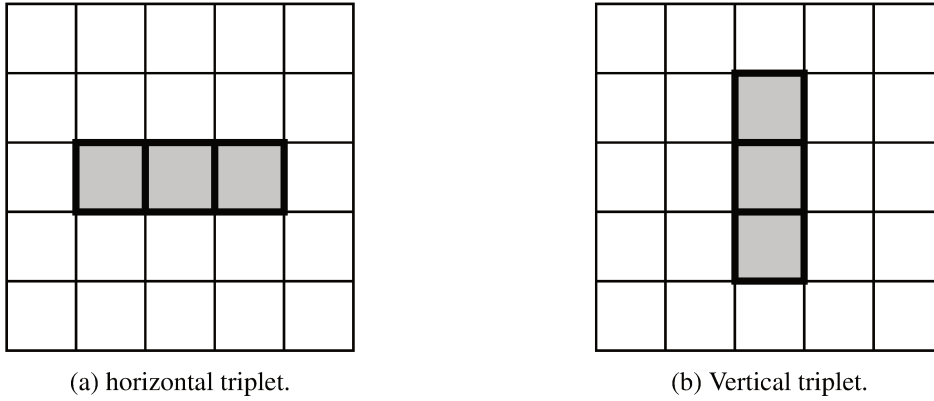


Figure B.4: Representation of triplets as a 3-neighbor mapping.

Usually, the logarithm of this cost function is maximized which gives the same solution since the logarithm is an increasing function:

$$\widehat{\theta}_{\text{MAP}}(\mathbf{x}^{(te)}) = \operatorname{argmax}_{d \in [[0,9]]} \left[\log(\Pr(\theta = d)) + \log(\Pr(\mathbf{x} = \mathbf{x}^{(te)} \mid \theta = d)) \right].$$

As we saw previously, the two terms of this sum are known. One is known by counting the occurrences of each digit in the training dataset and the other is obtained via PCTF3D. For a rank $R = 1$ model, this classifier can be simplified.

$$\widehat{\theta}_{\text{MAP}}(\mathbf{x}^{(te)} = (i_1, \dots, i_M)) = \operatorname{argmax}_{d \in [[0,9]]} \left[\log(\widehat{p}_d) + \sum_{m=1}^M \log(\widehat{a}_{dri_m}^{(m)}) \right].$$

However, for $R > 1$ it is not possible to simplify the expression of the likelihood like in the previous equation. This can lead to numerical issues, as products of 256 very small factors are needed to implement this classifier. To overcome this issue, we used the following classifier which is not the MAP classifier:

$$\widehat{\theta}_{\text{1D}}(\mathbf{x}^{(te)} = (i_1, \dots, i_M)) = \operatorname{argmax}_{d \in [[0,9]]} \left[\log(\widehat{p}_d) + \sum_{m=1}^M \log \left(\sum_{r=1}^R \widehat{\lambda}_{dr} \widehat{a}_{dri_m}^{(m)} \right) \right].$$

This classifier is noted $\widehat{\theta}_{\text{1D}}$ because this classifier uses the 1D marginals estimations to proceed to the estimation of a new label. We compare those two classifiers with the Bayes classifier that

considers an independent model for the joint distribution:

$$\widehat{\theta}_{\text{ind}}(\mathbf{x}^{(te)} = (i_1, \dots, i_M)) = \underset{d \in \llbracket 0, 9 \rrbracket}{\text{argmax}} \left[\log(\widehat{p}_d) + \sum_{m=1}^M \log(\widehat{h}_{di_m}^{(m)}) \right].$$

For the Bayes classifier $\widehat{\theta}_{\text{ind}}$, the term $\widehat{\mathbf{h}}_d^{(m)}$ is the estimation of the m -th 1D marginal for the d -th digit. Unlike for the classifier $\widehat{\theta}_{\text{MAP}}$, this estimation is not provided by PCTF3D but rather by counting the occurrences of each pixel value:

$$h_{di_m}^{(m)} = \frac{1}{N} \text{Card} \{ n \in \llbracket 1, N_{\text{tr}} \rrbracket \mid \Theta_n^{(\text{tr})} = d, x_{nm}^{(\text{tr})} = i_m \}. \quad (\text{B.3})$$

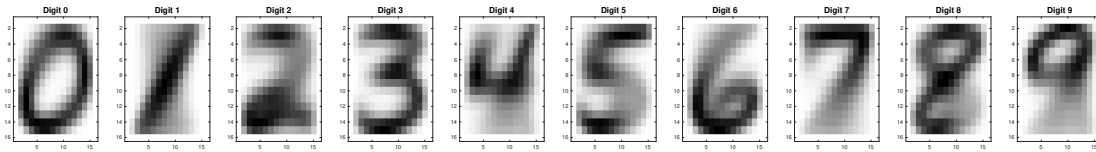
Results

The result of PCTF3D is a set of M factor matrices, each one representing one pixel obtained either with PCTF3D or Equation (B.3). Therefore, each 1D-marginal is then representing how a pixel is distributed. It is possible to plot each digit "feature image" y_d whose pixels are defined like the expectation value regarding the likelihood obtained with PCTF3D:

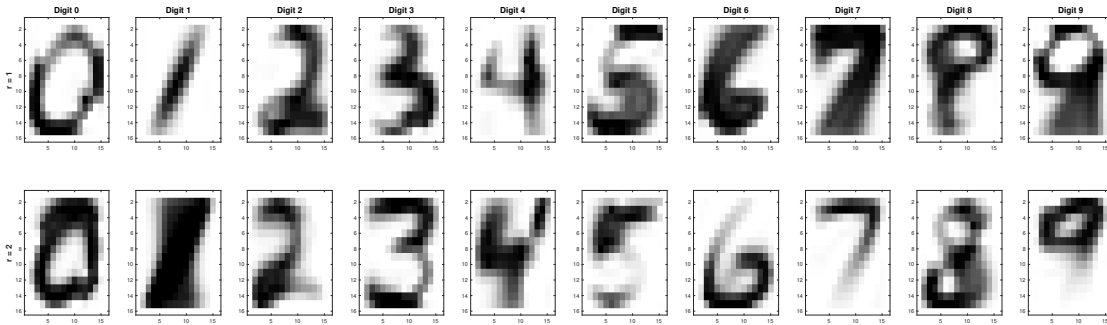
$$y_{dm} = \mathbb{E} \left[\text{Pr}(X^{(m)} \mid \theta = d) \right].$$

Note that for binary images ($I = 2$), the expression of $y_{dm} = \text{Pr}(X^{(m)} = 1 \mid \theta = d)$. The feature images are plotted for each rank-one term in Figure B.5 in the case of a Bayes and independent estimator and PCTF3D. We see that the PCTF3D estimation with random triplets results with artifacts in the feature images. However, both of PCTF3D estimations provide rank-one terms that can be interpreted as digit images.

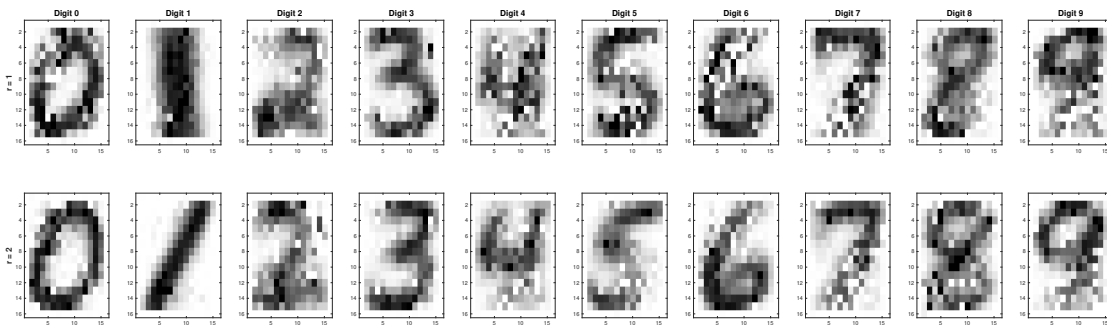
Regarding the accuracy of the classification, results can be found in Table B.2. The accuracy is defined as the percentage of good label estimation over N_{te} trials. This shows that PCTF3D have similar accuracy results compared to the Bayes classifier. However, as PCTF3D is heavier in terms of computational load, it means that PCTF3D may not be an adequate solution to the supervised classification problem.



(a) Bayes estimator (Equation (B.3)).



(b) PCTF3D with horizontal/vertical triplets.



(c) PCTF3D with random triplets.

Figure B.5: Feature images for the different likelihood estimators.

Table B.2: Accuracy of MNIST classification using PCTF3D and a Bayes classifier.

PCTF3D with pixel triplets		PCTF3D with random triplets		Bayes
$R = 1$	81.42%	$R = 1$	82.82%	83.67%
$R = 2$	83.47%	$R = 2$	83.56%	
$R = 3$	83.81%	$R = 3$	83.69%	
$R = 4$	83.54%	$R = 4$	83.69%	
$R = 5$	83.73%	$R = 5$	83.66%	
$R = 6$	83.83%	$R = 6$	83.69%	

Appendix C

Presentation of PCTF3D in the context of flow cytometry

Presentation of the algorithm PCTF3D

Author:

- Philippe Flores (<https://philippeflores.github.io>)
- e-mail: flores.philipe@gmail.com,
- github: github.com/philippeflores/fcm_ctflowhd,

PCTF3D stands for Partially Coupled Tensor Factorization of 3D Marginals. This algorithm estimates the probability density function (PDF) of random variables. In this demonstrative document, we present PCTF3D in the case of a flow cytometry experiment. This function is part of the flow cytometry data analysis package for MATLAB called CTFlowHD (github.com/philippeflores/fcm_ctflowhd).

For more information on the algorithm or its application to flow cytometry data analysis, please refer to the following reference:

- *Coupled tensor factorization for flow cytometry data analysis*, **Philippe Flores**, Guillaume Harlé, Anne-Béatrice Notarantonio, Konstantin Usevich, Maud d’Aveni, Stéphanie Grandemange, Marie-Thérèse Rubio and David Brie; In 2022 IEEE 32nd International Workshop on Machine Learning for Signal Processing (MLSP), 2022.

As mentioned before, the purpose of this script is to apply PCTF3D to a flow cytometry file. It will proceed to the following steps:

Table of contents

Load filename.....	1
Selection of variables.....	2
Definition of variable spaces.....	2
Choice of subset of triplets and computation of 3D marginals.....	3
PCTF3D optimization algorithm.....	4
Visualization of PCFT3D results.....	5

Load filename

Let us consider the flow cytometry file (.FCS) named `file_CCT.fcs`. The acronym "_CCT" stands for *Compensated, Cleaned and Transformed* which corresponds to the 3 preprocessing steps that have been applied before the application of PCFT3D. This preprocessing steps were applied using the R package named `flowAI` (doi.org/10.1093/bioinformatics/btw191):

- Gianni Monaco, Hao Chen, Michael Poidinger, Jinmiao Chen, João Pedro de Magalhães, Anis Larbi, *flowAI: automatic and interactive anomaly discerning tools for flow cytometry data*, Bioinformatics, Volume 32, Issue 16, August 2016, Pages 2473–2480.

The file `file_CCT.fcs` is then loaded into the workspace.

```
filename = "file_CCT";  
[X,fcsHdr,strFile] = loadFilename(filename);
```

This function outputs two different variables:

1. The data matrix X . This variable contains the data obtained during the flow cytometry experiment. If we note M' the number of variables measured in the flow cytometry experiment and N the number of cells, X is an array of size $N \times M'$.
2. The structure `fcsHdr` that contains the metadata related to the flow cytometry experiment. For example, it contains the label for each cytometry variable.

```
size(X)
```

```
ans = 1x2
      53881      14
```

For our example, the file `file_CCT.fcs` features $M' = 14$ variables and $N = 53,881$ cells.

Selection of variables

From the set of M' variables, we focus our study to a particular set of flow cytometry variables. We decided to select only fluorescence variables for further analysis. To help us in our choice, we define `strVariables` that stores variable labels of `fcsHdr`.

```
strVariables = {fcsHdr.par.marker}'; disp(strVariables)
```

```
{'FSC-H'      }
{'FSC-A'      }
{'SSC-H'      }
{'SSC-A'      }
{'CFSE FITC-H'}
{'CFSE FITC-A'}
{'TCR PE-H'   }
{'TCR PE-A'   }
{'CTV BV421-H'}
{'CTV BV421-A'}
{'MHCII BV510-H'}
{'MHCII BV510-A'}
{'FSC-Width'  }
{'Time'       }
```

```
indVar = [6 8 10 12];
[indVar,M,strLabel] = loadIndVar(X,indVar,fcsHdr);
```

Finally, it is possible to define the observation matrix that will permit to estimate the PDF of random variables in `indVar`. This matrix is of size $N \times M$ where M is equal to the number of variables in `indVar`. This observation matrix is denoted \mathbf{X} .

Definition of variable spaces

The PDF of random variables in `indVar` is a continuous function. PCFT3D estimated a discrete model of the PDF. Therefore, we need to define for each random variable a set of bin centroids on which the PDF will be estimated.

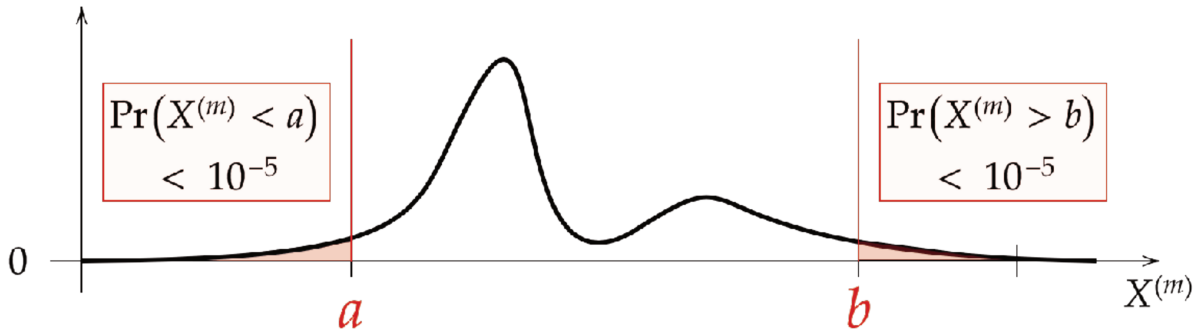
First, the number of bins for each dimension must be chosen by the user. This value is denoted I (or I) and must be chosen such that $I^3 \ll N$. Indeed, to ensure that 3D marginals are estimated with good accuracy, the number of values of a 3D histogram ($I \times I \times I = I^3$) must remain small compared to the number of cells N .

```
I = 20;
```

By choosing $I = 20$, a 3D histogram will contain $I^3 = 8,000$ values which is more than 6 times below the number of cells.

Then, the I centroids are defined by using the observation matrix. Because we decide to have a uniform discretization of centroids, we only need to find the centroid for both edges which are noted $[a, b]$. The values of a and b are defined such that $\Pr(X^{(m)} < a) < 10^{-5}$ and $\Pr(X^{(m)} > b) < 10^{-5}$ (see the figure below).

```
t = supportDist(X, indVar, I);
```



Choice of subset of triplets and computation of 3D marginals

The principle of PCTF3D is to couple the information of 3D marginals in the form of 3D histograms. But instead of considering all possible triplets, PCTF3D only uses a subset of 3D histograms to obtain an estimate of the PDF. In this section, the subset of triplets of variables also named a coupling and denoted by \mathcal{T} is going to be defined. There are several ways to define a coupling. We call a coupling strategy the different strategies that permits to define such a coupling. For more information on that matter, please refer to the reference mentioned in the introduction.

We propose the following strategies in the CTFlowHD package:

Coupling Strategy	strStrategy	Number of triplets T
'+1'	'+1'	M
'+2'	'+2'	$\lfloor \frac{M}{2} \rfloor$
All possible triplets	'full'	$\binom{M}{3}$
Random selection	'rng'	Set by user
Balanced selection	'bal'	Set by user

```
strStrategy = 'full';
[calT,T] = createTriplets(M,indVar,strStrategy);
```

After the strategy and the coupling \mathcal{T} is defined, 3D marginals are computed. Therefore, the resulted structure `dataMarg` contains T order-3 arrays of size $I \times I \times I$.

```
tic, dataMarg = computeMarginals(X,indVar,t,I,calT); timeMarg = toc;
```

```
Computing the marginals... (4 over 4) Done.
```

PCTF3D optimization algorithm

PCTF3D couples the set of 3D histograms $\{\mathcal{H}^{(jkl)} \mid \{j,k,l\} \in \mathcal{T}\}$ to obtain an estimation of the full histogram noted \mathcal{H} . However, PCTF3D does not directly output \mathcal{H} . Indeed, because of the curse of dimensionality, \mathcal{H} is defined by I^M values and this value can be prohibitive for experiments with large values of M . To circumvent the curse of dimensionality, the PDF noted $p(X^{(1)}, \dots, X^{(M)})$ is modelled with a naive Bayes model (NBM) that introduces latent variable L . This variable is taking R discrete values. The principle of the NBM is that the PDF is independent conditionnally with L :

$$p(X^{(1)}, \dots, X^{(M)}) = \sum_{r=1}^R \Pr(L = r) \prod_{m=1}^M p(X^{(m)} | L = r).$$

When discretized, this model can be represented by M factor matrices stored in the variable `y` and a loading vector `lambda`. The m -th factor matrix (`y{m}`) is of size $I \times R$ and stores the estimation of the set of densities $\{p(X^{(m)} \mid L = r) \mid r \in \{1, \dots, R\}\}$ in the form of 1-dimensional histograms. Concerning the loading vector `lambda`, it corresponds to the set of probabilities of the different latent variable states $\{\Pr(L = r) \mid r \in \{1, \dots, R\}\}$.

Therefore, the number of NBM components must be fixed by the user.

```
R = 30;
```

Along with the parameter `R`, the PCTF3D optimization algorithm has optional parameters:

- `y0` : a set of factor matrices that will be used as PCTF3D initialization (random by default),
- `lambda0` : an initialization for the loading vector (random by default),
- `T1` : the maximal number of outer iterations (default value is 1,000),
- `T2` : the maximal number of inner iterations (default value is 20),
- `eps` : stopping criterion for inner iterations (see the reference for more details) (default value is 10^{-9}),
- `tolNewY` : stopping criterion for outer iterations (if factors aren't evolving) (default value is 10^{-6}).

```
tic, [y,lambda,cost] = PCTF3D(dataMarg,R,'T1',2000,'T2',50);
timeA0ADMM = toc;
timePCTF3D = timeA0ADMM+timeMarg;
fprintf("\nComputation time A0-ADMM : %fs\n",timeA0ADMM)
```

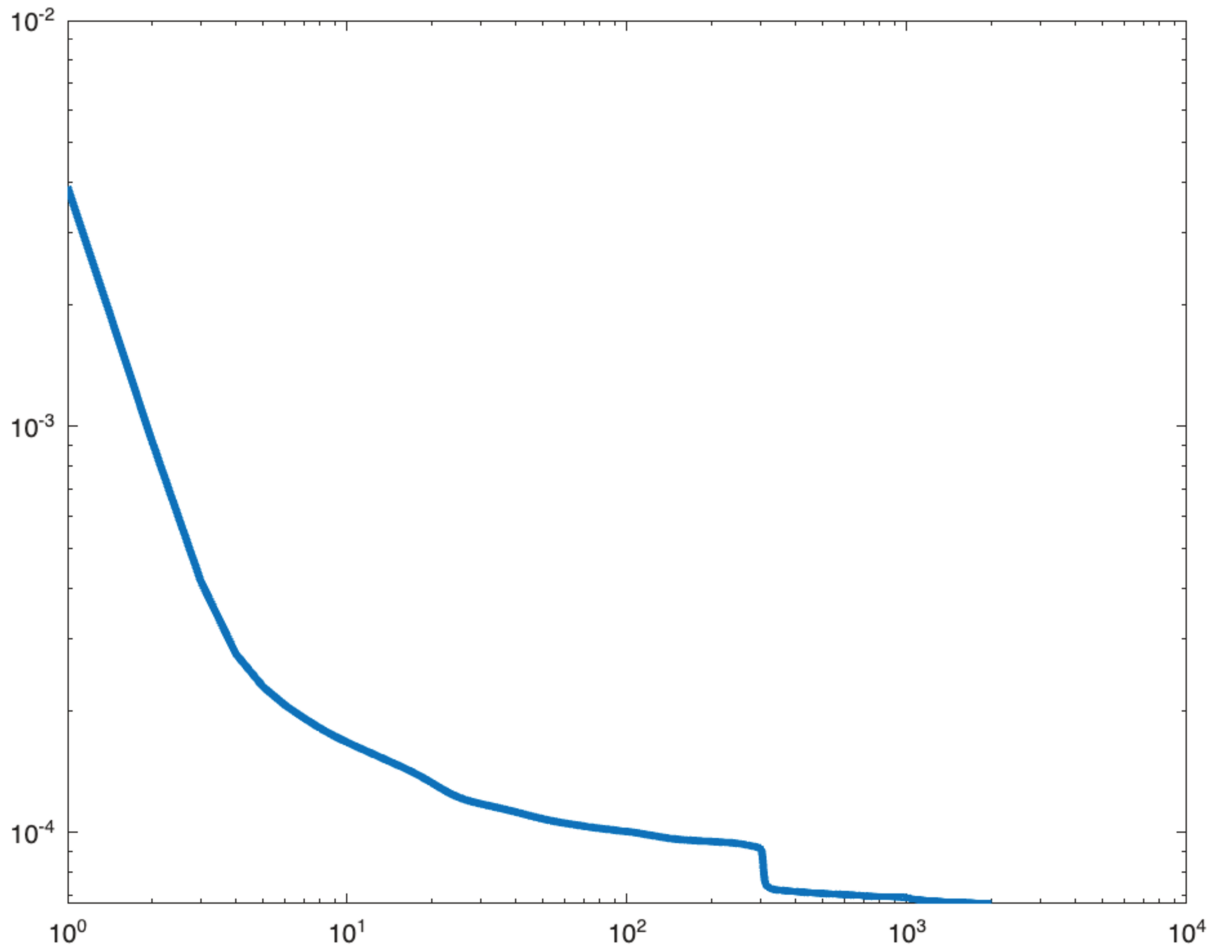
Computation time A0-ADMM : 38.120439s

```
fprintf("\nComputation time PCTF3D : %fs\n",timePCTF3D)
```

Computation time PCTF3D : 38.181485s

The PCTF3D algorithm also outputs the cost function value after the initialization and after each outer iteration.

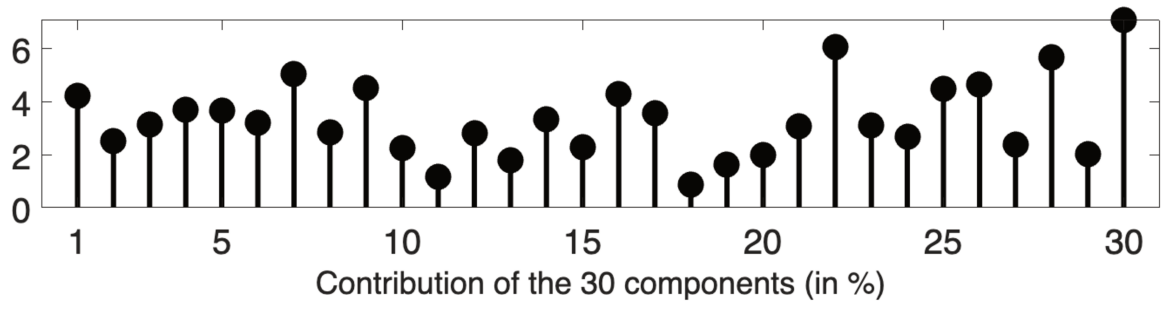
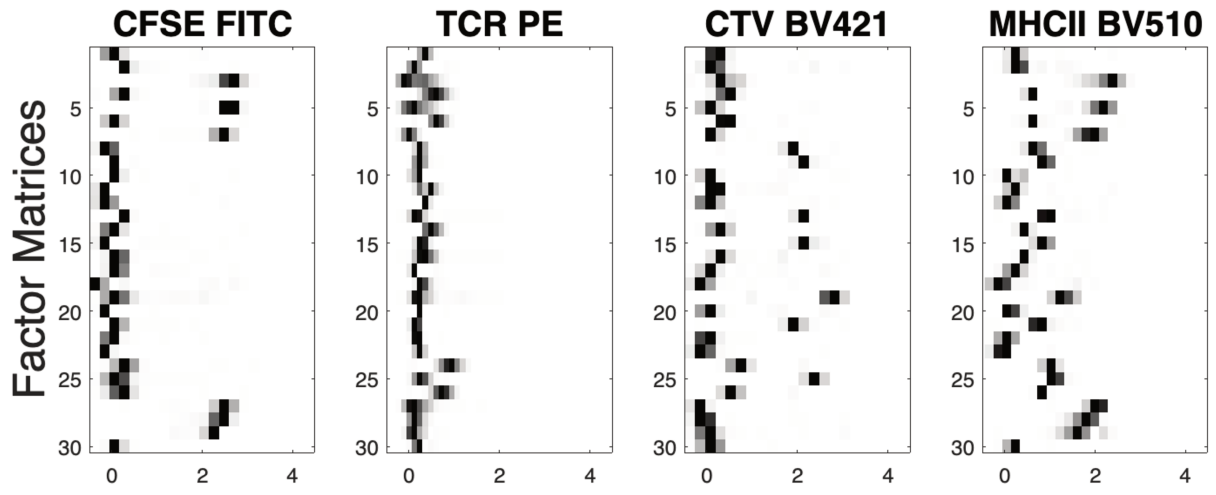
```
figure  
loglog(cost,'LineWidth',3)
```



Visualization of PCFT3D results

To visualize the NBM obtained by PCTF3D, we propose various visualization tools in the CTFlowHD package. In this demonstration, we use a raw visualization of NBM components. For complex visualization tools that provide clustering of NBM components, please refer to the script CTFlowHD_Visualizations of the CTFlowHD package.

```
[y,lambda] = plot_NBMblack(y,lambda,t,indVar,strLabel);
```

Bibliography

- [1] Evrim Acar, Daniel M. Dunlavy, Tamara G. Kolda, and Morten Mørup. Scalable tensor factorizations for incomplete data. *Chemometrics and Intelligent Laboratory Systems*, 106(1):41–56, 2011. Publisher: Elsevier. (page 90).
- [2] Aysun Adan, Günel Alizada, Yağmur Kiraz, Yusuf Baran, and Ayten Nalbant. Flow cytometry: basic principles and applications. *Critical reviews in biotechnology*, 37(2):163–176, 2017. Publisher: Taylor & Francis. (pages vii, 14).
- [3] Nima Aghaeepour, Greg Finak, FlowCAP Consortium, Dream Consortium, Holger Hoos, Tim R. Mosmann, Ryan Brinkman, Raphael Gottardo, and Richard H. Scheuermann. Critical assessment of automated flow cytometry data analysis techniques. *Nature methods*, 10(3):228–238, 2013. Publisher: Nature Publishing Group US New York. (pages 27, 29).
- [4] El-ad David Amir, Kara L. Davis, Michelle D. Tadmor, Erin F. Simonds, Jacob H. Levine, Sean C. Bendall, Daniel K. Shenfeld, Smita Krishnaswamy, Garry P. Nolan, and Dana Pe'er. viSNE enables visualization of high dimensional single-cell data and reveals phenotypic heterogeneity of leukemia. *Nature biotechnology*, 31(6):545–552, 2013. Publisher: Nature Publishing Group US New York. (pages 28, 29, 136).
- [5] Claus A Andersson and Rasmus Bro. The N-way Toolbox for MATLAB. *Chemometrics and Intelligent Laboratory Systems*, 52(1):1–4, August 2000. Number: 1. (pages 8, 127).
- [6] Boster antibody and ELISA experts. Flow Cytometry Basics, FACS Principle. How Does Flow Cytometry Work, <https://www.bosterbio.com/protocol-and-troubleshooting/flow-cytometry-principle>. (page 11).

- [7] C. Bruce Bagwell. Hyperlog—A flexible log-like transform for negative, zero, and positive valued data. *Cytometry Part A*, 64A(1):34–42, March 2005. (page 23).
- [8] Barthel Barlogie, Martin N. Raber, Johannes Schumann, Tod S. Johnson, Benjamin Drewinko, Douglas E. Swartzendruber, Wolfgang Göhde, Michael Andreeff, and Emil J. Freireich. Flow cytometry in clinical cancer research. *Cancer research*, 43(9):3982–3997, 1983. Publisher: AACR. (page 12).
- [9] David Barnett and John T. Reilly. Quality control in flow cytometry. *Flow Cytometry: Principles and Applications*, pages 113–131, 2007. Publisher: Springer. (pages 13, 20, 20).
- [10] Etienne Becht, Leland McInnes, John Healy, Charles-Antoine Dutertre, Immanuel WH Kwok, Lai Guan Ng, Florent Ginhoux, and Evan W. Newell. Dimensionality reduction for visualizing single-cell data using UMAP. *Nature biotechnology*, 37(1):38–44, 2019. Publisher: Nature Publishing Group US New York. (page 29).
- [11] Richard Bellman and Robert Kalaba. On adaptive control processes. *IRE Transactions on Automatic Control*, 4(2):1–9, 1959. Publisher: IEEE. (pages 30, 56).
- [12] Grigoriy Blekherman and Zach Teitler. On maximum, typical and generic ranks. *Mathematische Annalen*, 362(3-4):1021–1031, August 2015. (page 41).
- [13] V. Boonyasombat. Degree sequences of connected hypergraphs and hypertrees. In *Graph Theory Singapore 1983*, volume 1073, pages 236–247. Springer Berlin Heidelberg, Berlin, Heidelberg, 1984. Series Title: Lecture Notes in Mathematics. (page 80).
- [14] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011. Publisher: Now Publishers, Inc. (page 74).
- [15] Paul Breiding, Fulvio Gesmundo, Mateusz Michałek, and Nick Vannieuwenhoven. Algebraic compressed sensing. *arXiv preprint arXiv:2108.13208*, 2021. (pages 41, 42, 42, 104).
- [16] Robert V. Bruggner, Bernd Bodenmiller, David L. Dill, Robert J. Tibshirani, and Garry P. Nolan. Automated identification of stratifying signatures in cellular subpopulations. *Proceedings of the National Academy of Sciences*, 111(26), July 2014. (page 27).

-
- [17] Dario Campana and Elaine Coustan-Smith. Detection of minimal residual disease in acute leukemia by flow cytometry. *Cytometry: The Journal of the International Society for Analytical Cytology*, 38(4):139–152, 1999. Publisher: Wiley Online Library. (page [12](#)).
- [18] A. V. Carrano, J. W. Gray, R. G. Langlois, K. J. Burkhart-Schultz, and M. A. Van Dilla. Measurement and purification of human chromosomes by flow cytometry and sorting. *Proceedings of the National Academy of Sciences*, 76(3):1382–1384, 1979. Publisher: National Acad Sciences. (page [12](#)).
- [19] Rhodri Ceredig, Antonius G. Rolink, and Geoffrey Brown. Models of haematopoiesis: seeing the wood for the trees. *Nature Reviews Immunology*, 9(4):293–300, 2009. Publisher: Nature Publishing Group UK London. (page [141](#)).
- [20] Bruce A. Chabner. General principles of cancer chemotherapy. In *Goodman and Gilman's The Pharmacological Basis of Therapeutics*, pages 1665–1770. Laurence L. Brunton, 12 edition, 2011. (page [141](#)).
- [21] Pratip K. Chattopadhyay, Todd M. Gierahn, Mario Roederer, and J. Christopher Love. Single-cell technologies for monitoring immune systems. *Nature immunology*, 15(2):128–135, 2014. Publisher: Nature Publishing Group US New York. (page [12](#)).
- [22] Pratip K. Chattopadhyay and Mario Roederer. Good cell, bad cell: Flow cytometry reveals T-cell subsets important in HIV disease. *Cytometry part A*, 77(7):614–622, 2010. Publisher: Wiley Online Library. (page [12](#)).
- [23] Joseph K. Chege, Mikus J. Grasis, Alla Manina, Arie Yeredor, and Martin Haardt. Efficient Probability Mass Function Estimation from Partially Observed Data. In *2022 56th Asilomar Conference on Signals, Systems, and Computers*, pages 256–262. IEEE, 2022. (pages [64](#), [94](#), [150](#)).
- [24] Sindhu Cherian, Ben D. Hedley, and Michael Keeney. Common flow cytometry pitfalls in diagnostic hematopathology. *Cytometry Part B: Clinical Cytometry*, 96(6):449–463, 2019. Publisher: Wiley Online Library. (page [22](#)).
- [25] Melissa Cheung, Jonathan J. Campbell, Liam Whitby, Robert J. Thomas, Julian Braybrook, and Jon Petzing. Current trends in flow cytometry automated data analysis software. *Cytometry Part A*, 99(10):1007–1021, October 2021. (pages [24](#), [28](#), [28](#), [136](#)).

- [26] Luca Chiantini and Giorgio Ottaviani. On generic identifiability of 3-tensors of small rank. *SIAM Journal on Matrix Analysis and Applications*, 33(3):1018–1037, 2012. Publisher: SIAM. (page [40](#), [40](#)).
- [27] Luca Chiantini, Giorgio Ottaviani, and Nick Vannieuwenhoven. An Algorithm For Generic and Low-Rank Specific Identifiability of Complex Tensors. *SIAM Journal on Matrix Analysis and Applications*, 35(4):1265–1287, January 2014. (pages [41](#), [41](#), [41](#), [104](#), [119](#)).
- [28] William S. Cleveland. *Visualizing data*. Hobart press, 1993. (page [50](#)).
- [29] Pierre Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):287–314, 1994. Publisher: Elsevier. (page [39](#)).
- [30] Pierre Comon. Tensors: a Brief Introduction. *Signal Processing Magazine, IEEE*, 31:44–53, May 2014. (page [35](#)).
- [31] Pierre Comon, Jos MF Ten Berge, Lieven De Lathauwer, and Josephine Castaing. Generic and typical ranks of multi-way arrays. *Linear Algebra and its Applications*, 430(11-12):2997–3007, 2009. Publisher: Elsevier. (pages [40](#), [103](#)).
- [32] Elaine Coustan-Smith, Jose Sancho, Frederick G. Behm, Michael L. Hancock, Bassem I. Razzouk, Raul C. Ribeiro, Gaston K. Rivera, Jeffrey E. Rubnitz, John T. Sandlund, and Ching-Hon Pui. Prognostic importance of measuring early clearance of leukemic cells by flow cytometry in childhood acute lymphoblastic leukemia. *Blood, The Journal of the American Society of Hematology*, 100(1):52–58, 2002. Publisher: American Society of Hematology Washington, DC. (page [12](#)).
- [33] Madeleine Cule, Richard Samworth, and Michael Stewart. Maximum likelihood estimation of a multi-dimensional log-concave density. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 72(5):545–607, 2010. (page [51](#)).
- [34] Maud d’Aveni, Anne B. Notarantonio, Allan Bertrand, Laura Boulangé, Cécile Pochon, and Marie T. Rubio. Myeloid-derived suppressor cells in the context of allogeneic hematopoietic stem cell transplantation. *Frontiers in immunology*, 11:989, 2020. Publisher: Frontiers Media SA. (page [146](#)).

-
- [35] Lieven De Lathauwer. Decompositions of a higher-order tensor in block terms—Part II: Definitions and uniqueness. *SIAM Journal on Matrix Analysis and Applications*, 30(3):1033–1066, 2008. ISBN: 0895-4798 Publisher: SIAM. (page [38](#)).
- [36] Arthur P. Dempster, Nan M. Laird, and Donald B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the royal statistical society: series B (methodological)*, 39(1):1–22, 1977. Publisher: Wiley Online Library. (page [49](#)).
- [37] Mathias Drton and Marloes H. Maathuis. Structure learning in graphical modeling. *Annual Review of Statistics and Its Application*, 4:365–393, 2017. Publisher: Annual Reviews. (page [57](#)).
- [38] John Duchi, Shai Shalev-Shwartz, Yoram Singer, and Tushar Chandra. Efficient projections onto the ℓ_1 -ball for learning in high dimensions. *Proceedings of the 25th International Conference on Machine Learning*, pages 272–279, July 2008. (pages [75](#), [76](#)).
- [39] Richard O. Duda and Peter E. Hart. *Pattern classification and scene analysis*, volume 3. Wiley New York, 1973. (pages [28](#), [49](#)).
- [40] B. Ellis, R. Gentleman, F. Hahne, N. Le Meur, Deepayan Sarkar, and M. Jiang. flowViz: Visualization for flow cytometry, 2023. (page [127](#)).
- [41] B. Ellis, P. Haaland, F. Hahne, N. Le Meur, N. Gopalakrishnan, J. Spidlen, M. Jiang, and G. Finak. flowCore: Basic structures for flow cytometry data. *R package version*, 1(0), 2019. (page [127](#)).
- [42] Greg Finak and Juan Manuel-Perez. Package ‘flowTrans’. *bioViews ImmunoOncology, FlowCytometry*, 2010. (pages [23](#), [127](#)).
- [43] Philippe Flores, Guillaume Harlé, Anne-Béatrice Notarantonio, Konstantin Usevich, Maud d’Aveni, Stéphanie Grandemange, Marie-Thérèse Rubio, and David Brie. Coupled tensor factorization for flow cytometry data analysis. In *2022 IEEE 32nd International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6. IEEE, 2022. (pages [8](#), [72](#), [72](#)).
- [44] Philippe Flores, Guillaume Harlé, Konstantin Usevich, Stéphanie Grandemange, and David Brie. Factorisation couplée de tenseurs pour l’analyse de données de cytométrie en

- flux. In *XXVIIIème Colloque francophone de traitement du signal et des images, GRETSI 2022*, Nancy, France, 2022. (pages 8, 72).
- [45] Philippe Flores, Konstantin Usevich, and David Brie. Identifiabilité de modèles tensoriels couplés pour l'estimation de loi de probabilité discrète. In *XXIXème Colloque francophone de traitement du signal et des images, GRETSI 2023*, Grenoble, France, 2023. (page 8).
- [46] Andrea Frosini, Christophe Picouleau, and Simone Rinaldi. On the degree sequences of uniform hypergraphs. In *Discrete Geometry for Computer Imagery: 17th IAPR International Conference, DGCI 2013, Seville, Spain, March 20-22, 2013. Proceedings 17*, pages 300–310. Springer, 2013. (pages 78, 87, 89).
- [47] Stephan Fuhrmann, Mathias Streitz, and Florian Kern. How flow cytometry is changing the study of TB immunology and clinical diagnosis. *Cytometry Part A*, 73(11):1100–1106, 2008. Publisher: Wiley Online Library. (page 12).
- [48] Edgard N. Gilbert and John Riordan. Symmetry types of periodic sequences. *Illinois Journal of Mathematics*, 5(4):657–665, 1961. Publisher: Duke University Press. (page 89).
- [49] John C. Gower and Gavin JS Ross. Minimum spanning trees and single linkage cluster analysis. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 18(1):54–64, 1969. Publisher: Wiley Online Library. (pages 28, 137).
- [50] Burkhard Greve, Reinhard Kelsch, Kristina Spaniol, Hans Theodor Eich, and Martin Götte. Flow cytometry in cancer stem cell analysis and separation. *Cytometry Part A*, 81(4):284–293, 2012. Publisher: Wiley Online Library. (page 12).
- [51] Mats Gyllenberg, Timo Koski, Edwin Reilink, and Martin Verlaan. Non-uniqueness in probabilistic numerical identification of bacteria. *Journal of Applied Probability*, 31(2):542–548, 1994. Publisher: Cambridge University Press. (page 49).
- [52] W. Göhde and W. Dittrich. Simultaneous impulse fluorometry of DNA and protein content of tumour cells. *Fresenius' Zeitschrift für analytische Chemie*, 252:328–330, 1970. Publisher: Springer. (page 12).

-
- [53] John A. Hartigan and Manchek A. Wong. Algorithm AS 136: A k-means clustering algorithm. *Journal of the royal statistical society. series c (applied statistics)*, 28(1):100–108, 1979. ISBN: 0035-9254 Publisher: JSTOR. (page 29).
- [54] Frank L. Hitchcock. The expression of a tensor or a polyadic as a sum of products. *Journal of Mathematics and Physics*, 6(1-4):164–189, 1927. Publisher: Wiley Online Library. (page 37).
- [55] Henry R. Hulett, William A. Bonner, Janet Barrett, and Leonard A. Herzenberg. Cell sorting: automated separation of mammalian cells as a function of intracellular fluorescence. *Science*, 166(3906):747–749, 1969. Publisher: American Association for the Advancement of Science. (page 11).
- [56] Mariya Ishteva. Tensors and Latent Variable Models. In Emmanuel Vincent, Arie Yeredor, Zbyněk Koldovský, and Petr Tichavský, editors, *Latent Variable Analysis and Signal Separation*, volume 9237, pages 49–55. Springer International Publishing, Cham, 2015. Series Title: Lecture Notes in Computer Science. (pages 58, 59).
- [57] Michael I. Jordan. Graphical Models. *Statistical Science*, 19(1):140 – 155, 2004. (pages 57, 58, 58).
- [58] K. Huang, N. D. Sidiropoulos, and A. P. Liavas. A Flexible and Efficient Algorithmic Framework for Constrained Matrix and Tensor Factorization. *IEEE Transactions on Signal Processing*, 64(19):5052–5065, October 2016. Number: 19. (pages 64, 73, 74).
- [59] Tamara Kolda and Brett Bader. Tensor Decompositions and Applications. *SIAM Review*, 51:455–500, August 2009. (pages 35, 38).
- [60] Tamara Gibson Kolda. Multilinear operators for higher-order decompositions. *United States*, 4 2006. (page 36).
- [61] Joseph B. Kruskal. Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics. *Linear algebra and its applications*, 18(2):95–138, 1977. Publisher: Elsevier. (pages 35, 39, 39).
- [62] Hugo Larochelle and Iain Murray. The neural autoregressive distribution estimator. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 29–37. JMLR Workshop and Conference Proceedings, 2011. (page 51).

- [63] Steffen L. Lauritzen. *Graphical models*, volume 17. Clarendon Press, 1996. (page 58, 58).
- [64] Michael Levine, David R. Hunter, and Didier Chauveau. Maximum smoothed likelihood for multivariate mixtures. *Biometrika*, 98(2):403–416, 2011. Publisher: Oxford University Press. (page 60).
- [65] Marshall A. Lichtman, Kenneth Kaushansky, Thomas J. Kipps, Josef T. Prchal, and Marcel M. Levi. *Williams Manual of Hematology*. McGraw-Hill, 2011. (page 140).
- [66] Kenneth Lo, Ryan Remy Brinkman, and Raphael Gottardo. Automated gating of flow cytometry data via robust model-based clustering. *Cytometry Part A*, 73A(4):321–332, April 2008. (page 127).
- [67] Charles F. Van Loan. The ubiquitous Kronecker product. *Numerical Analysis 2000. Vol. III: Linear Algebra*, 123(1):85–100, November 2000. Number: 1. (page 34).
- [68] Roger C. Lyndon. On Burnside’s problem. *Transactions of the American Mathematical Society*, 77(2):202–215, 1954. (page 87).
- [69] David Lähnemann, Johannes Köster, Ewa Szczurek, Davis J. McCarthy, Stephanie C. Hicks, Mark D. Robinson, Catalina A. Vallejos, Kieran R. Campbell, Niko Beerenwinkel, Ahmed Mahfouz, Luca Pinello, Pavel Skums, Alexandros Stamatakis, Camille Stephan-Otto Attolini, Samuel Aparicio, Jasmijn Baaijens, Marleen Balvert, Buys De Barbanson, Antonio Cappuccio, Giacomo Corleone, Bas E. Dutilh, Maria Florescu, Victor Guryev, Rens Holmer, Katharina Jahn, Thamar Jessurun Lobo, Emma M. Keizer, Indu Khatri, Szymon M. Kielbasa, Jan O. Korb, Alexey M. Kozlov, Tzu-Hao Kuo, Boudewijn P.F. Lelieveldt, Ion I. Mandoiu, John C. Marioni, Tobias Marschall, Felix Mölder, Amir Niknejad, Alicja Rączkowska, Marcel Reinders, Jeroen De Ridder, Antoine-Emmanuel Saliba, Antonios Somarakis, Oliver Stegle, Fabian J. Theis, Huan Yang, Alex Zelikovsky, Alice C. McHardy, Benjamin J. Raphael, Sohrab P. Shah, and Alexander Schönhuth. Eleven grand challenges in single-cell data science. *Genome Biology*, 21(1):31, February 2020. (pages 29, 29, 30).
- [70] James MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA, 1967. Issue: 14. (page 136).

-
- [71] Jean-Michel Marin, Kerrie Mengersen, and Christian P. Robert. Bayesian modelling and inference on mixtures of distributions. *Handbook of statistics*, 25:459–507, 2005. Publisher: Elsevier. (page [50](#)).
- [72] Katherine M. McKinnon. Flow cytometry: an overview. *Current protocols in immunology*, 120(1):5–1, 2018. Publisher: Wiley Online Library. (page [26](#)).
- [73] Gianni Monaco, Hao Chen, Michael Poidinger, Jinmiao Chen, João Pedro de Magalhães, and Anis Larbi. flowAI: automatic and interactive anomaly discerning tools for flow cytometry data. *Bioinformatics*, 32(16):2473–2480, 2016. Publisher: Oxford University Press. (pages [23](#), [127](#)).
- [74] N. Kargas, N. D. Sidiropoulos, and X. Fu. Tensors, Learning, and “Kolmogorov Extension” for Finite-Alphabet Random Vectors. *IEEE Transactions on Signal Processing*, 66(18):4854–4868, September 2018. Number: 18. (pages [40](#), [48](#), [64](#), [64](#), [64](#), [64](#), [66](#), [66](#), [66](#), [66](#), [67](#), [67](#), [67](#), [67](#), [69](#), [69](#), [81](#), [100](#), [105](#), [105](#), [105](#), [111](#), [113](#), [113](#), [121](#), [122](#), [149](#), [150](#)).
- [75] Robert D. Nowak. Distributed EM algorithms for density estimation and clustering in sensor networks. *IEEE transactions on signal processing*, 51(8):2245–2253, 2003. Publisher: IEEE. (page [49](#)).
- [76] Krzysztof Nowicki and Tom A. B. Snijders. Estimation and prediction for stochastic blockstructures. *Journal of the American statistical association*, 96(455):1077–1087, 2001. Publisher: Taylor & Francis. (page [49](#)).
- [77] Christa E. Osuna, So-Yon Lim, Jessica L. Kublin, Richard Apps, Elsa Chen, Talia M. Mota, Szu-Han Huang, Yanqin Ren, Nathaniel D. Bachtel, and Athe M. Tsibris. Evidence that CD32a does not mark the HIV-1 latent reservoir. *Nature*, 561(7723):E20–E28, 2018. Publisher: Nature Publishing Group UK London. (page [12](#)).
- [78] W. Roy Overton. Modified histogram subtraction technique for analysis of flow cytometry data. *Cytometry: The Journal of the International Society for Analytical Cytology*, 9(6):619–626, 1988. Publisher: Wiley Online Library. (page [24](#)).
- [79] Bokyoung Park, Keon Hee Yoo, and Changsung Kim. Hematopoietic stem cell expansion and generation: the ways to make a breakthrough. *Blood research*, 50(4):194, 2015. Publisher: Korean Society of Hematology. (page [141](#)).

- [80] David R. Parks, Mario Roederer, and Wayne A. Moore. A new “Logicle” display method avoids deceptive effects of logarithmic scaling for low signals and compensated data. *Cytometry Part A: The Journal of the International Society for Analytical Cytology*, 69(6):541–551, 2006. Publisher: Wiley Online Library. (page [23](#)).
- [81] Karl Pearson. III. Contributions to the mathematical theory of evolution. *Proceedings of the Royal Society of London*, 54(326-330):329–333, 1894. Publisher: The Royal Society London. (page [49](#)).
- [82] Stephen P. Perfetto, Pratip K. Chattopadhyay, and Mario Roederer. Seventeen-colour flow cytometry: unravelling the immune system. *Nature Reviews Immunology*, 4(8):648–655, 2004. Publisher: Nature Publishing Group UK London. (pages [12](#), [27](#), [30](#)).
- [83] Yang Qi, Pierre Comon, and Lek-Heng Lim. Semialgebraic geometry of nonnegative tensor rank. *SIAM Journal on Matrix Analysis and Applications*, 37(4):1556–1580, 2016. Publisher: SIAM. (pages [38](#), [41](#), [41](#), [41](#), [41](#), [43](#), [43](#), [100](#), [118](#), [118](#), [118](#), [119](#)).
- [84] Peng Qiu, Erin F. Simonds, Sean C. Bendall, Kenneth D. Gibbs Jr, Robert V. Bruggner, Michael D. Linderman, Karen Sachs, Garry P. Nolan, and Sylvia K. Plevritis. Extracting a cellular hierarchy from high-dimensional cytometry data with SPADE. *Nature biotechnology*, 29(10):886–891, 2011. Publisher: Nature Publishing Group US New York. (pages [28](#), [137](#)).
- [85] Jonathan Richard, Maxime Veillette, Laurie-Anne Batraverse, Mathieu Coutu, Jean-Philippe Chapleau, Mattia Bonsignori, Nicole Bernard, Cécile Tremblay, Michel Roger, and Daniel E. Kaufmann. Flow cytometry-based assay to study HIV-1 gp120 specific antibody-dependent cellular cytotoxicity responses. *Journal of virological methods*, 208:107–114, 2014. Publisher: Elsevier. (page [12](#)).
- [86] Elina Robeva and Anna Seigal. Duality of graphical models and tensor networks. *Information and Inference: A Journal of the IMA*, 8(2):273–288, 2019. Publisher: Oxford University Press. (page [58](#)).
- [87] Javier O. Rodríguez, Signed E. Prieto, Catalina Correa, Carlos E. Pérez, Jessica T. Mora, Juan Bravo, Yolanda Soracipa, and Luisa F. Álvarez. Predictions of CD4 lymphocytes’ count in HIV patients from complete blood count. *BMC Medical Physics*, 13(1):1–6, 2013. Publisher: BioMed Central. (page [12](#)).

-
- [88] Mario Roederer. Compensation in flow cytometry. *Current protocols in cytometry*, 22(1):1–14, 2002. Publisher: Wiley Online Library. (pages [18](#), [19](#), [20](#)).
- [89] Mario Roederer, Adam Treister, Wayne Moore, and Leonore A. Herzenberg. Probability binning comparison: a metric for quantitating univariate distribution differences. *Cytometry: The Journal of the International Society for Analytical Cytology*, 45(1):37–46, 2001. Publisher: Wiley Online Library. (page [24](#)).
- [90] Elis A. Rosa, Silvia R. Lanza, Carlos R. Zanetti, and Aguinaldo R. Pinto. Immunophenotyping of Classic Murine Myeloma Cell Lines Used for Monoclonal Antibody Production. *Hybridoma*, 31(1):1–6, February 2012. (pages [vii](#), [22](#)).
- [91] Hayley Rose-Inman and Damon Kuehl. Acute leukemia. *Hematology/Oncology Clinics*, 31(6):1011–1028, 2017. Publisher: Elsevier. (pages [140](#), [141](#)).
- [92] Yvan Saeys, Sofie Van Gassen, and Bart N. Lambrecht. Computational flow cytometry: helping to make sense of high-dimensional immunology data. *Nature Reviews Immunology*, 16(7):449–462, 2016. Publisher: Nature Publishing Group UK London. (pages [27](#), [27](#), [29](#)).
- [93] Deepayan Sarkar. *Lattice: Multivariate Data Visualization with R*. Springer New York, New York, NY, 2008. (page [127](#)).
- [94] David W. Scott. On optimal and data-based histograms. *Biometrika*, 66(3):605–610, 1979. Publisher: Oxford University Press. (page [51](#), [51](#)).
- [95] David W. Scott. Frequency polygons: theory and application. *Journal of the American Statistical Association*, 80(390):348–354, 1985. Publisher: Taylor & Francis. (page [51](#)).
- [96] David W. Scott. *Multivariate density estimation: theory, practice, and visualization*. John Wiley & Sons, 2015. (pages [55](#), [56](#)).
- [97] David W. Scott, Richard A. Tapia, and James R. Thompson. Kernel density estimation revisited. *Nonlinear Analysis: Theory, Methods & Applications*, 1(4):339–372, January 1977. (page [51](#)).
- [98] George E. Seidel Jr. Sexing mammalian sperm—intertwining of commerce, technology, and biology. *Animal Reproduction Science*, 79(3-4):145–156, 2003. Publisher: Elsevier. (page [12](#)).

- [99] Howard M. Shapiro. *Practical flow cytometry*. John Wiley & Sons, 2005. (pages [10](#), [12](#), [13](#)).
- [100] Nicholas D. Sidiropoulos and Rasmus Bro. On the uniqueness of multilinear decomposition of N-way arrays. *Journal of Chemometrics: A Journal of the Chemometrics Society*, 14(3):229–239, 2000. Publisher: Wiley Online Library. (page [40](#)).
- [101] Bernard W Silverman. *Density estimation for statistics and data analysis*. Routledge, 2018. (page [51](#)).
- [102] Josef Spidlen, Wayne Moore, David Parks, Michael Goldberg, Chris Bray, Pierre Bierre, Peter Gorombey, Bill Hyun, Mark Hubbard, Simon Lange, Ray Lefebvre, Robert Leif, David Novo, Leo Ostruszka, Adam Treister, James Wood, Robert F. Murphy, Mario Roederer, Damir Sudar, Robert Zigon, and Ryan R. Brinkman. Data File Standard for Flow Cytometry, version FCS 3.1. *Cytometry Part A*, 77A(1):97–100, January 2010. (page [16](#)).
- [103] Volker Strassen. Rank and optimal computation of generic tensors. *Linear algebra and its applications*, 52:645–685, 1983. Publisher: Elsevier. (page [41](#)).
- [104] Henko Tadema, Wayel H. Abdulahad, Coen A. Stegeman, Cees GM Kallenberg, and Peter Heeringa. Increased expression of Toll-like receptors by monocytes and natural killer cells in ANCA-associated vasculitis. *PloS one*, 6(9):e24315, 2011. Publisher: Public Library of Science San Francisco, USA. (pages [vii](#), [14](#)).
- [105] Prisca Theunissen, Ester Mejstrikova, Lukasz Sedek, Alita J. van der Sluijs-Gelling, Giuseppe Gaipa, Marius Bartels, Elaine Sobral da Costa, Michaela Kotrová, Michaela Novakova, and Edwin Sonneveld. Standardized flow cytometry for highly sensitive MRD measurements in B-cell acute lymphoblastic leukemia. *Blood, The Journal of the American Society of Hematology*, 129(3):347–357, 2017. Publisher: American Society of Hematology Washington, DC. (page [12](#)).
- [106] Joseph Trotter. Alternatives to Log-Scale Data Display. *Current Protocols in Cytometry*, 42(1), October 2007. (page [23](#)).
- [107] Ledyard R. Tucker. Implications of factor analysis of three-way matrices for measurement of change. *Problems in measuring change*, 15(122-137):3, 1963. Publisher: University of Wisconsin Press Madison. (page [38](#)).

-
- [108] Benigno Uria, Iain Murray, and Hugo Larochelle. RNADE: The real-valued neural autoregressive density-estimator. *Advances in Neural Information Processing Systems*, 26, 2013. (page 51).
- [109] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11), 2008. (pages 28, 136).
- [110] Harry L. Van Trees. *Detection, estimation, and modulation theory, part I: detection, estimation, and linear modulation theory*. John Wiley & Sons, 2004. (page 49).
- [111] Ravi Varadhan and Christophe Roland. Simple and Globally Convergent Methods for Accelerating the Convergence of Any EM Algorithm. *Scandinavian Journal of Statistics*, 35(2):335–353, 2008. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9469.2007.00585.x>. (pages 50, 64).
- [112] E. Verronèse, A. Delgado, Jenny Valladeau-Guilemond, G. Garin, S. Guillemaut, Olivier Tredan, Isabelle Ray-Coquard, Thomas Bachelot, A. N’Kodia, and C. Bardin-Dit-Courageot. Immune cell dysfunctions in breast cancer patients detected through whole blood multi-parametric flow cytometry assay. *Oncoimmunology*, 5(3):e1100791, 2016. Publisher: Taylor & Francis. (page 12).
- [113] Graham Vesey, Joe Narai, Nicholas Ashbolt, Keith Williams, and Duncan Veal. Detection of specific microorganisms in environmental samples using flow cytometry. In *Methods in cell biology*, volume 42, pages 489–522. Elsevier, 1994. (page 12).
- [114] MP Wand. Data-based choice of histogram bin width. *The American Statistician*, 51(1):59–64, 1997. (page 52).
- [115] Brent L. Wood. Acute Myeloid Leukemia Minimal Residual Disease Detection: The Difference from Normal Approach. *Current Protocols in Cytometry*, 93(1):e73, June 2020. (pages vii, 22).
- [116] CF Jeff Wu. On the convergence properties of the EM algorithm. *The Annals of statistics*, pages 95–103, 1983. Publisher: JSTOR. (page 50).
- [117] X. Fu, W. Ma, K. Huang, and N. D. Sidiropoulos. Robust volume minimization-based matrix factorization via alternating optimization. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2534–2538, March 2016.

Journal Abbreviation: 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). (page [67](#)).

- [118] Heddy Zola, Bernadette Swart, Ian Nicholson, Bent Aasted, Armand Bensussan, Laurence Boumsell, Chris Buckley, Georgina Clark, Karel Drbal, and Pablo Engel. CD molecules 2005: human cell differentiation molecules. *Blood*, 106(9):3123–3126, 2005. Publisher: American Society of Hematology. (page [10](#)).

Résumé en français

1 Introduction

1.1 Cytométrie en flux

La cytométrie en flux (CMF) permet de mesurer des propriétés biologiques à partir d'un flux de cellules individuelles. Cette technique est largement utilisée en immunologie dans l'étude de cellules sanguines, par exemple lors du diagnostic de maladies auto-immunes comme les leucémies. Le principe de la CMF est d'acquérir deux types de mesures pour chaque cellule : des paramètres géométriques comme la taille ou la granularité ; ou des mesures de fluorescence permettant de caractériser une cellule par son contenu génétique. Avec ces paramètres mesurés pour chaque cellule, le problème de l'analyse de données de CMF peut se résumer en trois tâches principales :

- Identifier une population de cellules à l'intérieur d'un volume de cellules ;
- Quantifier la proportion de d'une certaine population de cellules d'intérêt ;
- Explorer des données pour caractériser une nouvelle population de cellules.

Ces dix dernières années, la CMF permet de mesurer de plus en plus de paramètres conjointement (jusqu'à 50) et ce pour de plus en plus de cellules (jusqu'à des millions). Cette augmentation a montré les limites de l'analyse manuelle de données de CMF et a fait émerger des méthodes automatiques sous l'égide de la cytométrie computationnelle. Ces méthodes ne permettent pas d'analyser efficacement des jeux de données en grandes dimensions du fait de leur coût computationnel ; et cela autant en termes de nombre de dimensions qu'en termes de nombre de cellules.

1.2 Approches probabiliste et tensorielle

D'un point de vue traitement du signal, le problème d'analyse de données de CMF revient à un problème de clustering et/ou de classification. Dans ce travail, ce problème a été adressé sous l'angle de l'estimation de densités de probabilité et plus particulièrement les histogrammes.

En analyse de données, estimer la densité représente le Graal, c'est pourquoi il a fait l'objet de nombreux travaux. Cependant, estimer une densité est souvent considéré impossible en pratique du fait de la malédiction de la dimension : la complexité d'un problème augmente de manière exponentielle avec son nombre de variables. Dans le cas des histogrammes, il est possible de les estimer jusqu'à environ 3 voire 4 dimensions.

Pour contourner la malédiction de la dimension, nous avons développé une approche d'estimation de densité en grandes dimensions se basant sur 3 grands principes : (i) un modèle dont la complexité reste linéaire avec le nombre de dimensions, (ii) utiliser seulement un sous-ensemble

d’histogrammes de dimension partielle (3D), (iii) un algorithme basé sur les tenseurs avec des garanties théoriques de reproductibilité. Nous verrons que cette approche permet d’obtenir des caractéristiques pouvant être utilisées ultérieurement dans des tâches de clustering/classification.

L’objectif de cette thèse a été de développer un outil d’analyse de données de cytométrie en flux qui permette l’analyse d’un grand nombre de cellules et de paramètres. Cette méthode est basée sur l’estimation de la densité de probabilité conjointe des données par des approches tensorielles couplées.

2 Plan du manuscrit

Dans un premier temps, les trois premiers chapitres de ce document présentent le problème d’analyse de données de CMF ainsi que deux outils utilisés dans ce travail : les tenseurs et l’estimation de densité par histogrammes.

Le premier chapitre introduit les principes de la CMF. De plus, des méthodes représentatives de CMF computationnelle sont présentées ainsi que les défis qui sont adressés dans ce travail. Le second chapitre introduit le contexte tensoriel utilisé dans cette thèse. En particulier, la décomposition canonique polyadique (CPD) est définie et accompagnée de résultats classiques d’unicité. Ensuite, les moins classiques résultats d’unicité de modèles additifs polynomiaux sont présentés. Ces résultats seront utilisés dans un chapitre dédié à l’identifiabilité de modèles tensoriels couplés. Enfin, dans le chapitre 3, le problème d’estimation de densité de probabilité est introduit. Un focus est attiré sur les histogrammes et leur lien avec les tenseurs. Pour vaincre la malédiction de la dimension, nous supposons que la densité suit un Modèle Bayésien Naïf (MBN). L’idée clé est que le MBN est interprétable comme une CPD contraintes. Pour obtenir les paramètres de ce modèle, deux méthodes basées sur les tenseurs sont présentées et discutées.

Dans un second temps, les trois derniers chapitres du manuscrit présentent les principales contributions de cette thèse.

Le chapitre 4 présente une approche de faible complexité comparé à la littérature en utilisant un nombre limité d’histogrammes 3D via un algorithme de factorisation tensorielle couplée. Cette méthode est appelée FTPC3D pour Factorisation Tensorielle Partiellement Couplée de marginales 3D. Le choix des histogrammes couplés est alors crucial. Plusieurs stratégies sont proposées et étudiées sous l’angle des hypergraphes. En particulier, un algorithme de génération de couplage équilibré est proposé, garantissant que les variables sont représentées équitablement en termes d’occurrences. Nous concluons ce chapitre avec des expériences numériques visant à comparer les performances de FTPC3D avec la littérature.

Le chapitre 5 aborde l’étude de l’unicité du modèle tensoriel couplé introduit par FTPC3D par le biais de la recouvrabilité (*il existe un nombre fini de décompositions*) grâce à des outils de géométrie algébrique. Cela consiste en une re-paramétrisation du modèle et de l’étude du rang de son Jacobien. Un algorithme est alors proposé pour rechercher le rang maximal recouvrable fournissant une condition suffisante de recouvrabilité. Lorsqu’il est appliqué à des stratégies de couplage aléatoire et équilibré, la borne de recouvrabilité semble faire apparaître des *cas pathologiques* avec de faibles bornes uniquement dans le cas de couplages aléatoires. Ces cas ne semblent

pas apparaitre pour les couplages équilibrés, montrant l'intérêt pour cette stratégie. En considérant le modèle couplé comme un modèle polynomial additif, ce chapitre fournit aussi une condition générique d'identifiabilité pour un ensemble particulier de couplages. Cette nouvelle condition permet d'étendre les résultats antérieurs d'identifiabilité en matière de modèle tensoriel couplé.

Enfin, le chapitre 6 présente le package Couplage Tensoriel pour la cytométrie en Flux en grandes dimension ou CTFlowHD. Le cœur de ce package est l'algorithme FTFC3D qui, appliqué à des données de CMF, résulte en un modèle de faible complexité interprétable par les experts de CMF : les facteurs de la CPD. Ces facteurs sont alors les entrées de méthodes classiques de clustering/visualisation permettant des interprétations biologiques à bas coût computationnel. La versatilité de CTFlowHD est alors illustrée sur différents jeux de données réels où de multiples exemples de visualisations sont présentés. Dans un exemple biologique, CTFlowHD permet l'identification de populations de cellules sanguines utilisant des données en 8 dimensions.

Dans ce résumé, nous avons décidé de ne pas présenter en détail les deux premiers chapitres car ils ne sont pas essentiels à la compréhension de ce résumé.

3 Histogrammes et estimation de densités de probabilité

3.1 Définition et malédiction de la dimension

Soit X une variable aléatoire prenant des valeurs dans $\mathcal{I} = [a, b[$ et soit \mathbf{X} un vecteur d'observation contenant N réalisations de X . La densité de probabilité est la fonction p_X définie sur \mathcal{I} telle que pour $\iota \subset \mathcal{I}$:

$$\Pr(X \in \iota) = \int_{\iota} p_X(x) dx.$$

Le principe de l'histogramme est de séparer \mathcal{I} en I intervalles disjoints $\{\Delta_i\}_{i=1}^I$ tel que :

$$\Delta_i = \left[a + (b - a) \frac{i-1}{I}, a + (b - a) \frac{i}{I} \right[.$$

Definition 3.1. Vecteur histogramme – Le vecteur histogramme (ou histogramme théorique) est le vecteur \mathbf{h} de taille I défini par :

$$h_i = \Pr(X \in \Delta_i).$$

Definition 3.2. Histogramme empirique – L'histogramme empirique est une estimation de \mathbf{h} défini par (qui définit par construction le vecteur $\tilde{\mathbf{h}}$) :

$$\tilde{h}_i = \frac{1}{N} \text{Card} \{n \in \llbracket 1, N \rrbracket \mid x_n \in \Delta_i\}.$$

Pour que l'histogramme empirique soit un estimateur efficace, le nombre d'intervalles I est choisi en fonction du nombre de réalisations N :

$$I = \left\lfloor 3 \sqrt[3]{N} \right\rfloor.$$

Par analogie, on peut définir les mêmes notions dans le cas de M variables aléatoires $\{X^{(m)}\}_{m=1}^M$ définis sur $\mathcal{I}^{(m)}$ et une matrice d'observation \mathbf{X} de taille $N \times M$.

Definition 3.3. Tenseur histogramme – *Le tenseur histogramme ou (histogramme théorique) est le tenseur d'ordre M \mathcal{H} de taille $I \times \dots \times I$ dont les entrées sont définies par :*

$$h_{i_1 \dots i_M} = \Pr \left(X^{(1)} \in \Delta_{i_1}^{(1)}, \dots, X^{(M)} \in \Delta_{i_M}^{(M)} \right).$$

Definition 3.4. Tenseur histogramme empirique – *Le tenseur histogramme empirique $\tilde{\mathcal{H}}$ est une estimation de \mathcal{H} définies par :*

$$\tilde{h}_{i_1 \dots i_M} = \frac{1}{N} \text{Card} \left\{ n \in \llbracket 1, N \rrbracket \mid \mathbf{x}_{n,:} \in \Delta_{i_1}^{(1)} \times \dots \times \Delta_{i_M}^{(M)} \right\}.$$

De même, I est choisi en fonction de N et M avec la formule :

$$I = \left\lceil 3^{M+2} \sqrt{N} \right\rceil$$

Cette équation montre qu'avec un nombre limité d'observations, on ne peut choisir I qu'avec de très faibles valeurs. Par exemple, pour $M = 10$ variables, il faudrait $n = 7.10^9$ observations si on souhaite avoir $I = 20$ intervalles. C'est la manifestation de la malédiction de la dimension qui exhibe que la complexité d'un problème augmente de manière exponentielle avec son nombre de variables.

3.2 Modèle Bayésien naïf et décomposition tensorielle

Pour contourner la malédiction de la dimension, la densité est supposée suivre un Modèle Bayésien Naïf (MBN) qui introduit une variable latente L discrète à R états. Pour un tenseur histogramme, ce modèle a l'expression suivante :

$$h_{i_1 \dots i_M} = \sum_{r=1}^R \Pr(L = r) \prod_{m=1}^M \Pr \left(X^{(m)} \in \Delta_{i_m}^{(m)} \mid L = r \right).$$

Premièrement, ce modèle a un nombre de paramètres ($R(MI + 1)$) qui évolue linéairement avec M . Deuxièmement, le MBN peut être vu comme une décomposition tensorielle et plus particulièrement une CPD sous contraintes de *simplex* :

$$\mathcal{H} = \sum_{r=1}^R \lambda_r \prod_{m=1}^M \mathbf{a}_r^{(m)}$$

où le poids λ_r représente la probabilité $\Pr(L = r)$ et où le facteur $\mathbf{a}_r^{(m)}$ représente la probabilité conditionnelle $\Pr \left(X^{(m)} \in \Delta_{i_m}^{(m)} \mid L = r \right)$. Le rang de la décomposition R est alors le nombre d'états latents de L .

3.3 FTC3D : Factorisation Tensorielle Couplée de marginales 3D

Dans un article de [Kargas *et al.*, 2018], une méthode qui utilise les histogrammes de marginales d'ordre 3 à défaut de pouvoir estimer l'histogramme global. Ces histogrammes 3D, accessibles pour faibles valeurs de N , sont ensuite couplés par un algorithme AO-ADMM (*Alternating Optimization - Alternating Direction Method of Multipliers*) qui résoud le problème d'optimisation sous contraintes suivant :

$$\begin{aligned} \widehat{\lambda}, \widehat{\mathbf{A}}^{(1)}, \dots, \widehat{\mathbf{A}}^{(M)} &= \underset{\lambda, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)}}{\operatorname{argmin}} \sum_{1 \leq j < k < \ell \leq M} \left\| \widetilde{\mathcal{H}}^{(j k \ell)} - \llbracket \lambda; \mathbf{A}^{(j)}, \mathbf{A}^{(k)}, \mathbf{A}^{(\ell)} \rrbracket \right\|_F^2, \\ \text{t.q. } \lambda &\geq 0, \quad \mathbf{A}^{(1)} \geq 0, \quad \dots, \quad \mathbf{A}^{(M)} \geq 0, \\ \mathbb{1}_R^\top \lambda &= 1, \quad \mathbb{1}_I^\top \mathbf{A}^{(1)} = \mathbb{1}_R^\top, \quad \dots, \quad \mathbb{1}_I^\top \mathbf{A}^{(M)} = \mathbb{1}_R^\top. \end{aligned}$$

La complexité de cette méthode est liée au nombre de marginales 3D, c'est-à-dire $\binom{M}{3}$. Cette valeur peut devenir prohibitive lorsque M augmente.

4 FTPC3D : Factorisation Tensorielle Partiellement Couplée de marginales 3D

4.1 Principe

Si le nombre de marginales est un problème, est-ce possible d'estimer les facteurs du MBN avec seulement un nombre limité de marginales ? C'est le principe de la méthode présentée dans ce chapitre : Factorisation Tensorielle Partiellement Couplée de Marginales 3D ou FTPC3D. FTPC3D résout un problème similaire où seules les marginales considérées dans le sous-ensemble \mathcal{T} sont présentes :

$$\begin{aligned} \widehat{\lambda}, \widehat{\mathbf{A}}^{(1)}, \dots, \widehat{\mathbf{A}}^{(M)} &= \underset{\lambda, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(M)}}{\operatorname{argmin}} \sum_{\{j, k, \ell\} \in \mathcal{T}} \left\| \widetilde{\mathcal{H}}^{(j k \ell)} - \llbracket \lambda; \mathbf{A}^{(j)}, \mathbf{A}^{(k)}, \mathbf{A}^{(\ell)} \rrbracket \right\|_F^2, \\ \text{t.q. } \lambda &\geq 0, \quad \mathbf{A}^{(1)} \geq 0, \quad \dots, \quad \mathbf{A}^{(M)} \geq 0, \\ \mathbb{1}_R^\top \lambda &= 1, \quad \mathbb{1}_I^\top \mathbf{A}^{(1)} = \mathbb{1}_R^\top, \quad \dots, \quad \mathbb{1}_I^\top \mathbf{A}^{(M)} = \mathbb{1}_R^\top. \end{aligned}$$

Un algorithme est alors adapté à cette nouvelle approche, toujours basée sur un AO-ADMM.

4.2 Couplages ou stratégies de couplages

Dans FTPC3D, le choix des marginales dans le sous-ensemble \mathcal{T} , aussi appelé le couplage, est crucial. Choisir les éléments dans \mathcal{T} est équivalent à choisir des triplets de variables. Ceci permet d'introduire \mathcal{T} sous le formalisme des hypergraphes, et plus particulièrement des 3-graphes puisque tous les éléments de \mathcal{T} sont de cardinal 3. Ce formalisme permet d'introduire une condition suffisante pour qu'un couplage soit utilisé dans FTPC3D : l'hypergraphe \mathcal{T} doit être connexe.

Afin de choisir un couplage, nous définissons aussi des stratégies de couplage tels que la stratégie aléatoire où les triplets sont choisis aléatoirement tout en garantissant que \mathcal{T} est connexe.

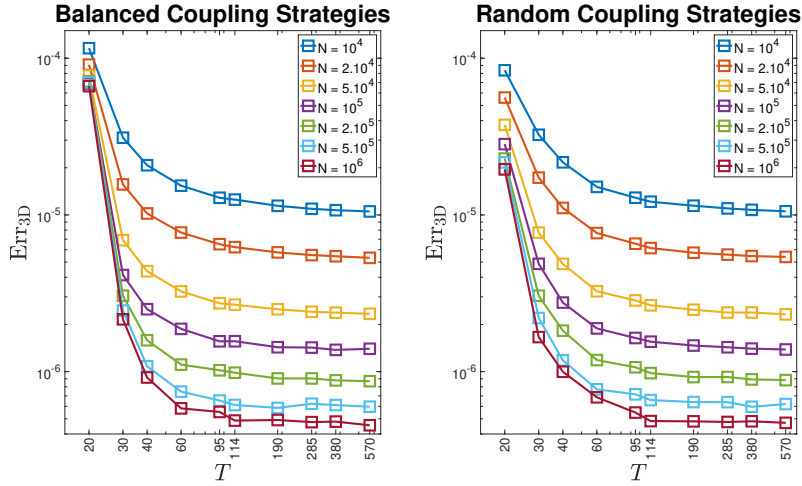


FIGURE 1 – Performance d’estimation en fonction du nombre de triplets T pour les stratégies de couplages aléatoire et équilibrée.

Il a été remarqué que, pour les couplages aléatoires, certaines variables sont représentées moins que d’autres. Pour pallier ce problème, un algorithme basé sur les mots de Lyndon permettent de définir une stratégie de couplage dite équilibrée puisque toutes les variables sont représentées de manière équivalente en termes d’occurrences.

4.3 Évaluation des performances de FTPC3D

Pour comparer FTPC3D, plusieurs expériences numériques ont été réalisées. L’une d’entre elles a été de mesurer les performances d’estimation sur des données synthétiques dans 3 cas : le cas du couplage total donc FTC3D, le couplage aléatoire et le couplage équilibré. Ces deux dernières étant réalisées pour différents nombres de triplets de variable.

Les résultats de cette étude sont fournies en Figure 1. Sur les deux graphes, l’estimation augmente avec T et les deux stratégies ont les mêmes performances. De plus, à partir d’une certaine valeur de T , on atteint un palier de performance. Comme ce palier correspond à la valeur d’erreur pour FTC3D, cela veut dire que l’on peut estimer avec des performances similaires à plus faible coût.

5 Unicité de modèles tensoriels couplés

5.1 Problème

La question de l’unicité est la suivante : à partir d’un couplage \mathcal{T} et d’un ensemble de marginales $\{\mathcal{H}^{(j k \ell)} \mid \{j, k, \ell\} \in \mathcal{T}\}$, pour quelles valeur de rang R peut-on retrouver les facteurs de la CPD ? Cette question est étudiée dans le chapitre 5 sous l’angle de deux notions fondamentales : la recouvrabilité et l’identifiabilité. La recouvrabilité garantit qu’il y a un nombre fini de décomposition tandis que l’identifiabilité permet de garantir qu’il existe une décomposition aux ambiguïtés de permutation et d’échelle près.

5.2 Recouvrabilité et étude de Jacobien

Pour étudier la recouvrabilité, un algorithme est proposé permettant de déterminer, dans une situation donnée (à M , I et \mathcal{T} fixés), le rang maximal recouvrable. Cet algorithme est basé sur le calcul du Jacobien de la paramétrisation de notre modèle couplé. Il permet alors dans une situation donnée de garantir une condition suffisante de recouvrabilité.

Cet algorithme a été appliqué à différentes valeurs de nombre de triplets pour les stratégies aléatoire et équilibrée. Dans le cas aléatoire, la borne est, pour un nombre non négligeable de cas appelés défectueux, plus faible que pour le cas équilibré. Par exemple, ces cas peuvent se présenter si une variable n'est contenue qu'à une occurrence dans le couplage. Pour ces cas défectueux, la structure du Jacobien limite alors le rang de celui-ci et donc la borne de recouvrabilité. Comme le couplage équilibré garantit par définition que toutes les variables sont représentées équitablement, les cas défectueux ne sont pas présents pour cette stratégie, montrant son intérêt.

5.3 Identifiabilité du modèle couplés

Dans le cas général, obtenir une condition suffisante d'identifiabilité dans le cas général est très difficile. Cependant, dans ce chapitre, nous fournissons une condition suffisante d'identifiabilité pour une stratégie de couplage appelée du produit Cartésien, puisqu'elle revient à effectuer le produit entre une tripartition des variables.

La contribution de ce travail a alors été de prendre en considération les contraintes de somme égale à 1 sur les facteurs de la décomposition. Cela a été fait entre autres grâce à la construction d'un ensemble semi-algébrique ouvert contenu dans notre espace de contraintes.

Les résultats obtenus avec cette nouvelle condition d'identifiabilité excède les résultats déjà existant en la matière, notamment ceux de FTC3D. De plus, les garanties d'identifiabilité obtenues pour le modèle de FTFC3D sont valables pour FTC3D.

6 CTFlowHD : un package pour l'analyse de données de cytométrie en flux

6.1 Présentation du package

Nous présentons dans le dernier chapitre un nouveau package disponible en ligne pour la cytométrie appelé *Coupled Tensor factorizations for Flow cytometry in High Dimensions* ou CTFlowHD. Ce package se décompose en deux grandes parties : application de FTFC3D à un jeu de données de CMF ; clustering et visualisations des facteurs de la CPD.

À la différence des méthodes existantes d'analyse de données de CMF, CTFlowHD applique des méthodes de clustering/visualisation sur R termes de rang 1, ce qui représente un coût computationnel très faible comparé à la littérature. De plus, CTFlowHD est très versatile puisqu'à partir du moment où FTFC3D est appliqué et que les facteurs sont obtenus, il est possible d'appliquer en parallèle plusieurs méthodes de clustering/visualisation, permettant une exploration à diverses échelles d'un jeu de données. Parmi les visualisations disponibles, nous proposons un clustering

hiérarchique des facteurs, un affichage de type t-SNE, un affichage selon un arbre couvrant de poids minimal ou encore sous la forme de reconstruction de marginales 2D.

6.2 Application à des données réelles à 8 variables

Dans le cadre d'un jeu de données à 8 dimensions contenant un mélange de cellules hématopoïétiques, CTFlowHD a permis d'identifier plusieurs populations de cellules telles que les Lymphocytes T, les cellules souches, les macrophages et les granulocytes. Pour ces populations, la quantification de leur proportion est en accord avec une expertise manuelle, fournissant de belles perspectives pour la méthode.

Conclusion

Dans cette thèse, nous avons développé une méthode d'analyse de données de cytométrie en flux en voyant ce problème comme une estimation de densité de probabilité en grandes dimensions. Nous avons alors proposé FTPC3D, un algorithme réduisant la complexité comparée à la littérature et permettant de résoudre un peu plus le problème de la malédiction de la dimension. Cette approche a introduit la notion de couplage et d'hypergraphe, qui a été longuement étudiée par le prisme de l'unicité de notre modèle couplé. Enfin, nous avons proposé l'outil CTFlowHD, qui utilise FTPC3D comme base, qui permet une analyse versatile et rapide de données de cytométrie en flux.

Estimation of probability density functions in high-dimensions with low-rank tensor models: application to flow cytometry

Flow cytometry (FCM) is a widely-used technique for blood cell analysis (i.e. in immunology) which measures cell biomarker fluorescence, allowing to identify cell populations. Manual analysis is somewhat subjective and time consuming, especially when the number of parameters increases. This motivated the development of unsupervised methods which, however, hardly scale to high-dimensional datasets.

In this work, this problem is addressed as a joint density estimation which is considered impossible in practice because of the curse of dimensionality. To break the curse of dimensionality, a naive Bayes model (NBM) is considered and viewed as a constrained tensor model whose factors can be estimated through the coupled tensor factorization of 3D marginals (CTF3D).

An algorithm that couples only a subset of 3D marginals is proposed (PCTF3D). To choose the subset of considered marginals, two coupling strategies are explored: randomly chosen triplets and balanced couplings where variables are represented evenly. PCTF3D shows similar performances with CTF3D while reducing the computational burden.

PCTF3D introduced a new constrained coupled model whose uniqueness properties are studied. This results in new recoverability bounds of the coupled model that apply to fully and partially coupled models. Defective cases which only arise for random couplings shows the interest of balanced couplings. Also, a sufficient identifiability condition is proved using the algebraic structure of the model constraints.

Finally, a new computational framework called *CTFlowHD* whose core is based on PCTF3D is proposed. The FCM data is then reduced to a set of NBM factors. *CTFlowHD* also provides several tools to visualize and cluster NBM factors, including clustering/visualization methods widely used in the FCM community. *CTFlowHD* is available online for use.

Keywords: flow cytometry, histograms, naive Bayes model, coupled tensor decompositions, hypergraphs.

Estimation de densités de probabilité en grandes dimensions par modèles tensoriels de rang faibles : application à la cytométrie en flux

La cytométrie en flux (CMF) est une technique d'analyse de cellules sanguines largement répandue (par exemple en immunologie) qui mesure la fluorescence de biomarqueurs cellulaires, permettant d'identifier des populations de cellules. L'analyse manuelle de ces données est subjective et prend du temps, en particulier lorsque le nombre de paramètres augmente. Cela a poussé le développement de méthodes non supervisées qui ne s'adaptent pas encore aux volumes de données disponibles.

Dans ce travail, ce problème est abordé sous la forme d'une estimation jointe de la densité, considérée comme impossible en pratique du fait de la malédiction de la dimension. Pour la contourner, un modèle de Bayésien naïf (MBN) est considéré et vu comme un modèle tensoriel contraint dont les facteurs peuvent être estimés par une factorisation tensorielle couplée de marginales 3D (CTF3D).

Un algorithme ne couplant qu'un sous-ensemble de marginales 3D est proposé (PCTF3D). Pour choisir le sous-ensemble de marginales considérées, deux stratégies de couplage sont explorées : des triplets choisis au hasard et des couplages équilibrés où les variables sont représentées en nombres égaux. PCTF3D a des performances similaires à celles de CTF3D tout en réduisant les temps de calcul.

PCTF3D a introduit un nouveau modèle couplé contraint dont les propriétés d'unicité sont étudiées. Il en résulte de nouvelles limites de recouvrabilité qui s'appliquent aux modèles entièrement et partiellement couplés. Les cas défectueux n'apparaissant que pour les couplages aléatoires montrent l'intérêt des couplages équilibrés. En outre, une condition suffisante d'identifiabilité est prouvée en exploitant la structure algébrique des contraintes du modèle.

Enfin, un nouvel outil se basant sur PCTF3D est proposé et appelé *CTFlowHD*. Les données de CMF sont donc réduites à un ensemble de facteurs d'un MBN. *CTFlowHD* fournit également plusieurs outils pour visualiser et regrouper les facteurs du MBN, y compris des méthodes de visualisation/clustering déjà utilisées dans la communauté de la CMF. *CTFlowHD* est disponible en ligne pour utilisation.

Mots-clés: Cytométrie en flux, histogrammes, modèle Bayésien naïf, décompositions tensorielles couplées, hypergraphes.