



HAL
open science

Stochastic modeling of aggregation and flocculation processes in chemistry

Daniel Paredes Moreno

► **To cite this version:**

Daniel Paredes Moreno. Stochastic modeling of aggregation and flocculation processes in chemistry. Modeling and Simulation. Université Paul Sabatier - Toulouse III, 2017. English. NNT: 2017TOU30368 . tel-02009796

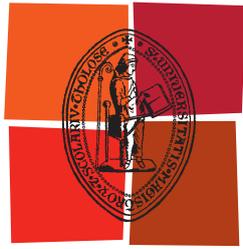
HAL Id: tel-02009796

<https://theses.hal.science/tel-02009796>

Submitted on 6 Feb 2019

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 Toulouse

THÈSE

En vue de l'obtention du

DOCTORAT DE L'UNIVERSITÉ DE TOULOUSE

Délivré par :

Université Toulouse 3 Paul Sabatier (UT3 Paul Sabatier)

Présentée et soutenue par :

Daniel Eduardo Paredes Moreno

Le vendredi 27 octobre 2017

Titre :

Modélisation Stochastique de Processus d'Agrégation en Chimie

ED MITT : Domaine Mathématiques : Mathématiques appliquées

Unité de recherche :

Institut de Mathématiques de Toulouse (IMT)

Directeur(s) de Thèse :

Fabrice GAMBOA - Université Paul Sabatier

Rapporteurs :

Jose Rafael LEON RAMOS - Universidad de la Republica-Uruguay
Jean-François DUPUY, Institut de Recherche Mathématique de Rennes

Autre(s) membre(s) du jury :

Léa COT - INSA-Toulouse
Bernard BERCU - Institut de Mathématiques de Bordeaux
Carole SAUDEJAUD - Laboratoire de Génie Chimique de Toulouse

*To my wife, Bárbara Montoya.
Te amo mucho vida mía.*

Agradecimientos

Sin duda alguna, este trabajo no habría podido realizarse sin la ayuda de muchas personas a las que quisiera expresar unas palabras en este espacio.

En primer lugar, quiero dar gracias a Dios Todopoderoso por habernos permitido alcanzar este proyecto, habernos traído a Toulouse Francia a mi esposa y a mi y a pesar de todas las dificultades encontradas, habernos tenido siempre bajo Su protección en todo momento. Muchas gracias Señor.

Jamas voy a poder agradecer lo suficiente a mi esposa Bárbara Montoya. Tu has sido la luz que Dios puso en mi camino para guiarme en todo momento desde que te conocí. Sin tu amor, dedicación, cuidados, apoyo y sacrificio este proyecto no se habría llevado a cabo jamas. Me inspiras siempre a ser mejor desde tu ejemplo y consejo. Eres mi motivo para seguir adelante. Te dedico este trabajo a ti, mi amor, que estuviste en todo momento a mi lado ayudándome de todas las formas posibles. Tu amor es el agua que necesito para vivir y que me da aliento para continuar siempre. Aunque hubo momentos muy difíciles, ha sido un sueño vivir todo esto a tu lado, tu lo haces mágico. Te amo. Mi corazón es tuyo.

Quiero agradecer a mi tutor, el Profesor Fabrice Gamboa, quien supo guiarme siempre de la mejor manera en este proyecto y quien, ante las dificultades que encontré en este proyecto, me apoyó para solventarlas mucho mas allá de su rol como director. Muchas gracias Fabrice.

Agradezco también a los distinguidos miembros del jurado evaluador: José Leon, Jean-François Dupuy, Léa Cot, Bernerd Bercu y Carole Saudejaud que con sus siempre acertados comentarios mejoraron sustancialmente este trabajo.

A todos los miembros del ENSIACET que trabajaron conmigo en el proyecto CARACAS: Pascal Floquet, Christine Frances, Alain Line y Léa Guérin. Gracias por todos los aportes, ayuda y discusiones llevados a cabo durante todo este proyecto. Siempre fueron de muchísima ayuda.

Quiero agradecer a mi familia: A mis padres Irene y Daniel por siempre apoyarme en todo momento, ser siempre un ejemplo para mi. A mis hermanas Irene y Danireé por siempre tener una palabra de aliento. A mis suegros

Maria del Pilar y Santiago y a mis cuñados Monica, Manuel, Vanesa, Andy y Adriana. Todos ustedes han sido un apoyo muy importante para nosotros. Gracias familia.

A nuestra familia en Francia: Norbert y Mireille Dedieu y toda la familia Dedieu. Gracias por acogernos como su familia, todo el cariño y amistad que nos han brindado es un regalo de Dios. Dios los bendiga siempre.

A tantos amigos que de una u otra manera nos ayudaron a superar tantos obstáculos. La lista es larga y temo a olvidar a alguno, les quiero decir de todo corazón, gracias. Especialmente a Diomy y Familia, Marvic y Yohnatan, Gaby y Franklin y familia, Franklin Camacho, Vincent Flores, Vincent Chin, Tamal, Maikol y Laura, Raphael, Kevin y Sophiane, Mélanie, Malika y Benoit (bureau 206). A l'Eglise de Saint-Agne, ASSOFRAVEN entre muchos otros. Muchas gracias a todos los que nos ayudaron.

A la Universidad de Los Andes, Mérida Venezuela, por darme apoyo institucional para poder venir a realizar este trabajo, aunque no pudieran ayudar financieramente. A la Université Paul Sabatier por darme la oportunidad de vivir este sueño.

Résumé

Modelization Stochastique de processus d'agrégation et floculation in Chimie

Auteur: Daniel Eduardo Paredes Moreno.

Directeur: Fabrice Gamboa.

Date et lieu de la soutenance: Le 27/10/2017 à l'Institut de Mathématiques de Toulouse.

Discipline: Mathématiques appliquées.

Mots-clés : Modélisation Stochastique, Méthode de la quadrature des moments, Extrapolation Minimale Généralisée, Filtre Étendu de Kalman, Estimateur de Moindres Carrés non-linéaire.

Résumé

Nous concentrons notre intérêt sur l'Équation du Bilan de la Population (PBE). Cette équation décrit l'évolution, au fil du temps, des systèmes de particules colloïdales en fonction de sa fonction de densité en nombre (NDF) où des processus d'agrégation et de rupture sont impliqués. La NDF dépend des propriétés physiques et morphologiques des particules formant le système. Nous avons également étudié la représentation du PBE comme une équation différentielle en termes des moments de la NDF. La description au fil du temps des systèmes, à la fois en termes de la NDF et en termes d'un ensemble fini des moments standards de la DNF, est pertinent dans des disciplines comme la simulation de la mécanique des fluides numérique et a des applications comme dans le traitement de l'eau.

Plusieurs recherches ont été destinées à trouver un ensemble de variables pertinentes dans la description de l'évolution du système et de maintenir un degré d'interprétation. Dans (Vlieghe 2014), plusieurs expériences ont été effectuées en utilisant de la Bentonite comme matière, et la granulométrie a été effectuée pour la population de particules initiale ainsi que pour les populations résultant des différentes conditions hydrodynamiques. En utilisant

des techniques exploratoires comme l'analyse en composantes principales, le partitionnement de données et l'analyse discriminante, nous avons étudié la formation de groupes et l'importance relative de ces variables dans la formation des ces groupes. Nous avons utilisé ce schéma d'analyse pour la population initiale de particules ainsi que pour les populations résultantes sous différentes conditions hydrodynamiques.

L'étude d'un ensemble fini de moments standard de la NDF est pertinente pour connaître plusieurs aspects physiques du système de particules. Dans les recherches récentes, il existe de nombreuses méthodes développées pour résoudre ce problème. L'une de ces méthodes est la méthode en quadrature des moments (QMOM) qui utilise une application de l'algorithme Produit-Différence. Nous avons étudié l'Extrapolation Minimale Généralisée (GME) afin de récupérer une mesure discrète non-négative, étant donnée un ensemble fini de ses moments standard. De plus, nous avons étudié l'utilisation de la PBE en fonction des moments de la NDF, et les méthodes QMOM et GME, afin de récupérer l'évolution, d'un ensemble fini de moments standard de la NDF.

La PBE est une équation intégro-différentielle impliquant la NDF ainsi que des noyaux représentant la fréquence d'agrégation et de rupture. Ces noyaux dépendent également d'un vecteur de paramètre. Afin de trouver une approximation numérique de la solution de la PBE, nous avons proposé un schéma de discrétisation. Nous avons utilisé trois cas où la solution analytique est connue (Silva et al. 2011) afin de comparer la solution théorique à l'approximation trouvée avec le schéma de discrétisation. Nous avons comparé l'approximation numérique à la NDF empirique estimée en utilisant les données expérimentales. Nous avons utilisé les noyaux et les paramètres identifiés dans cette recherche.

Il est intéressant d'estimer les paramètres apparaissant dans la modélisation des processus d'agrégation et de rupture impliqués dans la PBE. Nous avons proposé une méthode pour estimer les paramètres impliqués dans ces processus, en utilisant l'approximation numérique trouvée à travers le système de discrétisation, ainsi que le Extended Kalman Filter, car les données expérimentales donnent lieu à la fonction de distribution volumique des particules, mais la NDF n'est pas disponible. La méthode estime itérativement les paramètres à chaque instant du temps, en utilisant un estimateur de Moindres Carrés. Nous avons produit plusieurs simulations utilisant les noyaux identifiés. Nous avons également utilisé des données expérimentales réelles obtenues à partir de microparticules de latex (Gerin 2016) en appliquant notre méthode pour estimer le vecteur de paramètres dans ces cas.

Abstract

Stochastic modeling of aggregation and flocculation processes in Chemistry

Author: Daniel Eduardo Paredes Moreno.

Supervisor: Fabrice Gamboa.

Date and place of the defense: Le 27/10/2017 in the Institut de Mathématiques de Toulouse.

Discipline: Applied Mathematics.

Key words : Stochastic Modelization, Quadrature Method of Moments, Generalized Minimal Extrapolation, Extended Kalman Filter, non-linear Least-Squares Estimator.

Abstract

We center our interest in the Population Balance Equation (PBE). This equation describes the time evolution of systems of colloidal particles in terms of its number density function (NDF) where processes of aggregation and breakage are involved. The NDF depends on physical and morphological properties (size, volume, circularity...) of the particles forming the system. We studied also the representation of the PBE as a differential equation in terms of the moments of the NDF. The description of the time evolution of systems, both in terms of the NDF and in terms of a finite set of standard moments of the NDF, is relevant in disciplines like computational fluid dynamics simulation and has applications like in water treatment.

Several researches have investigated about what are the most important variables in order to better describe the behavior of the system in terms of the NDF through the PBE ? The interest of this is to find a set of physical and morphological variables relevant in the description of the evolution of the system and keeping a degree of interpretability. In (Vlieghe 2014), several experiments were done using Bentonite as material, and the granulometry was done for the initial particle population as well as for the populations

resulting at different hydrodynamic conditions. Using exploratory techniques like principal component analysis, cluster analysis and discriminant analysis, we investigated the formation of groups using the available variables and the relative importance of these variables in the formation of the groups. We used this scheme of analysis for the initial population of particles as well as in the resulting populations under different hydrodynamics conditions.

The study of a finite set of standard moments of the NDF is relevant for knowing several physical aspects of the system of particles. In recent research, there are many methods developed in order to solve this problem. One of these methods is the Quadrature Method of Moments (QMOM) which uses an application of the Product-Difference algorithm. We studied the Generalized Minimal Extrapolation (GME) in order to recover a discrete non-negative measure given a finite set of its standard moments. Also, we studied the use of the PBE in terms of the moments of the NDF, and the QMOM and GME methods, in order to recover the time evolution of a finite set of standard moments of the NDF.

The PBE is an integro-differential equation involving the NDF as well as kernels representing the frequency of aggregation and breakage. Those kernels also depend on a vector of parameter. In order to find a numerical approximation to the solution of the PBE, we proposed an discretization scheme. We used three cases where the analytical solution is known (Silva et al. 2011) in order to compare the theoretical solution to the approximation found with the discretization scheme. We compared the numerical approximation to the empirical NDF estimated using the experimental data from the experiments in (Vlieghe 2014). We used the kernels and parameters identified in this research.

It is of interest to estimate the parameters appearing in the modelisation of the aggregation and breakage processes involved in the PBE. We proposed a method for estimate the parameters involved in those processes, using the numerical approximation found through the discretization scheme, as well as the Extended Kalman Filter, because usually experimental data results in the volume distribution function of the particles, but the NDF is not available. The method estimates iteratively the parameters at each time, using an Least Square Estimator. We produced several simulations using the kernels identified in (Vlieghe 2014) for different initial parameter estimation. We also used real experimental data obtained from latex microparticles (Gerin 2016) applying our method for estimate the parameter vector in these case.

Contents

1	Model Formulation	1
1.1	Introduction	1
1.2	Formulation of the model	4
1.2.1	Introduction	4
1.2.2	Formulation of the model	5
1.3	Moments equations for aggregation and breakage	15
1.4	Preliminary Empirical Data	17
1.4.1	Introduction	17
1.4.2	About the data sets	18
1.4.3	Materials and methods	18
1.4.4	Metodology	21
1.4.5	Results of the univariate description	25
1.4.6	Bentonite Experimental Data	65
1.4.7	Conclusions	75
2	Quadrature of Moments Method	77
2.1	Introduction	77
2.2	Formulation of the model	78
2.3	Basis Pursuit	81
2.4	Exact reconstruction using Generalized Minimal Extrapolation	85
2.4.1	Reconstruction of a cone	88
2.5	Exact reconstruction of the nonnegatives measures	89
2.5.1	Markov systems	90
2.6	BLASSO	92
2.7	The Quadrature Method of Moments	95
2.7.1	The Product Difference Algorithm	95
2.8	Examples of Implementation of QMOM	96
2.8.1	Theoretical moments of the case of Silva 2010	96

3	Numerical Resolution of PBE	103
3.1	Introduction and theoretical considerations	103
3.1.1	The framework of Population Balance	104
3.1.2	Trapezoidal rule for regular and geometric grid	111
3.1.3	Simulations using the Discretized PBE	121
4	Parameter Estimation via EKF and OLS	125
4.1	Introduction	125
4.2	Estimation of the Number Distribution using Extended Kalman Filter	126
4.2.1	The Extended Kalman Filter	126
4.3	Least Squares Estimators for PBE parameters	129
4.3.1	General Information	129
4.4	Results	138
4.4.1	Simulation using the particular case (Silva 2010)	138
4.4.2	Simulation from the case of (Vlieghe, 2016) theoretical study	140
4.4.3	Simulation from the case of Vlieghe 2016	143
4.4.4	Simulation using the EKF and LSE in a simulated case	144
4.4.5	Parameter estimation of parameters values of Guerin datasets	144
5	Conclusion	145
5.1	Conclusions of the descriptive analysis of the data.	146

List of Figures

1.1	Histograms in natural scale for the Size properties	27
1.2	Histograms in logarithmic scale for the Size properties	29
1.3	Box-plots in logarithmic scale for the Size properties	30
1.4	Histograms for the Shape properties	31
1.5	Box-plots for the Shape properties	32
1.6	Frequency polygons comparing the two data sets obtained under 30 rpm. Each size variable is represented under a logarithmic transformation. The aim is to evaluate the replicability of the experimental results. The two distributions are almost overlapped	34
1.7	36
1.8	Frequency polygons comparing the two data sets obtained under 50 rpm. Each size variable is represented under a logarithmic transformation. The aim is to evaluate the replicability of the experimental results. Again, the two distributions are almost overlapped	37
1.9	Frequency polygons comparing the two data sets obtained under 50 rpm. Each shape variable is represented. The aim is to evaluate the replicability of the experimental results. As before, the two distributions are almost overlapped as in the precedent set of variables.	38
1.10	Frequency polygons comparing the two data sets obtained under 70 rpm. Each size variable is represented under a logarithmic transformation. The aim is to evaluate the replicability of the experimental results. Again, the two distributions are almost overlapped	40
1.11	Frequency polygons comparing the two data sets obtained under 70 rpm. Each shape variable is represented. The aim is to evaluate the replicability of the experimental results. As before, the two distributions are almost overlapped as in the precedent set of variables.	41

1.12	Frequency polygons comparing the two data sets obtained under 90 rpm. Each size variable is represented under a logarithmic transformation. The aim is to evaluate the replicability of the experimental results. Again, the two distributions are almost overlapped	43
1.13	Frequency polygons comparing the two data sets obtained under 90 rpm. Each shape variable is represented. The aim is to evaluate the replicability of the experimental results. As before, the two distributions are almost overlapped as in the precedent set of variables.	44
1.14	Frequency polygons comparing the populations obtained under 30, 50, 70, and 90 rpm. Each size variable is represented under a logarithmic transformation. In general, we observed an increase of the size of the flocks when the level of rpm increase. The group of small flocks also decrease when the rpm increase. This tendency is observed in all size variables	46
1.15	Box and whiskers comparing the populations obtained under 30, 50, 70, and 90 rpm. Each size variable is represented under a logarithmic transformation. We observe that the median augment when the level of rpm does. Also, the number of flocks with large size augment when the speed of mixing does.	47
1.16	Frequency polygons comparing the populations obtained under 30, 50, 70, and 90 rpm. Each shape variable is represented. The roundness of the flocks decrease when the level of rpm augment. Also, lower is the level of rpm, the more elongated are the flocks. When the level of rpm augment, more irregular are the flocks in terms of concavity and roughness.	48
1.17	Box and whiskers comparing the populations obtained under 30, 50, 70, and 90 rpm. Each Shape variable is represented. For Circularity, the quantity of individuals marked as outliers is larger at 30 rpm than the others levels of speed of mixing. About Aspect ratio, the distribution of the measures are more separated from 1 than at lower levels of rpm. The presence of outliers with low values of these variables is bigger when the level of rpm decrease. For the shape variables Circularity, Convexity and Solidity, the groups of individuals marked as outliers behaves like another population in terms of those variables	50
1.18	Scatterplots and correlations coefficients between all shape and size variables	66
1.19	Scatterplots and correlations coefficients between size variables	66

1.20	Scatterplots and correlations coefficients between shape variables	67
1.21	CPA variables	68
1.22	CPA variables	69
1.23	CPA variables	70
1.24	CPA variables	71
1.25	CPA variables	72
1.26	Dendogram for all the variables	73
1.27	PCA. Individuals	73
2.1	First 6 standard moments ofr the case of (Silva, 2011) (strength line) and the estimation with QMOM (dotted line)	99
2.2	First 6 standard moments ofr the case of (Silva, 2011) (strength line) and the estimation with QMOM (dotted line)	100
2.3	First 6 standard moments ofr the case of (Silva, 2011) (strength line) and the estimation with QMOM (dotted line)	101
4.1	Behavior of the Extended Kalman Filter recovering the number density. a) Theoretical state. b) Simulated state with noise; c) State recovered by EKF. d) Absolute error; e) Error in the last state recovered	141
4.2	Behavior of the discretized model using the kernels and parameters identified in Vlieghe 2016. a) Experimental volume distribution function data. b) Simulated state using the kernels identified by Vlieghe 2016. c) Simulated measure from the state obtained using the discretized model.	142
4.3	Behavior of the discretized model using the kernels and parameters identified in Vlieghe 2016. a) Experimental volume distribution function data. b) Empirical state identified by Vlieghe 2016. c) Recovered state using the kernels identified by Vlieghe 2016 and the EKF. d) Absolute error	143

List of Tables

1.1	Properties of Size	21
1.2	Properties of Shape	21
1.3	Univariate Descriptive Statistics of Size Variables	26
1.4	Univariate Descriptive Statistics of Shape Variables	26

Chapter 1

Preliminary empirical data analysis and model formulation

1.1 Introduction

Analysis of a particulate system seeks to synthesize the behavior of the *population* of particles and its environment from the behavior of single particles in their local environments. The population is described by the **density** of a suitable extensive variable, usually the *number* of particles, but sometimes (with better reason) by other variables such as the *mass* or *volume* of particles [Ram00].

Population balances are essential to scientists and engineers of widely variety of disciplines. In the application of population balances, one is more interested in the distribution of particle populations and their effect on the system behavior. Another feature of this systems is that they contain particles which are continually being created and destroyed by processes such as particles breakage and agglomeration. The particles of interest can be distinguished by its *internal* and *external* coordinates. The internal coordinates provide quantitative characterization of its distinguishing traits other than its localization while the external coordinates denote the location of the particles in the physical space. The joint space of internal and external coordinates will be called *particle state space*. These coordinates can be either discrete or continuous [Mar+03].

One application where we can find such kind of systems is in water treatment, specifically in the problem of removing colloidal particles. Turbidity (colloidal particles) originates from clay, microscopic organisms, municipal waters, color compounds, and organic water. A colloidal particle ranges in size from 1-500 nanometers (nm) or millimicrons (μm). Their small size

results in an extremely slow settling time and causes them to pass through a typical filter. For example, a colloidal particle with 1nm diameter has a settling velocity of 3 meters per million years. Because of the health and aesthetic issues they must be removed. [Saf]

Colloidal particles are removed from drinking water by flocculation. Chemicals (coagulants) are added to the water to allow the colloidal particles to agglomerate or come together to form floc. Different chemicals or coagulants may be added depending on the characteristics of the colloidal particles to be removed.

Hydrophobic colloidal particles remain suspended in solution because they repel each other due to the like charges they have adsorbed from solution and that remain on their surfaces. The magnitude of the particles repulsive forces is called the zeta potential. Destabilization occurs with the addition of chemicals that reduce the repulsive forces between the particles or lower the zeta potential. Destabilization takes advantage of the natural attraction between any two masses known as Van der Waals force. Without Destabilization, the particles repel each other and do not agglomerate into large particles that can settle out of solution [Cab11].

There are several mechanisms that allow the coagulant to destabilize colloids. In one, the coagulant increases the ionic strength of the water and compresses the thickness of the charged layer around each colloidal particle through repulsion. This allows the Van der Waal forces to draw the particles together to form a floc. Some coagulants directly neutralize the surface charge of the colloidal particles allowing the Van der Waals forces to predominate and pull the particles together. In another mechanism, the coagulant removes the colloidal particles by sweeping or bridging them into a precipitate mass.

Coagulants are divalent or trivalent cations (they have a +2 or +3 charge, respectively) or polymers. The most common trivalent cations usually aluminum and iron salts. The most common trivalent coagulants include aluminum sulfate or alum, $Al_2(SO_4)_3$, ferrous sulfate, $FeSO_4$. Lime, $Ca(OH)_2$, is the most common divalent coagulant. Polymers can also be used as coagulants, but they are more commonly used as coagulant aids. They result in larger flocs that are tougher, or less likely to disintegrate [Saa],[TD06].

These coagulants and an aid can be used singly or in combination to treat water. Sometimes additional turbidity (clay) is also added to allow for more rapid flock formation. To select the proper coagulant and dose, depend on the physical and chemical characteristics of water including pH, alkalinity, organic content, and original turbidity. However, the selection and dose must be optimized through jar testing.

The most common problems associated with coagulation are weak flocks

that do not stay together long enough to settle completely or flocks that settle poorly. Coagulant aides may be added to reduce or eliminate these problems. Depending on the specific characteristics of the water, it is not necessary to always use an aid. The addition of a coagulant aid may also reduce the amount of coagulant that is required. Coagulant aids may be nonionic or anionic polymers, sodium aluminate, activated silica, clay, acids, or alkalis.

To achieve coagulation/flocculation, 3 basins are used. The first is rapid mixing, used to mix the coagulant with the water. Next, the colloidal solids are allowed to flock together to form heavier, larger solids. These solids are settled out in a basin. Alternatively, all three steps can be achieved in one unit, a flocculating clarifier.

Rapid mixing is used to contact the coagulant with the water to be treated, typically using mechanical mixing or, less often, a hydraulic jump. Complete blending should occur in 10-30 seconds. In general, a square tank is the best.

Once flocculation is complete, the water flows into a sedimentation basin to permit the flock to settle out. The solids are wasted through the bottom and the water flows out of the basin via weirs.

The discussion thus far described the three units needed for coagulation/flocculation. It uses a system of three independent tanks. An alternative is a flocculator-clarifier which is also known as a solids contact unit or an up-flow tank. Mixing, flocculation, and sedimentation all occur in one tank. These units require less space than separate tanks but offer less operating flexibility.

Once the coagulant and any desired coagulation aids have been blended using rapid mixing, flocculation begins. During flocculation, the water is slowly agitated to allow the colloidal particles to bump into each other and agglomerate into larger and heavier flocks. Mixing can be achieved with paddle flocculation, flat blade turbines, and vertical turbine mixers. Based on experience, it has been determined that the optimal flocculation system should consist of a minimum of three tanks, in series, or three sections within one tank. In either case, there should be progressively slower mixing throughout the tank or tanks. The water should travel at a velocity of 0.5-1.5 feet/minute and at least 30 minutes should be allowed for flocculation.

The U.S. EPA has recently developed a protocol for enhanced coagulation-flocculation that calls for jar testing to optimize the coagulation-flocculation process. This is designed to maximize pathogen removal by maximizing turbidity removal. As a result, the aesthetic quality is also improved. The best coagulant, coagulant aid, chemical dosage, mixing speed, and flocculation time can be evaluated through a laboratory jar test. In jar test, a paddle

stirrer is used to blend and slowly agitate multiple samples. Rapid mixing occurs at a speed of 60-80 rpm for 1 minute and slow agitation occurs at a speed of 30 rpm for approximately 15 minutes. Part of the purpose of the jar testing, however, is to optimize these values.

The removal of colloidal solids from potable water is critical for two reasons. Some colloidal solids are pathogens, including protozoa such as cryptosporidium. Larger microorganisms are difficult to disinfect. Further, colloidal solids also can shelter attached pathogens making disinfection difficult. If very clear water is not produced, there is a probability that pathogens are present. Colloidal solids may also be organic materials that can be precursors for the formation of disinfection byproducts that can be carcinogens. Disinfection byproducts are synthesized when chlorine is added. They are regulated down to very low levels. Minimizing the organic materials minimizes the production of these dangerous byproducts

1.2 Formulation of the model

1.2.1 Introduction

There are a wide variety of applications in physics and chemistry that involve systems of particles in their nature. More precisely, systems of micro-particles are important in the study of flocculation and coagulation processes, as long as in other processes like crystallization. Such kind of processes are found in important industrial applications like treatment of water for human beings and pharmaceutical applications.

The description of population of particles and its evolution in time has been formulated through the **Population Balance Equation (PBE)**. The PBE is an integral-differential equation involving of the number distribution of particles in function of one or several of its morphological properties. These properties allow to describe each particle in terms of its size (or mass) and its shape (geometrically and about its surface). The time evolution of the population particles is formulated through the integral-differential equation involving the derivative of the number distribution with respect to time and the different processes which have an effect on the formation (or destruction) of particles.

The solution of this integral-differential equation is essential for Computational Dynamics Simulations of processes as flocculation. The solution of the PBE results in the evolution in time of the number distribution or, the evolution in time of some moments of this distribution allowing to characterize the particles population.

This section is divided into three parts. At first, we are going to explain the usual model formulation using the PBE having the volume of the particle as main property. After that, in the second subsection, this model will be presented in terms of the size of the particle which is used frequently due its direct physic interpretation. We are going to consider this last model in order to present the time evolution of the moments of the number distribution from the PBE in third subsection.

1.2.2 Formulation of the model

We are concerned with systems consisting of particles dispersed in an environmental phase which we shall refer to as the continuous phase. The particles may interact between themselves as well as with the continuous phase. Such behavior may vary from particle to particle depending upon a number of "properties" that may be associated with the particle. Continuous variables may be encountered more frequently in population balance analysis. The external coordinates denoting the position vector of (the centroid of) a particle describing continuous motion through space represent continuous variables.

The temporal evolution of the particulate system: We shall regard time as varying continuously and inquire into the rate of change of the particle state variables. It is convenient to deal with continuous variables in this regard. A fundamental assumption here is that the rate of change of state of any particle is a function only of the state of the particle in question and the local continuous phase variables. Thus we exclude the possibility of direct interactions between particles, although indirect interaction between particles via the continuous phase is indeed accounted for because of the dependence of particle behavior on the "local" continuous phase variables. In order to enable such a local characterization of the continuous phase variables, it is necessary to assume that the particles are considerably smaller than the length scale in which the continuous phase quantities vary. The continuous phase variables may be assumed to satisfy the usual transport equations with due regard to interaction with the particulate phase. Thus, such transport equations will be coupled with the population balance equation.

Particle state vector We are concerned with particle phase variables that are continuous. In general, the choice of the particle state is determined by the variable needed to specify:

- The rate of change of those of direct interest to the application, and
- The birth and death processes.

The particle state may generally be characterized by a finite dimensional vector.

- External coordinate $\mathbf{r} \equiv (r_1, r_2, r_3)$ denote the position (of the centroid) of the particle.
- Internal coordinates $\mathbf{x} \equiv (x_1, \dots, x_d)$ representing d different quantities associated with the particle.

The particle state vector (\mathbf{x}, \mathbf{r}) accounts for both internal and external coordinates. We shall let Ω_x represent the domain of internal coordinates and Ω_r be the domain of external coordinates, which is the set of points in physical space in which the particles are present. These domains may be bounded or may have infinite boundaries.

The particle population may be regarded as being randomly distributed in the particle state space, which include both external or internal coordinates.

Our concern will be about large populations, which will display relatively deterministic behavior because the random behavior of individual particles will be averaged out.

The continuous phase vector. The continuous phase variables may be collated into a finite c -dimensional vector field. The continuous phase variables affect the behavior of each particle.

We define a continuous phase vector. $Y(r, t) = [Y_1(r, t), \dots, Y_c(r, t)]$ which is clearly a function only of the external coordinate r and time t .

The evolution of this field in space and time is governed by the laws of transport and interaction with the particles.

In some applications, a continuous phase balance may not be necessary because interactions between the population and the continuous phase may not bring about any (or a substantial enough) change in the continuous phase. In such case, analysis of the population involves only the population balance equation.

The number density function. We postulate that there exist an average number density function defined on the particle state space,

$$E[n(x, r, t)] \equiv n_1(x, r, t)$$

with $x \in \Omega_x$ and $r \in \Omega_r$, where $E[n(x, r, t)]$ denote the expectation or the average of the actual number density $n(x, r, t)$, while $n_1(x, r, t)$ denotes the average number density. This definition implies that the average number of particles in the infinitesimal volume $dV_x dV_r$ (in the particle state space)

about the particle state (x, r) is $n_1(x, r, t) dV_x dV_r$. However, we will refer to particles in volume $dV_x dV_r$ about the particle state (x, r) .

The average number density $n_1(x, r, t)$ is assumed to be sufficiently smooth to allow differentiation with respect to any of its arguments as many times as may become necessary.

The (average) number density allows one to calculate the (average) number of particles in any region of particle state space. Thus, the (average) total number of particles in the entire system is given by

$$\int_{\Omega_x} dV_x \int_{\Omega_r} dV_r n_1(x, r, t)$$

where dV_x and dV_r are infinitesimal volume measures in the spaces of internal and external coordinates respectively.

The local (average) number density in physical space, i. e. the (average) total number of particles per unit volume of physical space, denoted $N(r, t)$ is given by

$$N(r, t) = \int_{\Omega_x} dV_x n_1(x, r, t).$$

Other densities such as volume or mass density may also be defined for the particle population. Thus, if $v(x)$ is the volume of the particle of internal state x , then the volume density may be defined as $v(x) f_1(x, r, t)$.

The volume fraction density $\phi(x, r, t)$ of a particle state is defined by

$$\phi(x, r, t) = \frac{1}{\Phi(r, t)} v(x) n_1(x, r, t)$$

where

$$\Phi(r, t) = \int_{\Omega_x} dV_x v(x) f_1(x, r, t)$$

the denominator above represents the total volume fraction of all particle. Similarly, mass fractions can also be defined. For the case of scalar interval state using volume, the volume fraction density of particles of volume v becomes

$$\phi(v, r, t) = \frac{vn_1(v, r, t)}{\Phi(r, t)}$$

where

$$\Phi(r, t) = \int_0^\infty vn_1(v, r, t) dv.$$

In contrast with number density, volume or mass denote the amount of dispersed phase material.

The rate of change of particle state vector We observed earlier that particle state might vary in time. We are concerned with smooth changes in particle state describable by some vector field defined over the particle state space both internal and external coordinates.

While changes of external coordinates refers to motion through physical space, that of internal coordinates refers to motion through an abstract property space (for example size).

We had collectively referred to them as convective processes for the reasons that they might be likened to physical motion.

It will be convenient to define "velocity" $\dot{R}(x, r, Y, t)$ for internal coordinates and $\dot{X}(x, r, Y, t)$ for external coordinates. These functions are assumed to be as smooth as necessary.

The velocity just defined may be random processes in space and time. Thus, $n_1(x, r, t) \dot{R}(x, r, Y, t)$ represents the particle flux through physical space and $n_1(x, r, t) \dot{X}(x, r, Y, t)$ is the particle flux through internal coordinate space.

The Population balance equation. Consider a population of particles distributed according to their size x which we shall take to be the mass of the particle and allow it to vary between 0 and ∞ .

The particles are uniformly distributed in space so that the number density is independent of external coordinates. Further, we assume for the present that the environment does not play any explicit role in particle behavior.

We let $\dot{X}(x, t)$ be the growth rate of the particle size x and let $n_1(x, t)$ denote the number density. All functions involved are assumed to be sufficiently smooth. Thus, we have the population balance equation

$$\frac{\partial n_1(x, t)}{\partial t} + \frac{\partial \dot{X}(x, t) n_1(x, t)}{\partial x} = 0. \quad (1.1)$$

In the above derivation, we did not take in account the birth and death of particles. To assess the rates of these contributions detailed modeling of breakage and aggregation processes will be needed. Let $h(x, t) dx$ the net rate of generation of particles in the size range x to $x + dx$, where the identity of $h(x, t)$ would depend on the models of breakage and aggregation. In this case, the population balance equation becomes

$$\frac{\partial n_1(x, t)}{\partial t} + \frac{\partial \dot{X}(x, t) n_1(x, t)}{\partial x} = h(x, t) \quad (1.2)$$

The preceding equation must be supplemented with initial and boundary conditions. The initial condition must clearly stipulate the distribution of particles in the particle state space.

The Population Balance Equation (PBE) is an equation that describes the evolution of one population of particles in colloidal systems. Changes in this kind of population are due to aggregation or breakage processes that can be seen as processes of birth and death. The formulation of PBE is traditionally made in terms of the particle's volume as Size property.

Thus, in the case of populations of particles in colloidal systems, the PBE is formulated modeling $h(x, t)$ in function of processes of birth due to aggregation $B_a(v; t)$, death due to aggregation $D_a(v; t)$, birth due to breakage $B_b(v; t)$, and death due to breakage $D_b(v; t)$. For the populations of particles in colloidal systems, the particle flux through size coordinate is constant so

$$\frac{\partial \dot{X}(x, t) n_1(x, t)}{\partial x} = 0.$$

Let's denote $\eta(v; t)$ the number density in function of the particle's volume as Size property. In the case considered, the PBE have the form

$$\frac{\partial \eta(v; t)}{\partial t} = B_a(v; t) - D_a(v; t) + B_b(v; t) - D_b(v; t). \quad (1.3)$$

This four terms at the right side of the equation are the corresponding processes of birth and death due to aggregation or breakage. This equation is expressed like

Definition 1.1. Population Balance Equation in terms of the particle's volume. The equation governing the evolution in time of the number distribution of a population of colloidal particles is known as Population Balance Equation in terms of the particle's volume, and it is defined as

$$\begin{aligned} \frac{\partial \eta(v; t)}{\partial t} &= B_a(v; t) - D_a(v; t) + B_b(v; t) - D_b(v; t) \\ &= \frac{1}{2} \int_0^v \phi(v - \epsilon, \epsilon) \eta(v - \epsilon; t) \eta(\epsilon; t) d\epsilon \\ &\quad - \eta(v; t) \int_0^\infty \phi(v, \epsilon) \eta(\epsilon; t) d\epsilon \\ &\quad + \int_v^\infty \psi(\epsilon) \rho(v/\epsilon) \eta(\epsilon; t) d\epsilon \\ &\quad - \psi(v) \eta(v; t), \end{aligned} \quad (1.4)$$

where

- $B_a(v; t) = \frac{1}{2} \int_0^v \phi(v - \epsilon, \epsilon) \eta(v - \epsilon; t) \eta(\epsilon; t) d\epsilon$: birth rate of particles with volume v by aggregation of little particles,

- $D_a(v; t) = \eta(v; t) \int_0^\infty \phi(v, \epsilon) \eta(\epsilon; t) d\epsilon$: death rate of particles with volume v by aggregation with other particles,
- $B_b(v; t) = \int_v^\infty \psi(\epsilon) \rho(v/\epsilon) \eta(\epsilon; t) d\epsilon$: birth rate of particles with volume v by breakage of big particles,
- $D_b(v; t) = \psi(v) \eta(v; t)$: death rate of particles with volume v by breakage into little particles.

and where

- $\eta(v; t)$: number density function using volume as coordinate,
- $\phi(v, \epsilon)$: aggregation kernel using volume as coordinate,
- $\psi(v)$: breakage kernel using volume as coordinate,
- $\rho(v/\epsilon)$: distribution function of fragments.

In some applications, it is interesting to express the PBE in terms of the Length of Diameter particle instead of the volume. Because of this, we are going to see how we can transform the PBE using the particle's volume as distribution variable to the PBE using the particle's size as distribution variable. We can see that this formulation is

Proposition 1.2. *The PBE in terms of particle's size coordinate. The PBE can be formulate in terms of the length coordinate like*

$$\begin{aligned}
\frac{\partial n(L, t)}{\partial t} &= B^a(L; t) - D^a(L; t) + B^b(L; t) - D^b(L; t) \\
&= \frac{L^2}{2} \int_0^L \frac{\beta\left(\frac{(L^3 - \lambda^3)^{1/3}}{L}, \lambda\right)}{(L^3 - \lambda^3)^{2/3}} n\left(\frac{(L^3 - \lambda^3)^{1/3}}{L}, t\right) n(\lambda, t) d\lambda \\
&\quad - n(L, t) \int_0^\infty \beta(L, \lambda) n(\lambda, t) d\lambda \\
&\quad + \int_L^\infty a(\lambda) b(L | \lambda) n(L, t) d\lambda \\
&\quad - a(L) n(L, t),
\end{aligned} \tag{1.5}$$

where where

- $B^a(L; t) = \frac{L^2}{2} \int_0^L \frac{\beta\left(\frac{(L^3 - \lambda^3)^{1/3}}{L}, \lambda\right)}{(L^3 - \lambda^3)^{2/3}} n\left(\frac{(L^3 - \lambda^3)^{1/3}}{L}, t\right) n(\lambda, t) d\lambda$: birth rate of particles with length L by aggregation of little particles,

- $D^a(L; t) = n(L, t) \int_0^\infty \beta(L, \lambda) n(\lambda, t) d\lambda$: death rate of particles with length L by aggregation with other particles,
- $B^b(v; t) = \int_L^\infty a(\lambda) b(L | \lambda) n(L, t) d\lambda$: birth rate of particles with length L by breakage of big particles,
- $D^b(v; t) = a(L) n(L, t)$: death rate of particles with length L by breakage into little particles.

and where

- $n(L; t)$: number density function using length as coordinate,
- $\beta(L, \lambda)$: aggregation kernel using length as coordinate,
- $a(L)$: breakage kernel using length as coordinate,
- $b(L/\lambda)$: distribution function of fragments.

Proof. To write the equation (3.1) in terms of the size particle instead of the volume coordinate, we need to assume that $v \propto L^3$. Therefore, we can propose the variable transformation

$$v = L^3 \qquad dv = 3L^2 dL,$$

and we can write the number density distribution in terms of particle's length $n(L; t) dL$ from the distribution in terms of the volume $\eta(v; t)$ like

$$\eta(v; t) dv = \eta(L^3; t) 3L^2 dL = n(L; t) dL,$$

also, we can write the aggregation kernel, breakage kernel and the distribution function of fragments in terms of the particle's length if we take $v = L^3$ and $\epsilon = \lambda^3$ like

- $\phi(v, \epsilon) = \phi(L^3, \lambda^3) = \beta(L, \lambda)$,
- $\psi(v) = \psi(L^3) = a(L)$,
- $\rho(v|\epsilon) = \rho(L^3|\lambda^3) 3L^2 = b(L|\lambda)$.

Using the factor $3L^2$ in both sides of the equation (3.1) and taking again $v = L^3$ and $\epsilon = \lambda^3$ we have, for the left side

$$\frac{\partial \eta(v; t)}{\partial t} = \frac{\partial \eta(L^3; t) 3L^2}{\partial t} = \frac{\partial n(L; t)}{\partial t},$$

and for the right side, we can write each process involved into the PBE respectively like:

- **Birth process due to aggregation.** We are going to define

$$B_a(v, t) 3L^2 = B_a(L^3, t) 3L^2 = B^a(L; t).$$

From the definition 3.1 we have

$$B_a(v; t) = \frac{1}{2} \int_0^v \phi(v - \epsilon, \epsilon) \eta(v - \epsilon; t) \eta(\epsilon; t) d\epsilon,$$

if we multiply this term by the factor $3L^2$ and we use the variable transformation

$$\begin{aligned} v &= L^3 & dv &= 3L^2 dL \\ \epsilon &= \lambda^3 & d\epsilon &= 3\lambda^2 d\lambda \end{aligned}$$

we obtain

$$\begin{aligned} B_a(L^3; t) 3L^2 &= \frac{3L^2}{2} \int_0^L \beta\left((L^3 - \lambda^3)^{1/3}, \lambda\right) \eta(L^3 - \lambda^3; t) \eta(\lambda^3; t) 3\lambda^2 d\lambda \\ &= \frac{3L^2}{2} \int_0^L \beta\left((L^3 - \lambda^3)^{1/3}, \lambda\right) \eta(L^3 - \lambda^3; t) \times \\ &= \frac{3(L^3 - \lambda^3)^{2/3}}{3(L^3 - \lambda^3)^{2/3}} n(\lambda; t) d\lambda \\ &= \frac{L^2}{2} \int_0^L \frac{\beta\left((L^3 - \lambda^3)^{1/3}, \lambda\right)}{(L^3 - \lambda^3)^{2/3}} n\left((L^3 - \lambda^3)^{1/3}; t\right) n(\lambda; t) d\lambda, \end{aligned}$$

then, we get the term corresponding the birth aggregation processes in function of the particle's length like

$$B^a(L; t) = \frac{L^2}{2} \int_0^L \frac{\beta\left((L^3 - \lambda^3)^{1/3}, \lambda\right)}{(L^3 - \lambda^3)^{2/3}} n\left((L^3 - \lambda^3)^{1/3}; t\right) n(\lambda; t) d\lambda. \quad (1.6)$$

- **Death processes due to aggregation.** Similarly, we are going to define

$$D_a(v, t) 3L^2 = D_a(L^3, t) 3L^2 = D^a(L; t).$$

From the definition 3.1 we have

$$D_a(v; t) = \eta(v; t) \int_0^\infty \phi(v, \epsilon) \eta(\epsilon; t) d\epsilon,$$

if we multiply this term by the factor $3L^2$ and we use the variable transformation

$$\begin{aligned} v &= L^3 & dv &= 3L^2 dL \\ \epsilon &= \lambda^3 & d\epsilon &= 3\lambda^2 d\lambda \end{aligned}$$

we obtain

$$\begin{aligned} D_a(L^3; t) 3L^2 &= \eta(L^3; t) 3L^2 \int_0^\infty \phi(L^3, \lambda^3) \eta(\lambda^3; t) 3\lambda^2 d\lambda \\ &= n(L; t) \int_0^\infty \beta(L; \lambda) n(\lambda; t) d\lambda. \end{aligned} \quad (1.7)$$

then, we get the term corresponding the death aggregation processes in function of the particle's length like

$$D^a(L; t) = n(L; t) \int_0^\infty \beta(L; \lambda) n(\lambda; t) d\lambda. \quad (1.8)$$

- **Birth processes due to breakage.** Similarly, we are going to define

$$B_b(v; t) 3L^2 = B_b(L^3; t) 3L^2 = B^b(L^3; t).$$

From the definition 3.1 we have

$$B_b(v; t) = \int_v^\infty \psi(\epsilon) \rho(v/\epsilon) \eta(\epsilon; t) d\epsilon$$

if we multiply this term by the factor $3L^2$ and we use the variable transformation

$$\begin{aligned} v &= L^3 & dv &= 3L^2 dL \\ \epsilon &= \lambda^3 & d\epsilon &= 3\lambda^2 d\lambda \end{aligned}$$

we obtain

$$\begin{aligned} B_b(v; t) 3L^2 &= 3L^2 \int_L^\infty \psi(\lambda^3) \rho(L^3/\lambda^3) \eta(\lambda^3; t) 3\lambda^2 d\lambda \\ &= \int_L^\infty \psi(\lambda^3) \rho(L^3/\lambda^3) 3L^2 n(\lambda; t) d\lambda \\ &= \int_L^\infty a(\lambda) b(L/\lambda) n(\lambda; t) d\lambda, \end{aligned}$$

then, we get the term corresponding the birth breakage processes in function of the particle's length like

$$B^b(L; t) = \int_L^\infty a(\lambda) b(L/\lambda) n(\lambda; t) d\lambda$$

- **Death processes due to breakage.** Finally, we are going to define

$$D_b(v; t) 3L^2 = D_b(L^3; t) 3L^2 = D^b(L^3; t).$$

From the definition 3.1 we have

$$D_b(v; t) = \psi(v) \eta(v; t), \quad (1.9)$$

if we multiply this term by the factor $3L^2$ and we use the variable transformation

$$v = L^3 \quad dv = 3L^2 dL$$

we obtain

$$D_b(v; t) 3L^2 = \psi(v) \eta(v; t) 3L^2 \\ a(L) n(L, t),$$

then, we get the term corresponding the death breakage processes in function of the particle's length like

$$D^b(L; t) = a(L) n(L, t). \quad (1.10)$$

Finally, we can formulate the PBE in terms of the particle's length like

$$\begin{aligned} \frac{\partial n(L, t)}{\partial t} &= B^a(L; t) - D^a(L; t) + B^b(L; t) - D^b(L; t) \\ &= \frac{L^2}{2} \int_0^L \frac{\beta\left((L^3 - \lambda^3)^{1/3}, \lambda\right)}{(L^3 - \lambda^3)^{2/3}} n\left((L^3 - \lambda^3)^{1/3}, t\right) n(\lambda, t) d\lambda \\ &\quad - n(L, t) \int_0^\infty \beta(L, \lambda) n(\lambda, t) d\lambda \\ &\quad + \int_L^\infty a(\lambda) b(L | \lambda) n(\lambda, t) d\lambda \\ &\quad - a(L) n(L, t) \end{aligned}$$

□

The Population Balance Equation in terms of particle's size will be used throughout the next chapters for modeling the time evolution of the number density function.

1.3 Moments equations for aggregation and break-age

For some applications, it is convenient to study the evolution in time of a finite number of standard moments of the number density function. In this case, the properties of the particle's population of interest are fully known as a function of a few set of standard moments. Although, an equation like the Population Balance Equation can be used for studying the time evolution of the standard moments.

In the precedent section 1.3, we introduced the Population Balance Equation describing the time evolution of the number density function in terms of the particle's size coordinate like it was shown in proposition 3.2 and equation (3.2). Also, from that equation (3.2) it is possible to find an equivalent PBE describing the time evolution of the standard moments of the number density in terms of the particle size coordinate as it can be seen in the next proposition.

Proposition 1.3. *The Population Balance Equation describing the time evolution of the standard moments of the number density in terms of the particle size coordinate can be expressed like*

$$\begin{aligned}
\frac{\partial m_k(t)}{\partial t} &= \frac{1}{2} \int_0^L n(\lambda; t) \int_0^\infty \beta(u, \lambda) n(u; t) (u^3 + \lambda^3)^{k/3} du d\lambda \\
&\quad - \int_0^\infty L^k n(L; t) \int_0^\infty \beta(L, \lambda) n(\lambda; t) d\lambda dL \\
&\quad + \int_0^\infty L^k \int_0^\infty a(\lambda) b(L/\lambda) n(\lambda; t) d\lambda dL \\
&\quad - \int_0^\infty a(\lambda) n(L; t) L^k dL
\end{aligned} \tag{1.11}$$

Proof. We begin with the transformations of moments [MVF03]

$$m_k(t) = \int_0^\infty n(L; t) L^k dL \tag{1.12}$$

If we introduce this transformation into the PBE in the form of equation

(3.2) we get

$$\begin{aligned}
\frac{\partial n(L; t)}{\partial t} &= B^a(L; t) - D^a(L; t) + B^b(L; t) - D^b(L; t) \\
\frac{\partial \int_0^\infty n(L; t) L^k dL}{\partial t} &= \int_0^\infty [B^a(L; t) - D^a(L; t) + B^b(L; t) - D^b(L; t)] L^k dL \\
\frac{\partial m_k(t)}{\partial t} &= \int_0^\infty B^a(L; t) L^k dL - \int_0^\infty D^a(L; t) L^k dL \\
&\quad + \int_0^\infty B^b(L; t) L^k dL - \int_0^\infty D^b(L; t) L^k dL
\end{aligned} \tag{1.13}$$

now, if we denote

1. $\overline{B}_k^a = \int_0^\infty B^a(L; t) L^k dL$
2. $\overline{D}_k^a = \int_0^\infty D^a(L; t) L^k dL$
3. $\overline{B}_k^b = \int_0^\infty B^b(L; t) L^k dL$
4. $\overline{D}_k^b = \int_0^\infty D^b(L; t) L^k dL$

then, we get

$$\frac{\partial m_k(t)}{\partial t} = \overline{B}_k^a - \overline{D}_k^a + \overline{B}_k^b - \overline{D}_k^b \tag{1.14}$$

Now, we can get the expressions for these new components

$$\begin{aligned}
\overline{B}_k^a &= \int_0^\infty B^a(L; t) L^k dL \\
&= \int_0^\infty \left[\frac{L^2}{2} \int_0^L \frac{\beta\left((L^3 - \lambda^3)^{1/3}, \lambda\right)}{(L^3 - \lambda^3)^{2/3}} n\left((L^3 - \lambda^3)^{1/3}; t\right) n(\lambda; t) d\lambda \right] L^k dL
\end{aligned}$$

if we use the change $u^3 = L^3 - \lambda^3$ and $dL = \frac{u^2}{L^2} du$ we get

$$\begin{aligned}
\overline{B}_k^a &= \int_0^\infty \left[\frac{L^2}{2} \int_0^L \frac{\beta(u, \lambda)}{u^2} n(u; t) n(\lambda; t) d\lambda \right] L^k dL \\
&= \frac{1}{2} \int_0^\infty \left[\int_0^L \beta(u, \lambda) n(u; t) n(\lambda; t) d\lambda \right] \frac{L^2}{u^2} L^k \frac{u^2}{L^2} du
\end{aligned}$$

and we have that $u \rightarrow 0$ as $L \rightarrow 0$, and $u \rightarrow \infty$ as $L \rightarrow \infty$. Then we can write

$$\begin{aligned}
\overline{B}_k^a(t) &= \frac{1}{2} \int_0^L n(\lambda; t) \int_{-\lambda}^{\infty} \beta(u, \lambda) n(u; t) (u^3 + \lambda^3)^{k/3} du d\lambda \\
\overline{B}_k^a(t) &= \frac{1}{2} \int_0^L n(\lambda; t) \int_0^{\infty} \beta(u, \lambda) n(u; t) (u^3 + \lambda^3)^{k/3} du d\lambda \\
\overline{D}_k^a(t) &= \int_0^{\infty} \left[n(L; t) \int_0^{\infty} \beta(L, \lambda) n(\lambda; t) d\lambda \right] L^k dL \\
\overline{D}_k^a(t) &= \int_0^{\infty} L^k n(L; t) \int_0^{\infty} \beta(L, \lambda) n(\lambda; t) d\lambda dL \\
\overline{B}_k^b(t) &= \int_0^{\infty} \left[\int_L^{\infty} a(\lambda) b(L/\lambda) n(\lambda; t) d\lambda \right] L^k dL \\
\overline{B}_k^b(t) &= \int_0^{\infty} L^k \int_0^{\infty} a(\lambda) b(L/\lambda) n(\lambda; t) d\lambda dL \\
\overline{D}_k^b(t) &= \int_0^{\infty} a(\lambda) n(L; t) L^k dL
\end{aligned} \tag{1.15}$$

□

1.4 Preliminary Empirical Data

1.4.1 Introduction

In the first part of this chapter we formulated the model describing the population of particles in terms of the PBE. This formulation involves the number distribution in function of one main property (size of volume). But often, this information is not enough for completely describe the population. The shape of the particle (for example, general geometric shape of shape of the surface of the particle) affects the hydrodynamic properties of the system and then affects the formation or destruction of particles in the system.

One better description of this kind of systems will involve a Population Balance Equation in terms of properties of size and shape. This kind of equations are not straightforward and often there are several morphological properties available and it is not clear which set of them describes better the population.

Several studies have been proposed in order to better understand how morphological properties can describe a population and how to incorporate this information to the PBE. One of this studies was ([Vli14]), in which

experimental data was obtained from a population of Bentonite particles using a Taylor-Couette reactor and several properties were measured using diffraction laser. In this study, the authors recovered data from one initial population. Then, they measured the properties in populations obtained using different hydrodynamic conditions. They proposed the reduction of the set of measures into two properties, one size property and one shape property. They use a Principal Components Analysis (PCA) for proposing this set of variables.

We are going to use this experimental data in order to study this morphological properties. We are going to do an analysis to the initial distribution of particles and data from several hydrodynamic conditions.

1.4.2 About the data sets

1.4.3 Materials and methods

In this section, we apply a set of statistical tools for analyzing the data-sets get from the experiments with the Bentonite. This data-sets were gotten from the different experiments under different hydrodynamic conditions, producing different populations of flocs. We distinguish between the initial population, and the populations changing the rate of speed of rotation.

In all cases, the granulometry by lasser diffraction are used to get a set of measures characterizing the morphology of flocs [Mor]. Those measures can be classified into Size measures and Shape measures.

We begin with a descriptive analysis for the Size measures and the Shape measures including univariate, bivariate and multivariate exploratory techniques. The objective is to describe the main characteristics of the population, including the presence of outliers. Also we want to observe the relations between the variables of the same group (Size variables and Shape variables), and the relation between the variables of different groups. We want to answer the following questions:

- What are the main characteristics of the populations of Bentonite flocs by each variable measured?
- Do we find outliers in the populations?
- What are the relationships between the variables in the same group of classification (Size and Shape)?
- What are the relationships between the variables of different groups?
- Does this relationships change in time?

- Do we find different groups of flocs in the same population? If we do find this groups, can we characterize those groups?

The experimental data are obtained by image's treatment. The images containing the flocs are taken in gray's scale. The MATLAB software used implements a smoothing step and a binarization step in order to get appropriate data. This software change the color in images from gray scale to black or with the binary scale. Also, the software fills the images of objects containing less than 10 pixels. The resulting images are appropriate for measuring the characteristics describing the flocs.

The image treatment of the experimental data sets will result in measures quantifying the morphology of flocs. Morphology is the set of characteristics related to the size and the shape of flocs, distinguishing the size properties measured in length units (or surface units) and the shape properties in general dimensionless.

Properties related to the size of the flocs

- Area and Circle Equivalent Diameter

The Area (A) is obtained automatically by counting the pixels belonging to the floc's image and measured in μm by convention. The Circle Equivalent Diameter (CED) is the diameter of the circle having an equivalent area A .

- Major and Minor axes

Taking in account the relative position of pixels, the software gets the length of the major (L) and minor (l) axes of the ellipse having the having the equivalent inertia.

- Hull Convex

It is the smallest polygon enclosing the image of the floc. Convex Hull is characterized by its area (A_{ch}) and its perimeter (P_{ch}). Those measures allow to calculate the solidity and the convexity of the shape, by comparison with the actual area and perimeter of the floc.

- Perimeter

In order to estimate the Perimeter of the floc's image, the sum of the distances between the centers of the pixels forming the contour of the image is calculated. Then a layer of pixels are added to the contour region and the sum of the distances between the centers of those pixels are computed. The Perimeter estimation (P) is the mean of these two measures.

- **Radius of Gyration**

The Radius of Gyration (Rg) is a characteristic size taking account the relative position of pixels. It is computed like the quadratic mean of distances between the pixels and the barycenter.

Properties related to the shape of the flocs

The properties describing the shape of the flocs were chosen such that they take the value of 1 when the shape is a circle.

- **Solidity and Aspect Ratio**

The solidity (S) compares the area of the floc A to the convex area A_{ch} . If $A = A_{ch}$ then solidity is equal to 1. Solidity is inferior to 1 in presence of concavities in the shape

$$S = \frac{A}{A_{ch}}.$$

The Aspect Ratio (AR) is the ratio between the length of the major L and minor axes l of the ellipse equivalent. This property takes the value of 1 if the ellipse equivalent is a circle, and it is inferior to 1 if the floc has an elongated shape

$$AR = \frac{l}{L}$$

- **Circularity and Convexity**

The circularity and the convexity take into account the perimeter of the floc. Both properties take values between 0 and 1, and it is 1 in the case of a circle.

The Circularity (C_I) is the ratio between the perimeter of the circle equivalent P_{ch} (having the same area) and the actual perimeter P

$$C_I = \frac{\sqrt{4\pi A}}{P}.$$

The Convexity (C_V) is the relation between the convex perimeter P_{ch} and the actual perimeter P . This property quantifies the presence of concavities and roughness in the flocs

$$C_V = \frac{P_{ch}}{P}.$$

Properties of Size

Property	Notation	Definition
Perimeter	A (μm)	Number of pixels in the image region.
Area	P (μm^2)	Cumulative distance between the pixels of the contour.
Hull Convex	A_{ch} (μm^2)	The smallest polygon enclosing the image region.
Hull Convex Perimeter	P_{ch} (μm)	
Circle Equivalent Diameter	CED (μm)	Diameter of the circle having area equal to $\sqrt{4A/\pi}$.
Length of the major and minor axes	L, l (μm)	Axes of the ellipse equivalent having the same inertia.
Radius of Gyration	Rg (μm)	mean distance between the pixels and the barycenter.

Table 1.1: Properties of Size

Properties of Shape

Property	Definition	Interpretation
Aspect Ratio	$AR = l/L$ (-)	Elongation of the global form.
Solidity	$S = A/A_{ch}$ (-)	Concavities.
Convexity	$C_V = P_{ch}/P$ (-)	Irregularities in the surface.
Circularity	$C_I = \sqrt{4\pi A}/P$ (-)	Concavities and irregularities.

Table 1.2: Properties of Shape

These characteristics measure then the similarity to the circle in all the scales. The convexity takes the same value if the figure is a circle or a square for example (in both cases, the convexity is 1). Although, the circularity does discriminate between a circle and a square, then circularity is more sensitive to various geometries.

1.4.4 Metodology

In order to answer those questions, we are going to use the following statistic techniques:

- Univariate descriptive statistic

- Exploratory bivariate statistics
- Exploratory multivariate statistics

Univariate descriptive statistic techniques

For describing each population considered, a set of classical univariate techniques are used. These techniques include graphical representation and main numerical statistical indexes. We use the graphics for having a global idea of the distribution of the population according to each variable. The numerical statistics are used to know one specific feature of interest of those populations like the average value, the dispersion or the asymmetry. Also, a combination of those graphical and numerical statistics allow to determine the presence of outliers or different groups in the same population.

Histogram and Box and whiskers plot ([HS15]). The boxplot is a graphical technique that displays the distribution of variables. It helps us to see the location, skewness, spread, tail length and outlying points. It is particularly useful in comparing different batches. The boxplot is a graphical representation of the Five Number Summary. In the Five Number Summary, we calculate the upper quartile Q_3 , the lower quartile Q_1 , the median M , the minimum and the maximum. The Q-spread, d_Q is defined as $d_Q = Q_3 - Q_1$. The outside bars

$$O_U = Q_3 + 1.5d_Q$$

$$O_L = Q_1 - 1.5d_Q$$

are the borders beyond which a point is regarded as an outlier. For the construction of the boxplot, we have:

1. Draw a box with borders (edges) at Q_1 and Q_3 (i.e. 50% of the data are in the box).
2. Draw the median as a solid line.
3. Draw "whiskers" from each end of the box to the most remote point that is not an outlier.
4. Show outliers as "points" depending on if they are outside of $[O_L; O_U]$.

The histograms are density estimates. A density estimate gives a good impression of the distribution of the data. In contrast to boxplots, density estimates show possible multimodality of the data. The idea is to locally represent the data density by counting the number of observations in a sequence

of consecutive intervals (bins) with origin x_0 . Let $B_j(x_0, h)$ denote the bin of length h which is the element of a bin grid starting at x_0 :

$$B_j(x_0, h) = [x_0 + (j - 1)h, x_0 + jh), j \in \mathbb{Z}$$

where, $[\cdot, \cdot)$ denotes a left closed and right open interval. If $\{x\}_{i=1}^n$ is an i.i.d. sample with density f , the histogram is defined as follows:

$$\hat{f}_h(x) = n^{-1}h^{-1} \sum_{j \in \mathbb{Z}} \sum_{i=1}^n \mathbf{I}\{x_i \in B_j(x_0, h)\} \mathbf{I}\{x \in B_j(x_0, h)\} \quad (1.16)$$

where in (1.16) the first indicator function $\mathbf{I}\{x_i \in B_j(x_0, h)\}$ counts the number of observations falling into bin $B_j(x_0, h)$. The second indicator function is responsible for "localizing" the counts around x . The parameter h is a smoothing or localizing parameter and controls the width of the histogram bins. One "optimal" h parameter for n observations is given by:

$$h_{opt} = \left(\frac{24\sqrt{\pi}}{n} \right)^{1/3}$$

Experimental distributions

In order to compute the empirical experimental distribution of each property the variation rang was divided into n_c classes. For shape properties, linear classes were constructed using

$$L_i = L_1 + (i - 1)(L_{n_c+1} - L_1)/n_c$$

and the abscissa representing each class, denoted by \bar{L}_i , was its arithmetic mean, for $i = 1, \dots, n_c$. For size properties geometrical classes were constructed using

$$L_i = L_1 \times \left(\exp \left(\ln \left(\frac{L_{n_c+1}}{L_1} \right) / n_c \right) \right)$$

and the abscissa representing each class, \bar{L}_i , was its geometric mean, for $i = 1, \dots, n_c$.

Then, The empirical experimental distribution of each property were computed counting the frequency (proportion of the population) in each class N_i , and this quantity was assigned to the abscissa \hat{L}_i . The frequencies are computed as (1.17) [Sal07]

$$N_i = \frac{\sum_{j=1}^{N_F} Q(j)}{\sum_{j=1}^{N_F} Z(j)} \quad (1.17)$$

where

$$Q(j) = \begin{cases} Z(j), & \text{if } L(j) \in [L_i, L_{i+1}[\\ 0, & \text{otherwise} \end{cases}$$

where $L(j)$ is the size of the j -th floc, for $j = 1, \dots, N_F$, and N_F is the total number of flocs considered. For a number distribution, the function $Z(j) = 1$ is taken and N_i represents the fraction in number of flocs belonging to the class i .

The empirical moments of the distribution m_k are computed using the density $n(l)$, where

$$n(l) = \frac{N_i}{L_{i+1} - L_i}, \quad \forall l \in [L_i, L_{i+1}[.$$

We can see the empirical moments as an approximation of the standard moments of the number distribution

$$\begin{aligned} m_k &= \int_0^{+\infty} x^k n(L) dL \\ &= \sum_i \int_{L_i}^{L_{i+1}} L^k n(L) dL \\ &\approx \sum_i \frac{N_i}{L_{i+1} - L_i} \int_{L_i}^{L_{i+1}} L^k dL \\ m_k &\approx \sum_i \frac{N_i}{L_{i+1} - L_i} \frac{(L_{i+1}^{k+1} - L_i^{k+1})}{k+1}. \end{aligned}$$

Simple Flocculation Experiments

A series of simple flocculation experiments (under controlled hydrodynamic conditions) were performed using Bentonite in a Taylor-Couette reactor. The measure of several morphological parameters was done using image analysis.

The experiments performed by [Vli14], describe the Bentonite flocculation (30 mgL^{-1} with aluminum sulfate at a concentration $3.5 \times 10^{-5} \text{ molL}^{-1}$) under four rotation speeds. Each experiment starts as follows:

- The reactor is filled to three quarters of its capacity with demineralized water and switched on with a rotation speed of the inner cylinder $N = 100 \text{ rpm}$. Then, the suspension is poured, followed by demineralized water until the reactor is almost full. Finally, the flocculant is added. The rotation speed is maintained at 100 rpm during 3 min . This step produces a population of initial aggregates, which is the initial population considered afterwards.

- Then, the speed rotation is set at the desired value (30, 50, 70, or 90 *rpm*) for 4 *h*.
- 70 photographs are taken every 2 *min* at the beginning of an experiment, then every 5 and 10 *min*.
- Only aggregate images composed of 10 pixels or more were taken into account, thus the smallest aggregates analyzed have a *CED* of 26 μm , *Rg* of 9 μm , and *L* of 28 μm .

1.4.5 Results of the univariate description

Now, the results of the univariate description of the data sets will be presented. It will be organized in three parts, each one corresponding to a different data set. In the first part will be described the initial population of Bentonite flocs obtained at a speed of mixing of 50 revolutions per minute (*rpm*). The second part will contain the results of the univariate description for different populations of flocs obtained at different hydrodynamic conditions (30, 50, 70, and 90 *rpm*). Finally, the third part will contain the description of several populations obtained under the same hydrodynamic conditions through the time. Several data set were available from different speed of mixing (30, 50, 70, 90, and 100 *rpm*).

Initial population of Bentonite flocs. In this part, the initial population of Bentonite flocs are described. We are going to begin with the univariate descriptive analysis. This will be a characterization of each size and shape variables in terms of numerical descriptive statistics and graphic description (histograms and box-plots), taking special attention to the outliers univariates.

The description was performed for data obtained for the speed of mixing of 50 *rpm*. In this data set were measured 36854 flocs. We describe the size and shape variables using histograms with the rectangles representing the density of floc's frequency divided by the amplitude of each class. We present the size measures in logarithmic scale, because the great variability in the flocs, as usually in granulometry analysis (granulometry).

The initial summary of numerical descriptive statistics are shown in the table 1.3 and 1.4

We can see in table 1.3 the description of the set of variables representing the size of the flock or "Size variables". We can observe that the range of variation in all the variables is wide. Also the the particles in the "50%

Univariate Descriptive Statistics of Size Variables for the Initial Population

Property	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Diameter	2.170	2.730	3.810	7.835	7.090	338.990
Length	2.22	3.71	5.28	11.41	10.10	718.56
Width	1.110	2.770	3.825	7.729	6.940	341.490
Perimeter	4.90	8.06	12.70	42.40	28.57	4645.41
Area Pixels	12	19	37	661	128	292647

Table 1.3: Univariate Descriptive Statistics of Size Variables

between the 1st and the 3rd quantile are distributed very near to the minimum values. This trend is very similar in each variable. We can remark the presence of very big flocks in the last 75% of the distribution.

Univariate Descriptive Statistics of Shape Variables

Property	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
HS Circularity	0.0320	0.3960	0.6200	0.6019	0.8230	1.0000
Convexity	0.2660	0.8610	0.9590	0.9066	1.0000	1.0000
Solidity	0.3040	0.8490	0.9460	0.9094	1.0000	1.0000
Elongation	0.0000	0.1580	0.2580	0.2718	0.3770	0.9300

Table 1.4: Univariate Descriptive Statistics of Shape Variables

The same kind of analysis was performed for the variables representing the shape of the particles or "Shape variables". By definition, these variables are limited to the range $[0, 1]$. We observe that for the Circularity, the flocks in the 50% between the 1st and the 3rd quantile are distributed in the center of the range. It could represent flocks with a non-regular surface but even so smooth. The mean value of the Circularity also indicates that the geometry shape of the flocks tend to a circle. In the case of Convexity and Solidity, their behavior is very similar. The 1st quantile is in both variables near to the value of 0.85. It means that the 75% of the flocks present values for these two variables very close to 1. It means that the flocks present little concavity and smooth surface. Finally, the Elongation presents values very close to 0. The 3rd quantile is approximately 0.38 and then, the 75% of the flocks have Elongation inferior to this value. Then, the flocks tend to have an elongated shape.

We use graphics like histograms and box-plots for having an idea of the distribution (proportion of individuals in each class) for the size variables and shape variables.

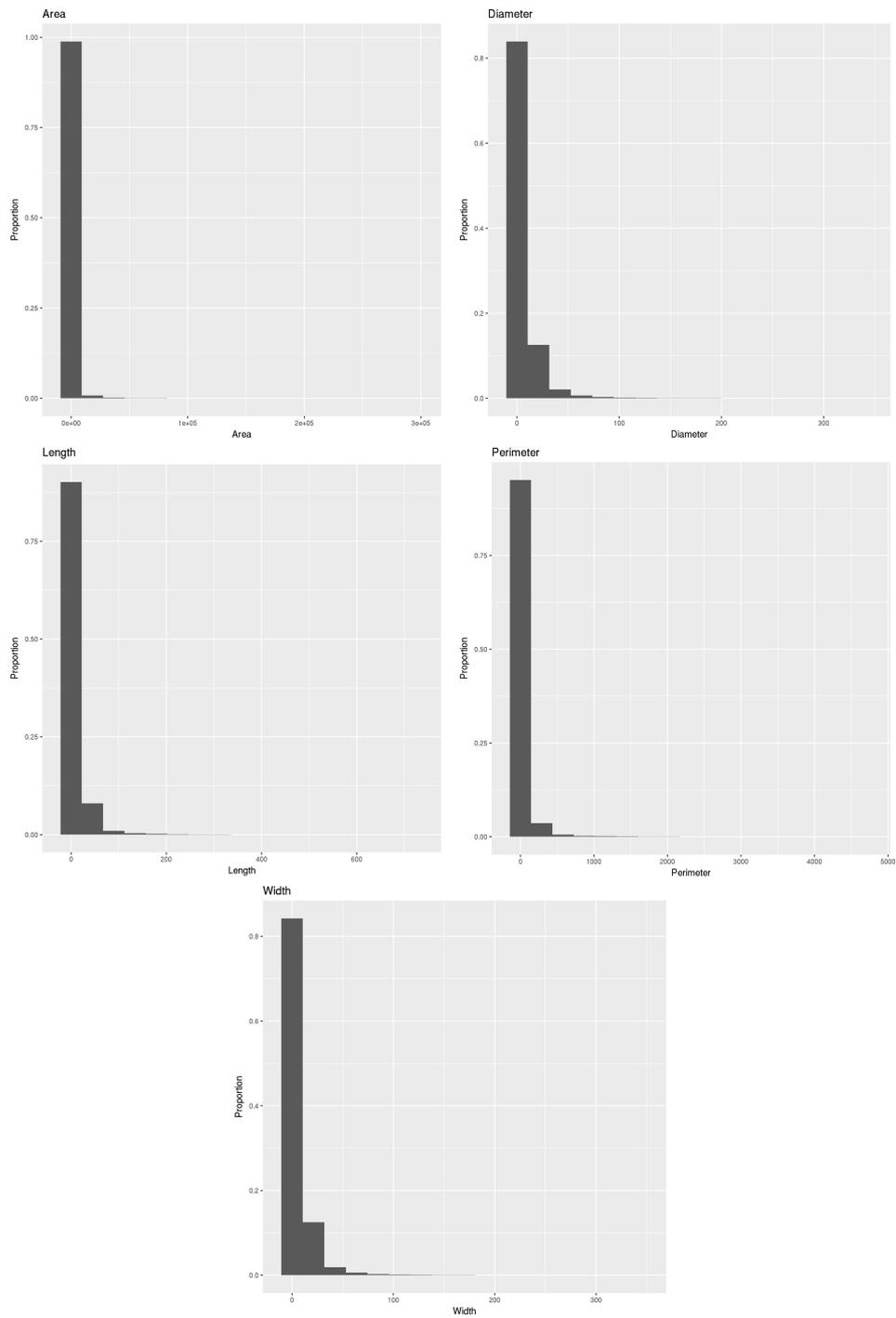


Figure 1.1: Histograms in natural scale for the Size properties

For the size variables, the properties are not well represented because the large variability of these measures. In figure 1.1 we can find the histograms for the five variables representing the size of the flocks. The first class of the histogram is grouping a large quantity of individuals and the graphic is not informative. In order to avoid this, the classes were constructed using a geometrical construction.

The variables representing the size of the flocks are transformed using a logarithmic transformation. This information is presented in the histograms in figure 1.2. We can observe that the behavior of all the variables is the same. The flocks are very asymmetrically distributed with positive skewness. This distribution is unimodal. The most of the flocks are grouped in the beginning of the range. There are a large percentage of individuals having small or average size. Also there exists a set of flocks having very large dimensions compared to those first described.

We complement the graphical analysis of the size variables with box and whiskers plots. Those plots are presented in the figure 1.3. In the box-plots we note that, for all the variables representing size, the 50% of the values between the 1st and the 3rd quantile are placed at the beginning of the range of the variables. The variability in this 50% of the flocks is small compared to the values out of the box. Also, there are large amount of flocks with elevated value for size measures. More than outliers, they behave like a different population of flocks. Those individuals are presented like blue dots in the graphics.

For the shape variables, we have constructed histograms and box-plots using natural scale for classes. The graphic representation is observed in figure 1.4. The shape variables take values between 0 and 1. We can observe in the case of Circularity, a negative skewness, where the measures are distributed more or less in the larger values of the property and they have an unimodal distribution. The Convexity and Solidity present a behavior very similar. Both distributions are asymmetric, with negative skewness, where the most of the flocks present values near to or equal to 1. The distribution is unimodal in both cases. The elongation presents a distribution asymmetric, with positive skewness, where the most part of the flocks are distributed more or less in the inferior part of the range.

In a similar manner, we complement the graphical analysis with box and whiskers plots. Those plots are presented in figure 1.5. In the case of Circularity, the box is placed oriented near to the larger values of the range. The size of the box is wide indicating large variability in that 50%. The behavior of the variables Convexity and Solidity is similar. In both cases, the box representing the 50% interquartile is placed at the bottom of the range. It means that, 75% of the individuals have values very close to 1. The

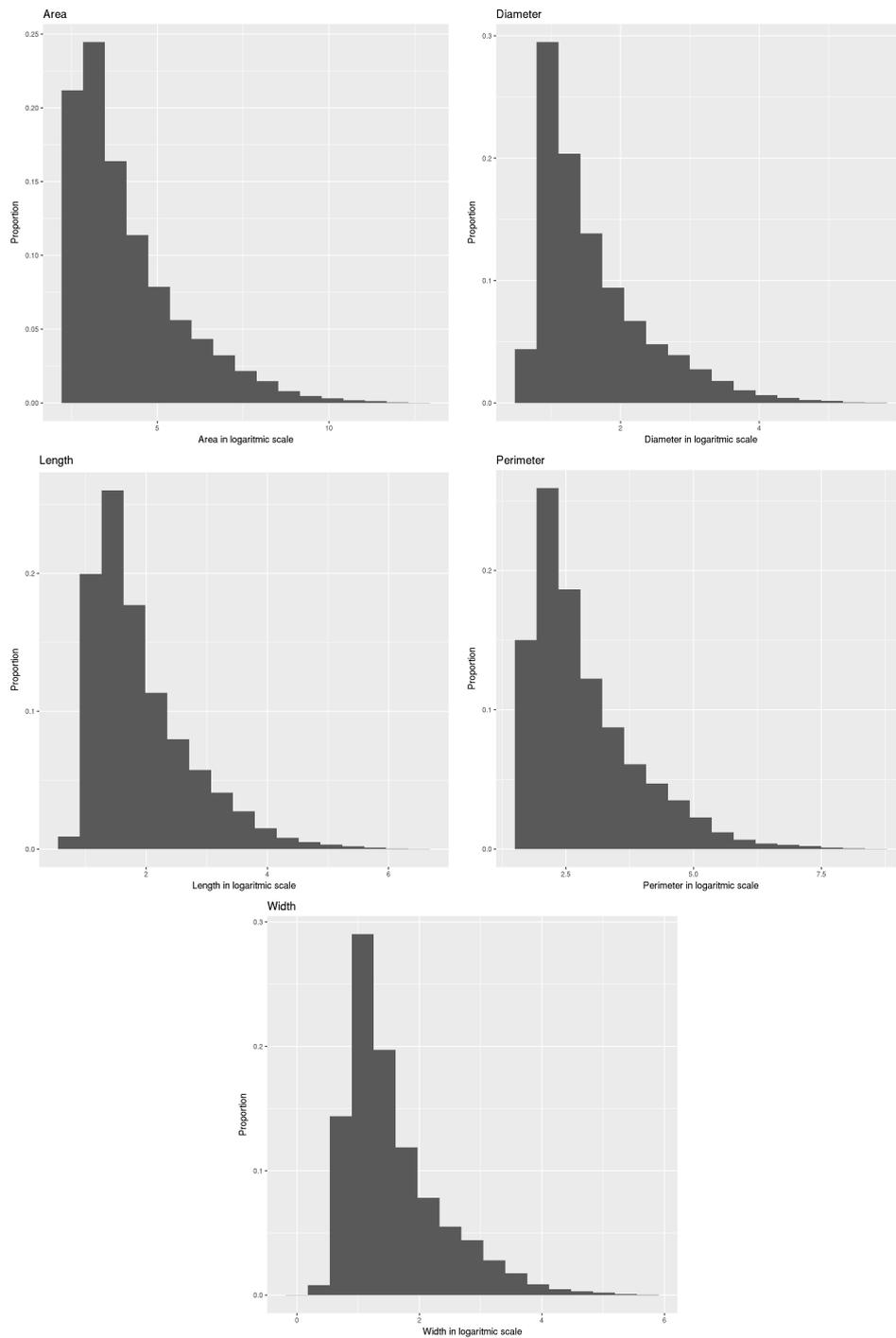


Figure 1.2: Histograms in logarithmic scale for the Size properties

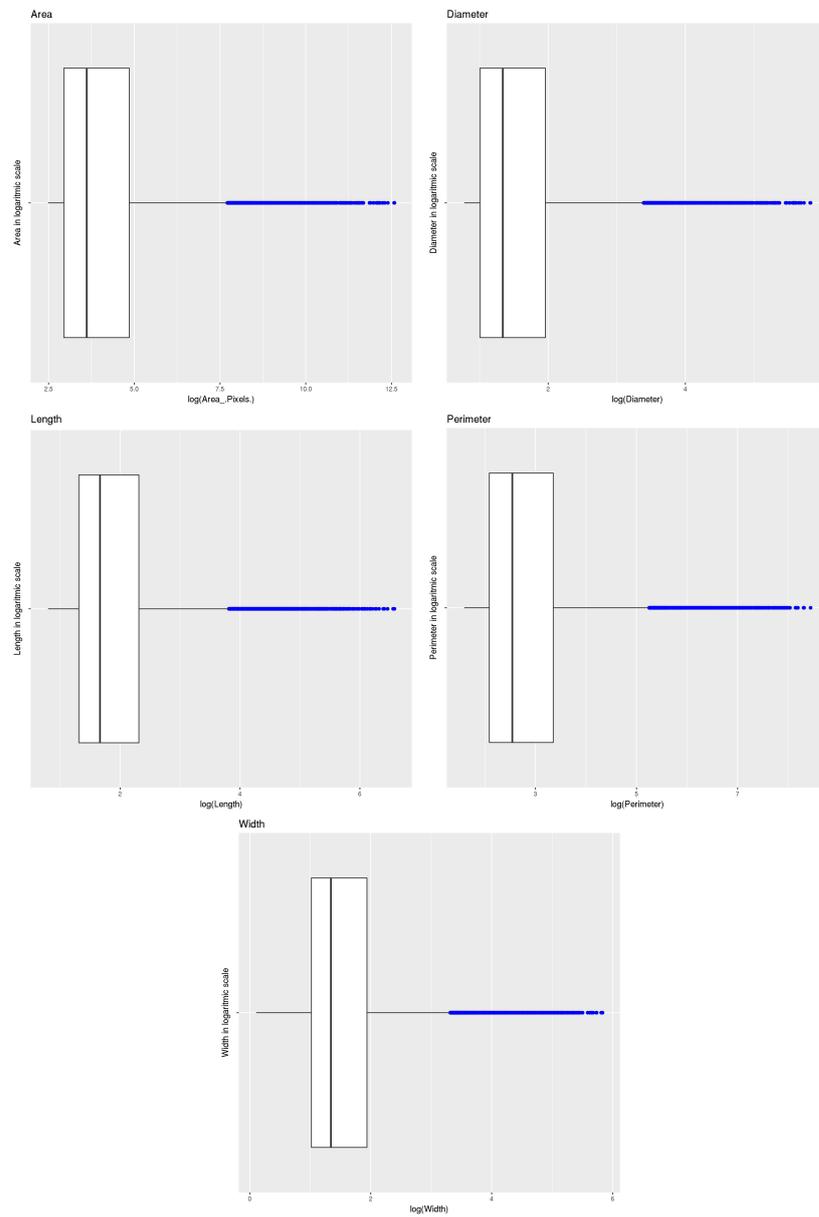


Figure 1.3: Box-plots in logarithmic scale for the Size properties

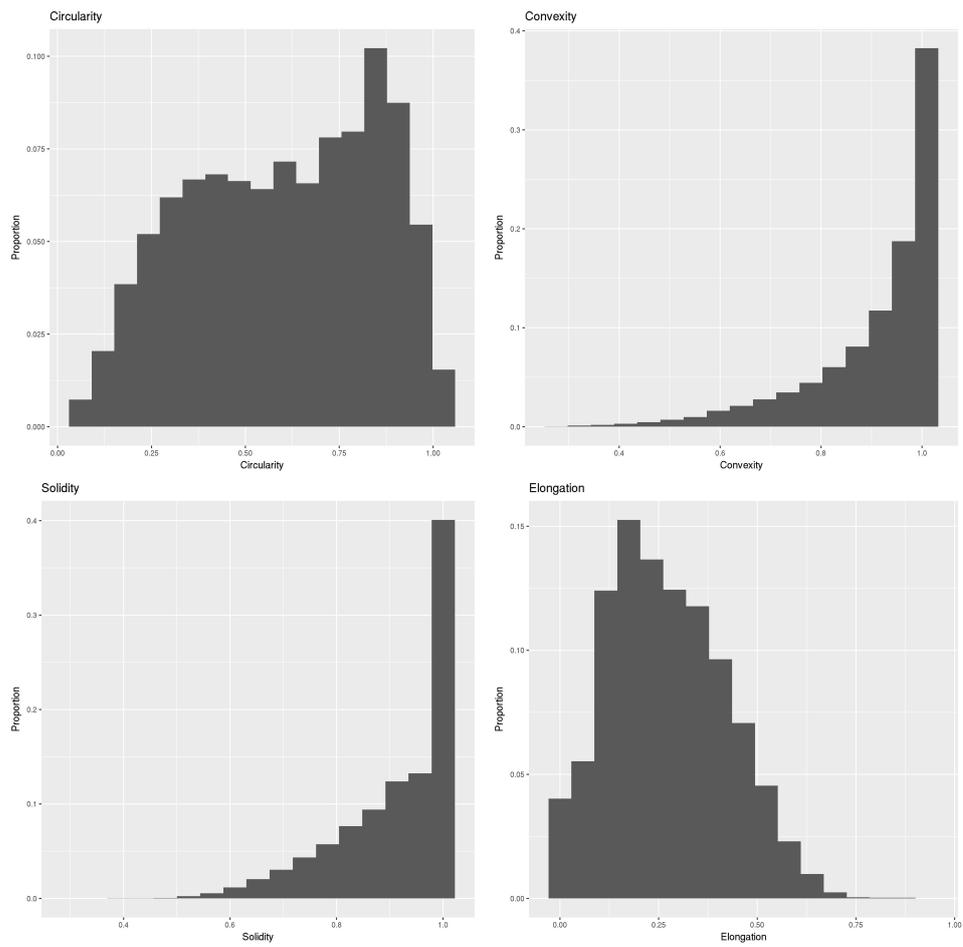


Figure 1.4: Histograms for the Shape properties

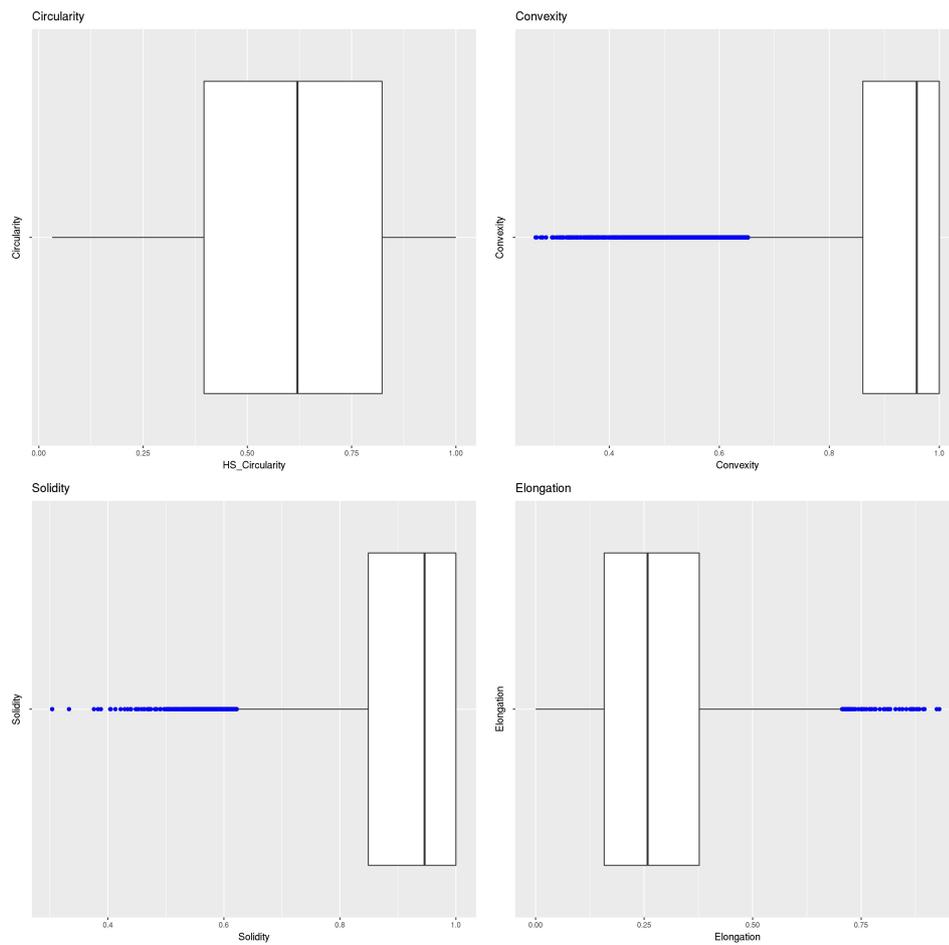


Figure 1.5: Box-plots for the Shape properties

variability in that 50% of the population is small. There are an important number of individuals with very small values of these two characteristics, indicating a very different shape between the population and those "outliers" which seem to have the behavior of another population.

Populations of Bentonite flocks under different hydrodynamic conditions. In this part, we analyze a set of measures from populations of particles obtained from different hydrodynamic conditions. Those hydrodynamic conditions are represented by the speed of mixing at 30, 50, 70, and 90 rpm. The data sets are compound by two populations obtained under the same hydrodynamic conditions and the same time of evolution of the process. Those data sets were obtained in order to evaluate the replicability of the results under the same experimental conditions using the Bentonite as material.

First, we will compare the univariate distributions of each pair of data sets, in order to see the differences in the experimental results. Then, the two datasets are gathered and used as the measures of one population. The different populations are compared with frequency polygons and box and whiskers for each variable, in order to see the differences between the measures of each population.

The variables measured in this case were:

- Representing particle's size: Area, Major Axis, Minor Axis, Perimeter.
- Representing particle's shape: Extent, Solidity, Aspect Ratio, Convexity, Circularity.

Population under 30 rpm. Following the univariate analysis descriptive of the data sets, we are going to compare the frequency polygons of the variables representing particle's size, under a logarithmic transformation, in order to compare the distribution of the two data sets one variable at time.

We can observe in figure 2.2 that the frequency polygons corresponding each data set are almost overlapped. It indicate that, in terms of the size variables point of view, the experimental results are replicable. The distributions of the population of flocks have, a similar behavior for the variables Area, Perimeter and Minor axis. The distribution presents positive skewness, with large amount of flocks with small size. Then, there are groups of flocks with more important size at the end of the distribution. In the case of the variable Mayor axis, the distribution presents a large amount of flocks of small size, but we can observe the presence of flocks of medium size. The distribution is almost bimodal because of these groups of flocks. In similar

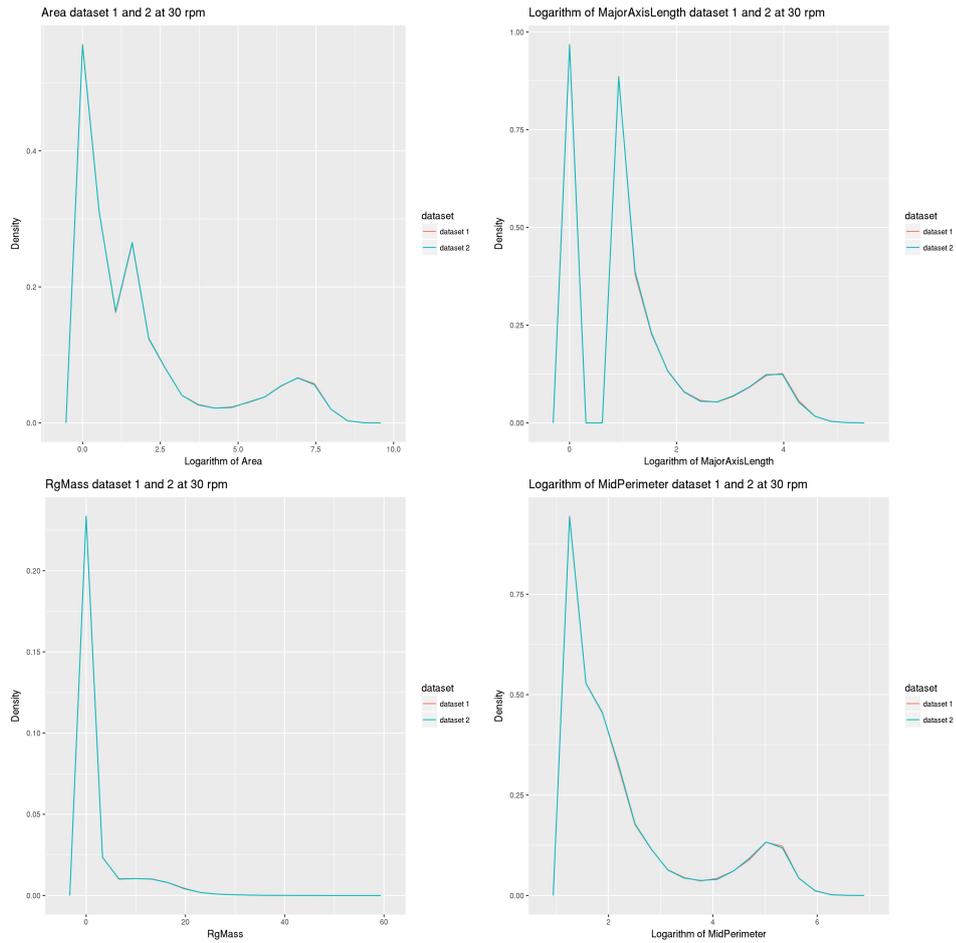


Figure 1.6: Frequency polygons comparing the two data sets obtained under 30 rpm. Each size variable is represented under a logarithmic transformation. The aim is to evaluate the replicability of the experimental results. The two distributions are almost overlapped

way, there are groups of flocks with more important size at the end of the distribution.

In a similar way, we are going to represent the frequency polygons for the shape. In the graphic, the two data sets obtained at 30 rpm are compared one variable at time.

For the variables representing the particle's shape, the frequency polygons of the two data sets are almost overlapped in each case, and are shown in figure 2.3. For the Circularity, the measures of the flocks are more concentrated in the values close to 1, indicating a geometric shape rounded but not as a perfect circle. The frequency polygons of the Aspect ratio shows that the most of the flocks takes 1 as value of this characteristic, indicating that those flocks have a contour or surface very smooth and with shape very rounded. Although, the presence of a group of flocks in the center of the distribution shows the presence of flocks with surface more irregular and shape more elongated. The distributions shown in the case of Convexity and Solidity are very similar. The most of the flocks take values close to or equal to 1. It indicates that the flocks do not present concavities or roughness important. A less important group of flocks are found more at the center of the distribution, showing the presence of concavities in some flocks.

Population under 50 rpm. In a similar way, we are going to compare the frequency polygons of the variables representing particle's size, under a logarithmic transformation, in order to compare the distribution of the two data sets one variable at time. These results were obtained under a speed of mixing of 50 rpm.

As before, we can observe in figure 1.8 that the frequency polygons corresponding each data set are almost overlapped. It indicate that, in terms of the size variables point of view, the experimental results are repicable. The distributions of the population of flocks have, a similar behavior for the variables Area, Perimeter and Minor axis. The distribution presents positive skewness, with large amount of flocks with small size. Then, there are groups of flocks with more important size, but in this case, the groups are more close to the group of small particles, in contrast with the case of 30 rpm. In the case of the variable Mayor axis, the distribution presents a large amount of flocks of small size, but we can observe the presence of flocks of medium size. The distribution is almost bimodal because of these groups of flocks. In similar way, there are groups of flocks with more important size, but again, this group of particles are closer to the small ones.

Now, we represent the frequency polygons for the shape. In the graphic, the two data sets obtained at 50 rpm are compared one variable at time.

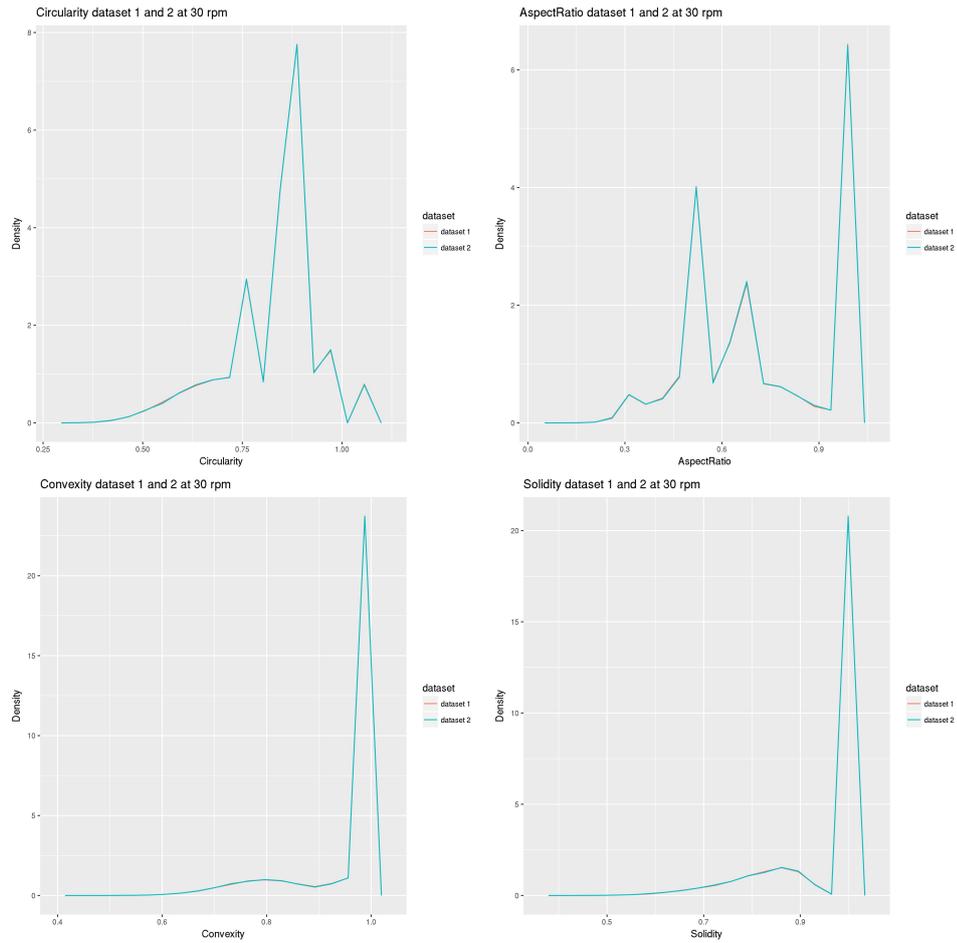


Figure 1.7: Frequency polygons comparing the two data sets obtained under 30 rpm. Each shape variable is represented. The aim is to evaluate the replicability of the experimental results. In this case, the two distributions are almost overlapped as in the precedent set of variables.

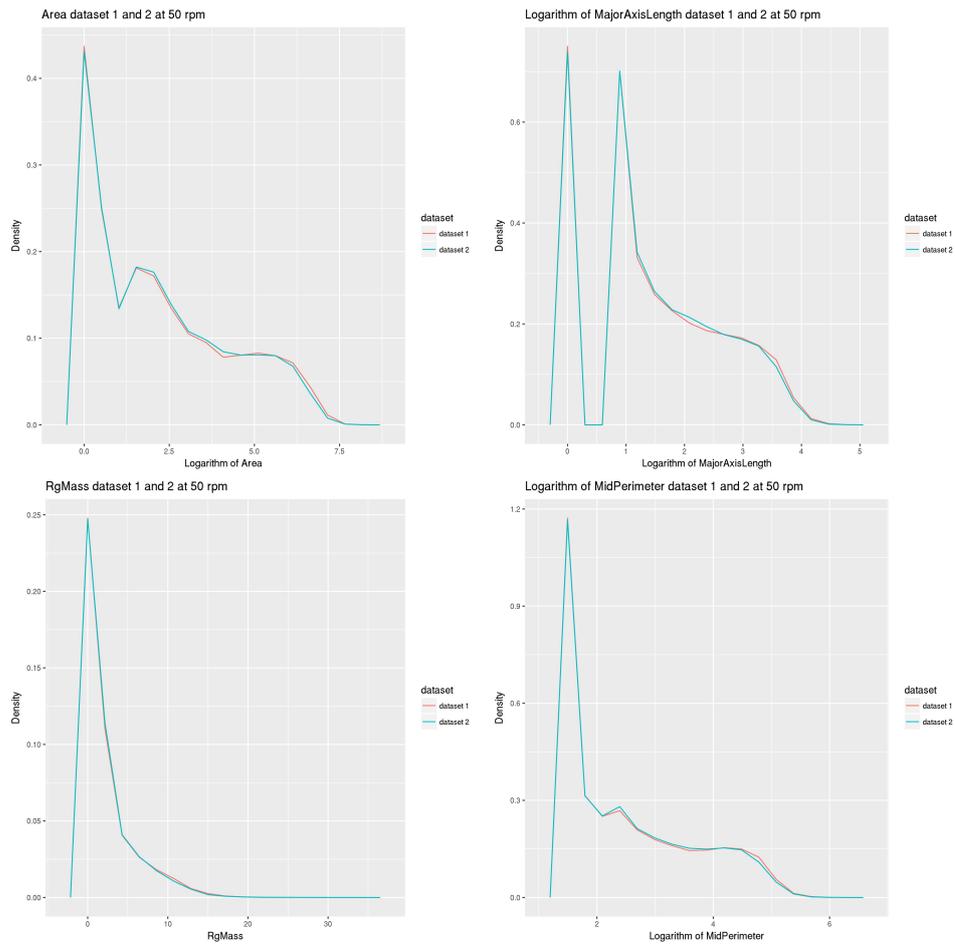


Figure 1.8: Frequency polygons comparing the two data sets obtained under 50 rpm. Each size variable is represented under a logarithmic transformation. The aim is to evaluate the replicability of the experimental results. Again, the two distributions are almost overlapped

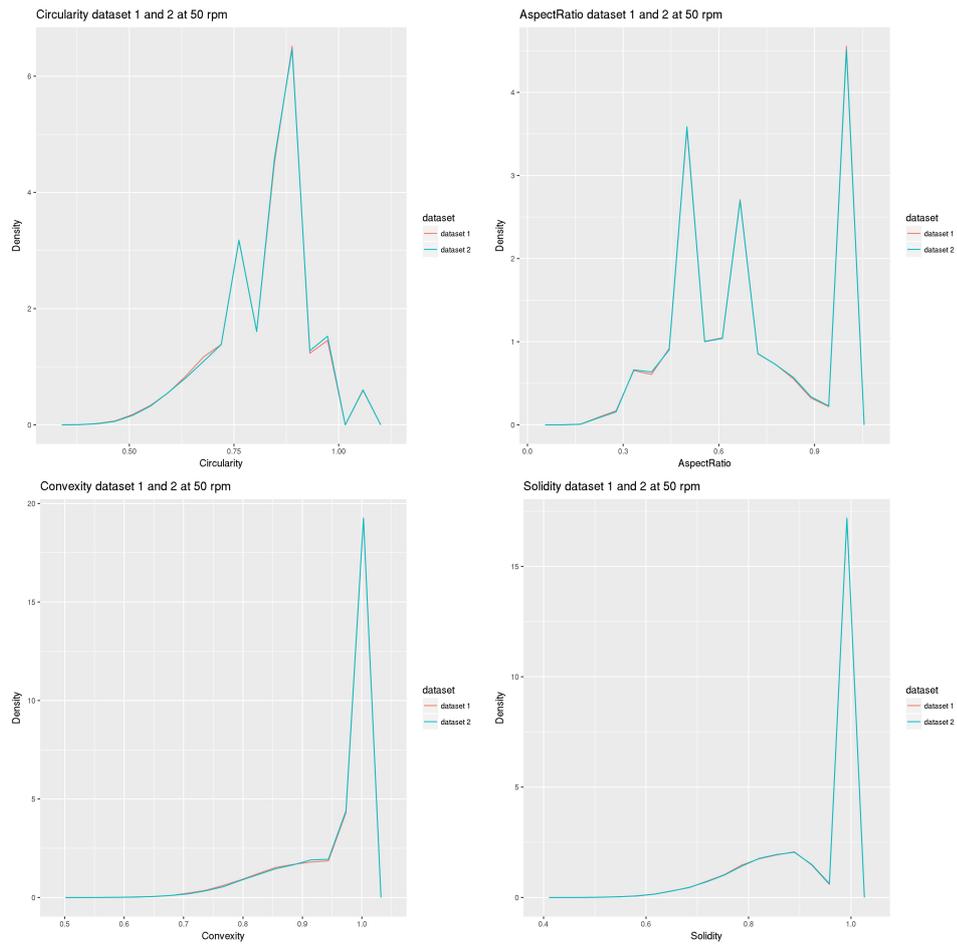


Figure 1.9: Frequency polygons comparing the two data sets obtained under 50 rpm. Each shape variable is represented. The aim is to evaluate the replicability of the experimental results. As before, the two distributions are almost overlapped as in the precedent set of variables.

For the variables representing the particle's shape, the frequency polygons of the two data sets are almost overlapped in each case, and are shown in figure 1.9. For the Circularity, the measures of the flocks are more concentrated in the values close to 1, indicating a geometric shape rounded but not as a perfect circle. The frequency polygons of the Aspect ratio shows that the most of the flocks takes 1 as value of this characteristic, indicating that those flocks have a contour or surface very smooth and with shape very rounded. Although, the presence of an important group of flocks in the center of the distribution shows the presence of flocks with surface more irregular and shape more elongated. The distributions shown in the case of Convexity and Solidity are very similar. The most of the flocks take values close to or equal to 1. It indicates that the flocks do not present concavities or roughness important. A less important group of flocks are found more at the center of the distribution, showing the presence of concavities in some flocks. These groups of flocks are closer to the value of 1 than the precedent case.

Population under 70 rpm. In a similar manner, we compare the frequency polygons of the variables representing particle's size, under a logarithmic transformation, in order to compare the distribution of the two data sets one variable at time. These results were obtained under a speed of mixing of 70 rpm.

As before, we can observe in figure 1.10 that the frequency polygons corresponding each data set are almost overlapped. It indicate that, in terms of the size variables point of view, the experimental results are repicable. The distributions of the population of flocks have, a similar behavior for the variables Area, Perimeter and Minor axis. The distribution presents positive skewness, with large amount of flocks with small size. Then, there are groups of flocks with more important size, but in this case, the groups are more close to the group of small particles, in contrast with the cases precedents. In the case of the variable Mayor axis, the distribution presents a large amount of flocks of small size, but we can observe the presence of flocks of medium size. The distribution is almost bimodal because of these groups of flocks. In similar way, there are groups of flocks with more important size, but again, this group of particles are closer to the small ones.

Now, we represent the frequency polygons for the shape. In the graphic, the two data sets obtained at 70 rpm are compared one variable at time.

For the variables representing the particle's shape, the frequency polygons of the two data sets are almost overlapped in each case, and are shown in figure 1.11. For the Circularity, the measures of the flocks are more con-

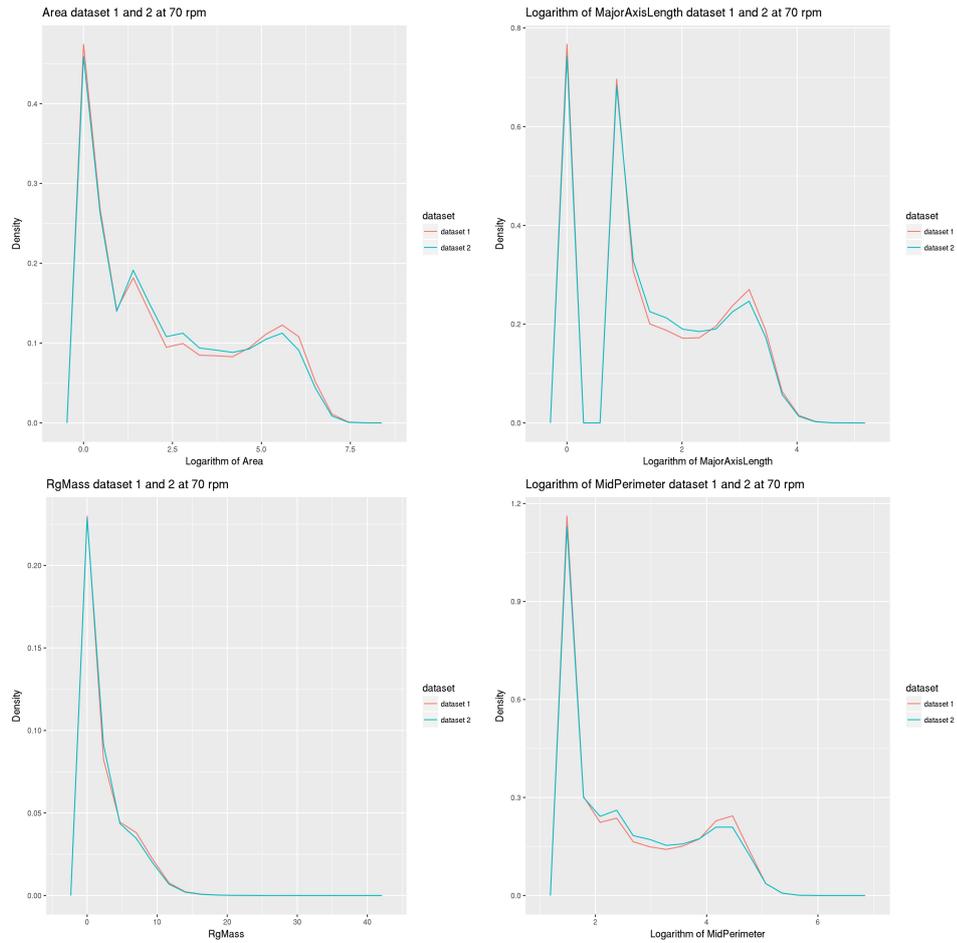


Figure 1.10: Frequency polygons comparing the two data sets obtained under 70 rpm. Each size variable is represented under a logarithmic transformation. The aim is to evaluate the replicability of the experimental results. Again, the two distributions are almost overlapped

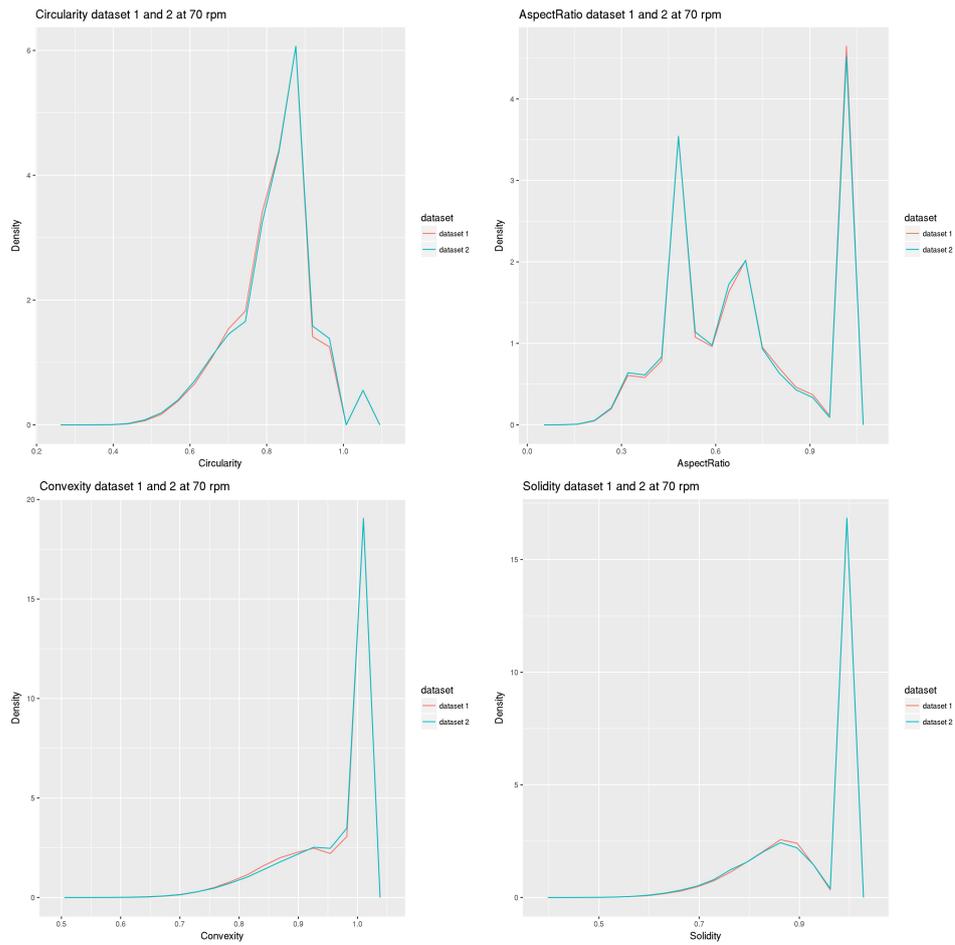


Figure 1.11: Frequency polygons comparing the two data sets obtained under 70 rpm. Each shape variable is represented. The aim is to evaluate the replicability of the experimental results. As before, the two distributions are almost overlapped as in the precedent set of variables.

concentrated in the values close to 1, indicating a geometric shape rounded but not as a perfect circle. The frequency polygons of the Aspect ratio shows that the most of the flocks takes 1 as value of this characteristic, indicating that those flocks have a contour or surface very smooth and with shape very rounded. Although, the presence of an important group of flocks in the center of the distribution shows the presence of flocks with surface more irregular and shape more elongated. The distributions shown in the case of Convexity and Solidity are very similar. The most of the flocks take values close to or equal to 1. It indicates that the flocks do not present concavities or roughness important. A less important group of flocks are found more at the center of the distribution, showing the presence of concavities in some flocks. These groups of flocks are closer to the value of 1 than the precedent case.

Population under 90 rpm. Finally, we compare the frequency polygons of the variables representing particle's size, under a logarithmic transformation, in order to compare the distribution of the two data sets one variable at time. These results were obtained under a speed of mixing of 90 rpm.

In a similar way, we can observe in figure 1.12 that the frequency polygons corresponding each data set are almost overlapped. It indicate that, in terms of the size variables point of view, the experimental results are repicable. The distributions of the population of flocks have, a similar behavior for the variables Area, Perimeter and Minor axis. The distribution presents positive skewness, with large amount of flocks with small size. Although, there are groups of flocks with more important size, but in this case, the groups are more close to the group of small particles and is more important in number, in contrast with the cases precedent. In the case of the variable Mayor axis, the distribution presents a large amount of flocks of small size, but we can observe the presence of an important number of flocks of medium size. In similar way, there are groups of flocks with more important size, but again, this group of particles are closer to the small ones.

Now, we represent the frequency polygons for the shape. In the graphic, the two data sets obtained at 90 rpm are compared one variable at time.

For the variables representing the particle's shape, the frequency polygons of the two data sets are almost overlapped in each case, and are shown in figure 1.13. For the Circularity, the measures of the flocks are more concentrated in the values close to 1, indicating a geometric shape rounded but not as a perfect circle. The frequency polygons of the Aspect ratio shows that the most of the flocks takes 1 as value of this characteristic, indicating

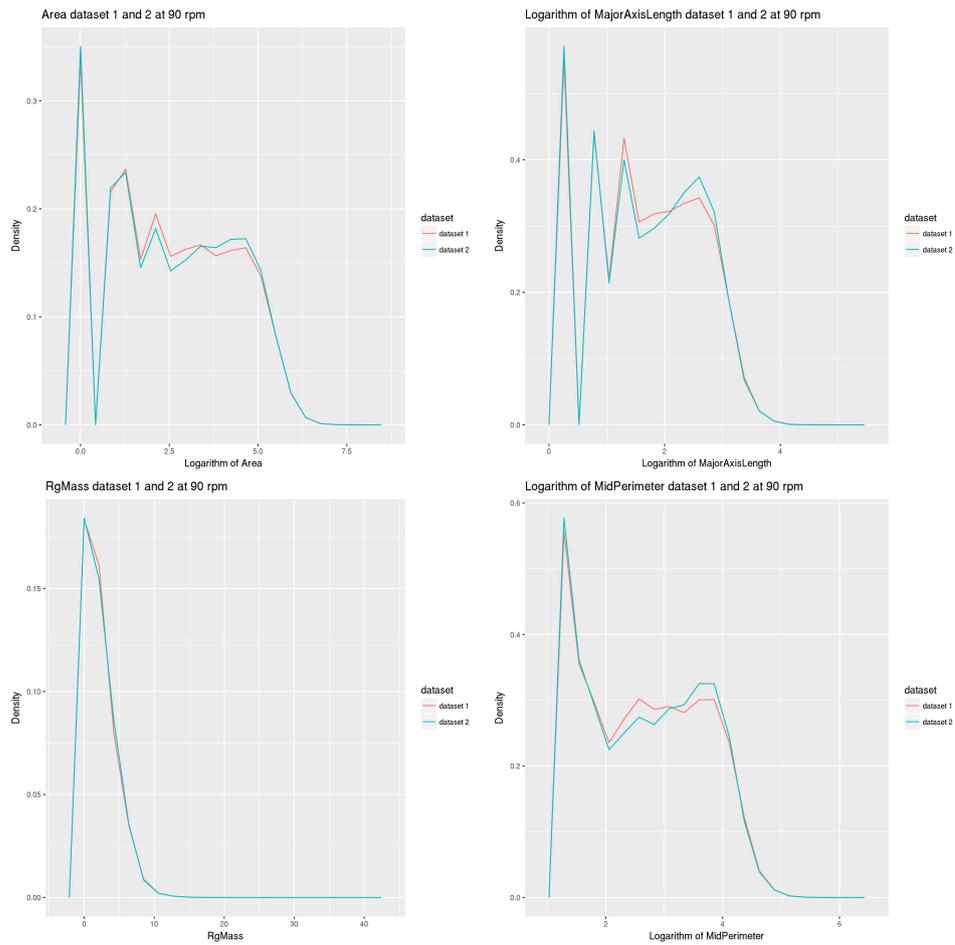


Figure 1.12: Frequency polygons comparing the two data sets obtained under 90 rpm. Each size variable is represented under a logarithmic transformation. The aim is to evaluate the replicability of the experimental results. Again, the two distributions are almost overlapped

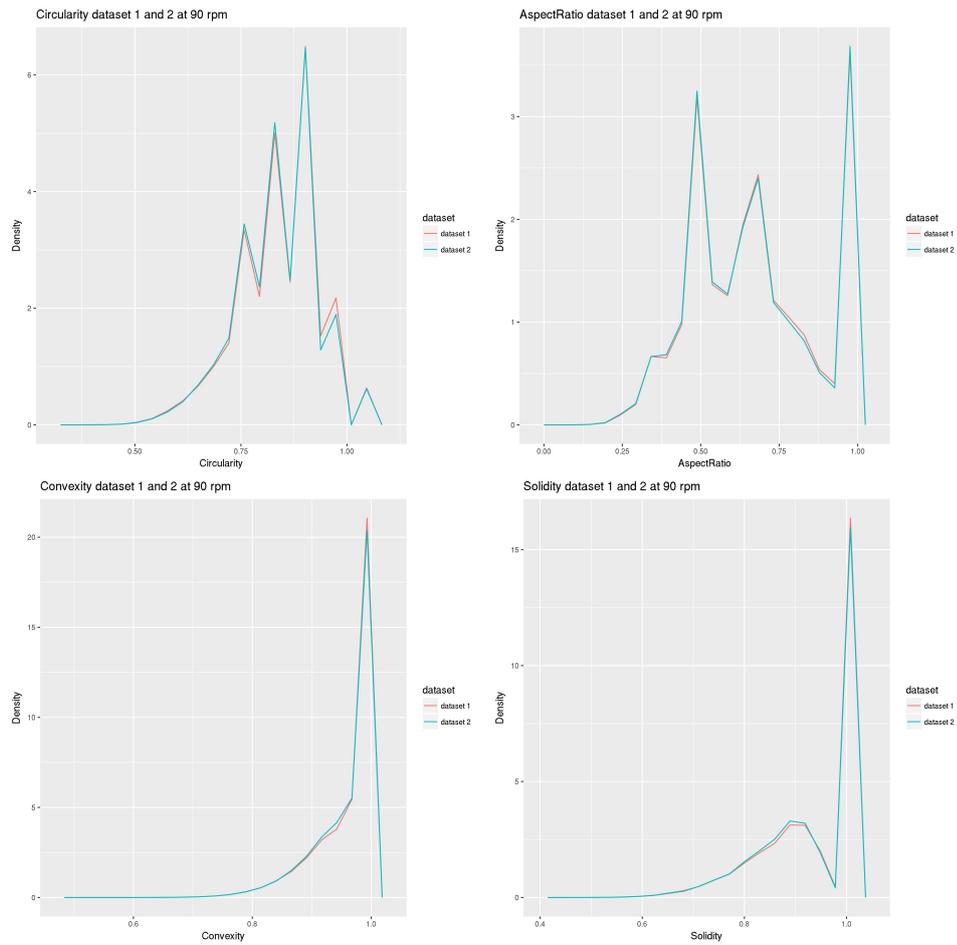


Figure 1.13: Frequency polygons comparing the two data sets obtained under 90 rpm. Each shape variable is represented. The aim is to evaluate the replicability of the experimental results. As before, the two distributions are almost overlapped as in the precedent set of variables.

that those flocks have a contour or surface very smooth and with shape very rounded. Although, the presence of an important group of flocks in the center of the distribution shows the presence of flocks with surface more irregular and shape more elongated. This group is more important in number than in the cases precedent. The distributions shown in the case of Convexity and Solidity are very similar. The most of the flocks take values close to or equal to 1. It indicates that the flocks do not present concavities or roughness important. A less important group of flocks are found more at the center of the distribution, showing the presence of concavities in some flocks. These groups of flocks are closer to the value of 1 than the precedent case.

Comparison of the univariate distributions under different levels of rpm. As we could observe in the precedent part, the two set of data are similar in all the variables, so they are going to form a single data set for each level of rpm. We are going to use again the frequency polygons and the box and whiskers plot for comparing the distributions of each variable under different hydrodynamics conditions.

We present first the variables representing the particle's size and then the variables representing the particle's shape.

In figure 1.14, we find the frequency polygons in the same graphic for the populations of particles obtained under different levels of rpm. The representation was done one variable at time. Here, for the size variables, we use a logarithmic transformation in order to better describe the data. We observed a positive skewness in all four variables. In general, there are an important number of flocks with small size. There are also a group of flocks of medium size in the center of the distribution. we observed an increase of the size of the flocks when the level of rpm increase. The group of small flocks also decrease when the rpm increase. This tendency is observed in all size variables.

If we compare the distribution of the size variables with the box and whiskers plot, like in figure 1.15, we observe that the 50% inter-quantile is always placed at the low values of the distribution. We can see also an augmentation in the variability of the measures when the level of rpm augment. The presence of large flocks in the lowest level of rpm (30 rpm) is remarked. The median of the distribution augment when the level of rpm does. Also, the number of flocks with large size augment when the speed of mixing does.

In a similar manner, the distributions of the variables representing the particle's shape are shown using frequency polygons in figure 1.16. In the graphs, we observe that the measures of Circularity of flocks tend to the values close to 1. It indicates that the geometric shape of the flocks tend to

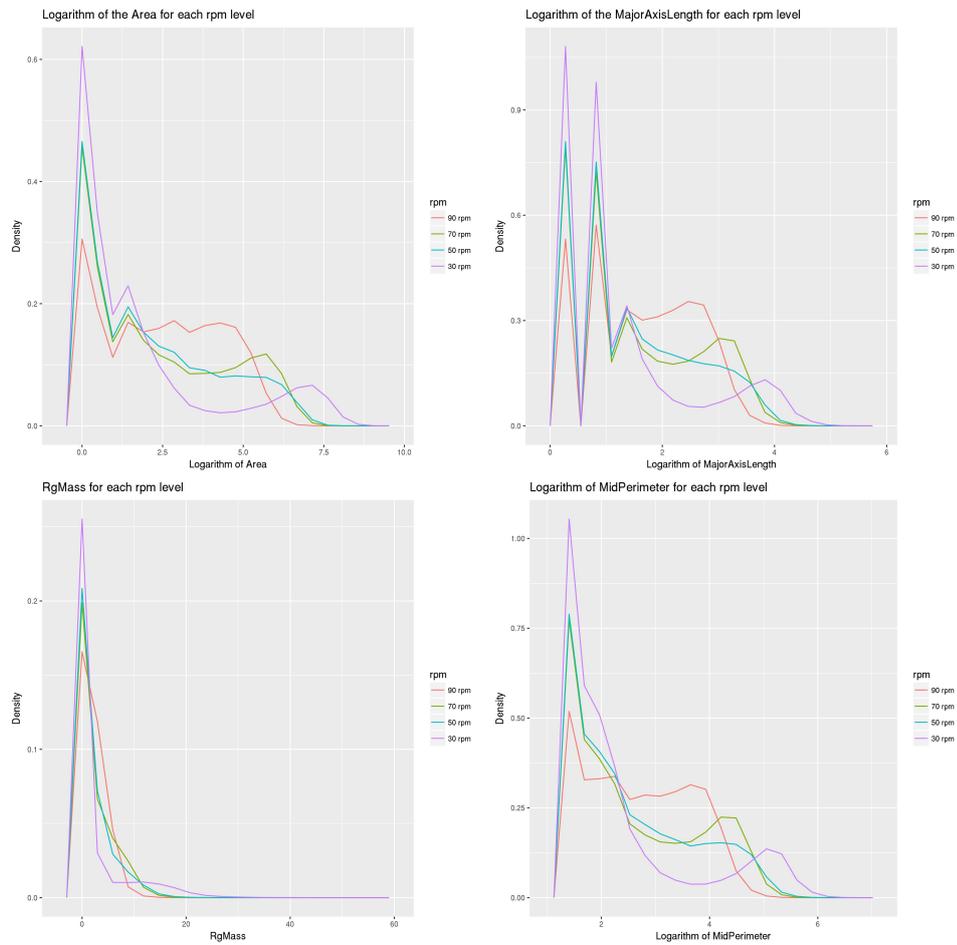


Figure 1.14: Frequency polygons comparing the populations obtained under 30, 50, 70, and 90 rpm. Each size variable is represented under a logarithmic transformation. In general, we observed an increase of the size of the flocks when the level of rpm increase. The group of small flocks also decrease when the rpm increase. This tendency is observed in all size variables

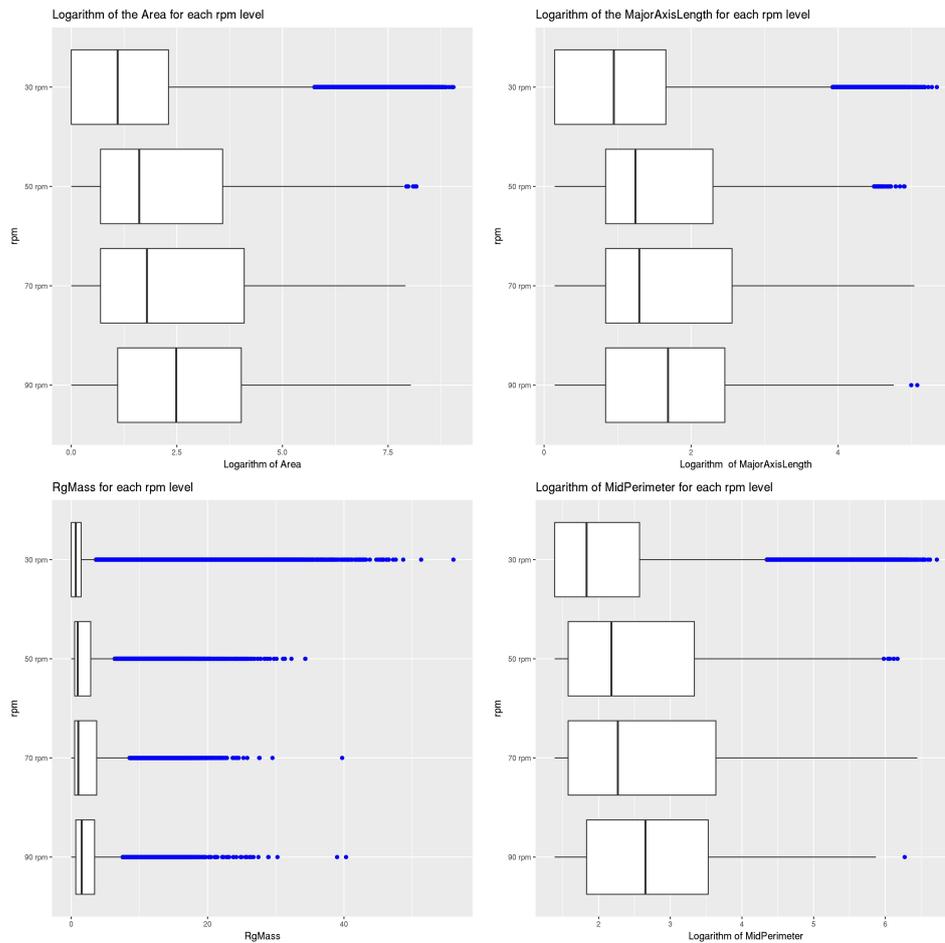


Figure 1.15: Box and whiskers comparing the populations obtained under 30, 50, 70, and 90 rpm. Each size variable is represented under a logarithmic transformation. We observe that the median augment when the level of rpm does. Also, the number of flocks with large size augment when the speed of mixing does.

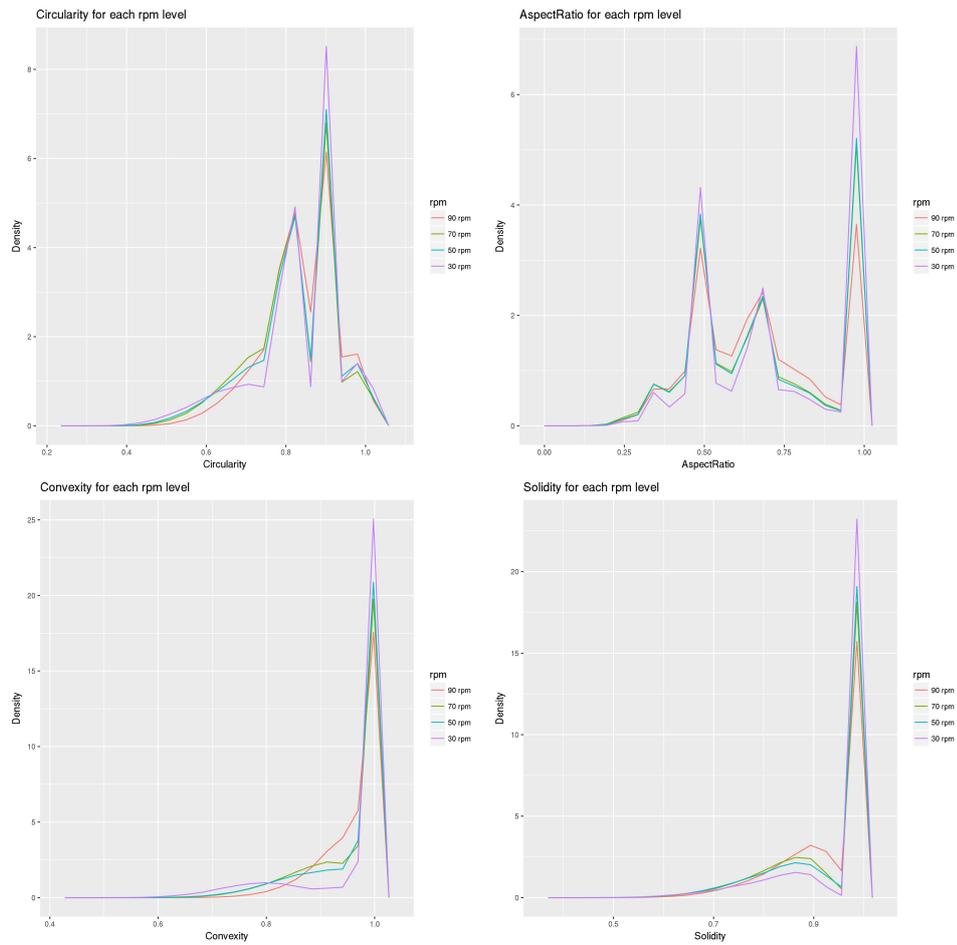


Figure 1.16: Frequency polygons comparing the populations obtained under 30, 50, 70, and 90 rpm. Each shape variable is represented. The roundness of the flocks decrease when the level of rpm augment. Also, lower is the level of rpm, the more elongated are the flocks. When the level of rpm augment, more irregular are the flocks in terms of concavity and roughness.

be rounded but they do not achieve the perfect sphere shape. The variability seem to decrease when the level of rpm augment. For the Aspect ratio, we find a large number of particles having a measure close to 1 for this characteristic. Also, there are groups of flocks taking values at the center of the distribution of this variable. It represents groups of flocks with an elongated form. Lower is the level of rpm, the more elongated are the flocks. In the case of Convexity and Solidity, the most of the flocks have values very close to 1 for both measures. It means that the flocks do not present important concavities of roughness. Although, there are a group of flocks taking values near to 0.9, representing flocks more irregular in shape. When the level of rpm augment, more irregular are the flocks in terms of concavity and roughness.

In the figure 1.17, the distribution of the variables representing the particle's shape are shown using box and whiskers plots. In the case of Circularity, the boxes representing the 50% inter-quantile, have a similar behavior in central tendency and in variability. All distributions present negative skewness, and the predominant values in the populations are close to 1. However, the quantity of individuals marked as outliers is larger at 30 rpm than the others levels of speed of mixing. It seems that flocks with less rounded shape are present in populations obtained under lower levels of rpm. About Aspect ratio, the distributions present negative skewness too, and the values are placed at the end of the distribution. The variability in this variable decrease when the level of rpm augment. For the level of 90 rpm, the distribution of the measures are more separated from 1 than at lower levels of rpm. For Convexity and Solidity, the behavior is of the distributions are similar in both cases. The distributions present a marked negative skewness and the most of the individuals are concentrated close to 1. The presence of outliers with low values of these variables is bigger when the level of rpm decrease. For the shape variables Circularity, Convexity and Solidity, the groups of individuals marked as outliers behaves like another population in terms of those variables.

Populations of Bentonite flocks under different hydrodynamic conditions and through time. In this part, we analyse a set of measures from populations of particles obtained from different hydrodynamic conditions. Those hydrodynamic conditions are represented by the speed of mixing at 30, 50, 70, 90, and 100 rpm. The data sets are compound of the measures of size and shape characteristics from populations of particles obtained under the same speed of mixing but at different instants of time.

We are going to analyze the distribution of each variable, comparing the

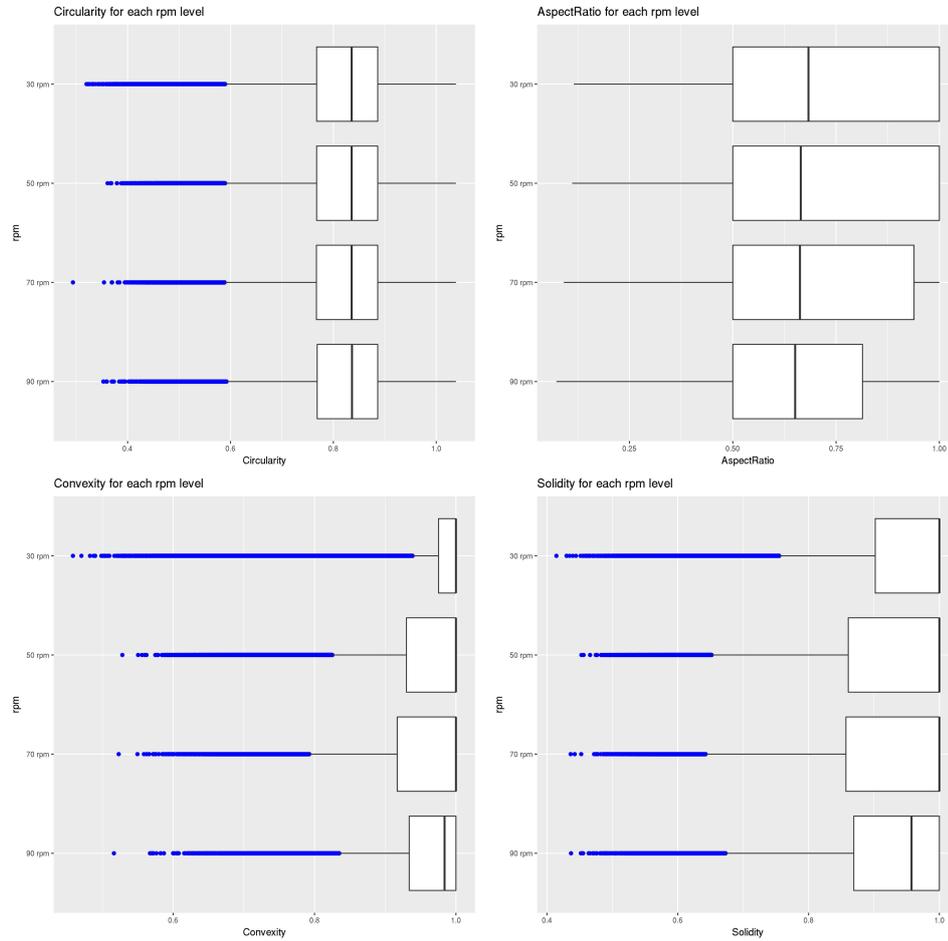


Figure 1.17: Box and whiskers comparing the populations obtained under 30, 50, 70, and 90 rpm. Each Shape variable is represented. For Circularity, the quantity of individuals marked as outliers is larger at 30 rpm than the others levels of speed of mixing. About Aspect ratio, the distribution of the measures are more separated from 1 than at lower levels of rpm. The presence of outliers with low values of these variables is bigger when the level of rpm decrease. For the shape variables Circularity, Convexity and Solidity, the groups of individuals marked as outliers behaves like another population in terms of those variables

frequency polygons of populations obtained at different instant of time. Also, we are going to compare the distributions using box and whiskers plot. The analysis is performed at each level of speed of mixing.

Population under 30 rpm.

Exploratory multivariate statistics

The Principal Components Analysis (PCA). The Principal Components Analysis or PCA has as main objective the reduction of dimension of data. That is, to describe precisely the values of p variables by a smaller subset $r < p$ of them, losing a little amount of information. Given n observations of p variables the PCA analyses if it is possible to represent accurately this information with a smaller number of variables build as a linear combination of the original ones.

This technique lets to represent optimally observations of an p -dimensional space in an space of low dimension. Further, it allows to transform the original variables, generally correlated, into new variables, making easier the interpretation of data.

This technique can be seen through different perspectives. The first one is the descriptive approach consisting of finding a subspace of dimension less than p , so that when we project the observations, they conserve the essential of its structure with the least possible distortion. Specifically, we search to conserve the relative distance in the original space between the observations. In an statistical approach, the technique looks for the observation projection verifying this property, minimizing the orthogonal distances, so that the new variables generating the new subspace are orthogonal. It also can be seen geometrically as if the data points cloud are in an ellipsoid where the best approximation is provided by its projection on the mayor axes of the ellipsoid. This is equivalent to minimizing the orthogonal distance.

Computation of the Principal Components. It can be shown that the r -dimensional space that better represents the observations is defined as the eigenvectors associated to the r largest eigenvalues of the Variance and Covariance matrix S . Those directions are named Principal Directions and they define the new variables named *Principal Components*.

Let's denote data matrix X defined by n observations of p variables. Generally, the matrix X has rang equals to p (and in consequence, also the matrix S). Then, there exists as many Principal Components as variables. This Principal Components are calculate using the characteristic roots $\lambda_1, \dots, \lambda_p$

of S through:

$$|S - \lambda I| = 0 \quad (1.18)$$

and they associated vectors are

$$(S - \lambda_i I) a_i = 0 \quad (1.19)$$

where the terms λ_i are positive reals, because the matrix S is symmetric and definite positive. If we have $p - r$ variables as linear combination of the rest variables, the matrix S is semi definite positive and we have r real positive characteristic roots and the other $p - r$ roots are zero.

Properties of the Principal Components. The Principal Components are new variables with the following properties:

- They conserve the initial variability. The sum of the variances of the Principal Components is equal to the sum of variances of the original variables, and the generalized variance of the Components is also equal to the generalized variance of the original variables.
- The proportion of explicated variance by a component is the quotient between its variance, the associated eigenvalue and the sum of the eigenvalues of S .
- The covariances between each Principal Component and one original variable X_i is given by the product of the coordinates of the eigenvector multiplied by its eigenvalue.
- The correlations between a Component and one original variable X_i is proportional to the coefficient of that variable in the definition of the Component and the standard deviation of the variable.

$$\text{Corr}(Z_i, X_j) = \frac{\lambda_i a_{ij}}{\sqrt{\lambda_i S_j^2}} = \frac{a_{ij} \sqrt{\lambda_i}}{S_j}.$$

- The r Principal Components provide the optimal linear prediction with r variables of the original set X .
- If we standardize the Principal Components, we obtain the multivariate standardization of the original data set.

CPA analysis of the Correlation matrix. Since the Principal Components are obtained maximizing the variance of the projection, if some variable has variance so much large than the others, the first Principal Component tends to coincide to that variable. When the variables are expressed in different measure units this property is not desirable. If we want to avoid that problem, it is convenient to standardize the variables before the CPA. In this way, we obtain the normed Principal Components or equivalently, we can obtain the normed Principal Components computing the eigenvalues and eigenvectors of the correlation matrix R .

The properties of the Principal Components extract from R are:

- $\sum_i \lambda_i^R = \text{trace}(R) = p$.
- The proportion of explained variability by each Component is

$$\frac{\lambda_i^R}{p}.$$

- The correlation between each Component Z_j and the original variable X_i is $a'_j \sqrt{\lambda_j}$, where $Z_j = X a_j$.

Interpretation of the Principal Components. If the variables have large positive correlations, the first Component is a weighted average of all variables. This kind of situation leads to an interpretation of the first Component as a "Size" Component. The others Components oppose groups of variables and they are often interpreted as "Shape" Components.

When the original variables are transformed using a logarithmic transformation, the Components can be represented often as ratios of geometric means of variables.

Selection of the number of Principal Components r Several approach can be regarded when we are going to select the dimension of the representation r . Those approaches include to select the number of Components explaining a required proportion of variance or to reject those Components which eigenvalue associated is less than a fixed quote, this quote is usually taken as the mean of the eigenvalues.

Graphic representation. In order to represent the individuals, often we use a projection into an space of dimension 2. This projection is directly calculate as the value of the Principal Component using the eigenvectors.

To represent the original variables, we use the correlation coefficient between them and the Principal Components as coordinates. The vector of

correlations between the first component and the original variables is given by

$$\sqrt{\lambda_1} a_1' D$$

where D is a diagonal matrix having in the principal diagonal the inverses of the standard deviations of each variable.

Selection of relevant size and shape properties using PCA

In order to properly choose one size variable and one shape variable to monitor, Principal Components Analysis (PCA) was used. The PCA performs an orthogonal transformation that provides a new set of uncorrelated variables called Principal Components (PC). These variables are the eigenvalues of the correlation matrix. The first principal component is the one that explains most of the variability of the original system.

A graphical interpretation of the principal components is obtained by plotting the circle of correlation, where the abscissa and ordinate of each point are the correlation coefficients between one property and respectively the first and second principal component. The correlation coefficient of the i^{th} PC and the j^{th} variable is

$$r(PC_i, j) = \sqrt{PC_i} V_i(j),$$

where V_i is the i^{th} eigenvector.

PCA performed with all the size properties and shape properties confirm that these two groups of variables are naturally related. Then, a PCA analysis was performed with the size variables as one group and shape variables as another group. Perimeter was the most important property related to the first PC and after that the *Rg*. Circularity and Convexity are the properties more related to the second PC.

Outlier identification in high dimension In the data coming from population of particles, we can find group of individuals that behave atypically. This kind of populations are formed by a large number of individual and also can be observed a large variation in the particles properties.

In order to study the presence of outlier individuals from a multivariate perspective, we are going to use a computational procedure proposed by (Filzmoser, et al. 2008) ([FMW08]) and the package *mvoutlier* ([FG17]) in the R language ([R C17]). This algorithm utilizes simple properties of principal components analysis to identify outliers in the transformed space.

The algorithm consists of two basic parts: a first phase that aims to detect location outliers, and a second phase that aims to detect scatter outliers.

Scatter outliers possess a different scatter matrix than the rest of the data, while location outliers are described by a different location parameter.

To start, the authors propose to robustly rescale or sphere each component using the coordinatewise median and the median absolute deviation (MAD), according to

$$x_{ij}^* = \frac{x_{ij} - \text{med}(x_{1j}, \dots, x_{nj})}{\text{MAD}(x_{1j}, \dots, x_{nj})}, \quad j = 1, \dots, p, \quad (1.20)$$

and where the median absolute deviation (MAD) defined for a sample $\{x_1, \dots, x_n\} \in \mathbb{R}$ is computed like

$$\text{MAD}(x_{1j}, \dots, x_{nj}) = 1.4826 \text{med}_j |x_j - \text{med}_i x_i|. \quad (1.21)$$

Other kind of estimators of location and scale are given by the class of S-estimators, Starting with the rescaled data x_{ij}^* , a weighted covariance matrix is calculated, and from this matrix, the eigenvalues and the eigenvectors are computed and hence a semi-robust principal component decomposition. Only those eigenvectors/values that contribute to at least 99% of the total variance; call this new dimension p^* . For the case $p \gg n$, this also solves the singularity problem since $p^* < n$. For the $p^* \times p^*$ matrix of eigenvectors V , the matrix of principal components is computed as

$$Z = X^*V, \quad (1.22)$$

where X^* is the matrix with the elements x_{ij}^* . The principal components are rescaled by the median and the MAD similar to (1.20),

$$z_{ij}^* = \frac{z_{ij} - \text{med}(z_{1j}, \dots, z_{nj})}{\text{MAD}(z_{1j}, \dots, z_{nj})}, \quad j = 1, \dots, p^*. \quad (1.23)$$

After the above pre-processing steps, the location outlier phase is initiated by calculating the absolute value of a robust kurtosis measure for each component according to:

$$w_j = \left| \frac{1}{n} \sum_{i=1}^n \frac{(z_{ij}^* - \text{med}(z_{1j}^*, \dots, z_{nj}^*))^4}{\text{MAD}(z_{1j}^*, \dots, z_{nj}^*)^4} - 3 \right|, \quad j = 1, \dots, p^*. \quad (1.24)$$

This quantity allows to assign weights to each component. The authors use relative weights $w_i / \sum_j w_j$. If no outliers are present in a given component, one expect the principal components to be approximately normally distributed (similar to the original data), yielding a kurtosis close to zero. Also, a robust Mahalanobis distance (denoted RD_i) is calculated using the

distance from the median (as scaled by the MAD), weighting each component according to the relative weights $w_i/\sum_j w_j$, with the kurtosis measure w_i defined in (1.24).

To finish the first phase of the algorithm, these robust Mahalanobis distances $\{RD_i\}$ are transformed according to

$$d_i = RD_i \frac{\sqrt{\chi_{p^*,0.5}^2}}{\text{med}(RD_1, \dots, RD_n)}, \quad i = 1, \dots, n, \quad (1.25)$$

where $\chi_{p^*,0.5}^2$ is the $\chi_{p^*}^2$ 50th quantile. The translated biweight function is used to assign weights to each observation and use these weights as a measure of outlyingness. The weights for each observation are calculated according to

$$w_{1i} = \begin{cases} 0, & d_i \geq c, \\ \left(1 - \left(\frac{d_i - M}{c - M}\right)^2\right)^2, & M < d_i < c, \\ 1, & d_i \leq M, \end{cases} \quad (1.26)$$

where $i = 1, \dots, n$, M is the $33\frac{1}{3}$ rd quantile of the distances $\{d_1, \dots, d_n\}$, and

$$c = \text{med}(d_1, \dots, d_n) + 2.5 \text{MAD}(d_1, \dots, d_n). \quad (1.27)$$

The second phase of the algorithm is similar to the first except that the kurtosis weighting scheme is not used. Principal components focuses on those directions that have large variance, so the algorithm search for scatter outliers in the semi-robust principal component space described before. That is, the algorithm search for the outliers in the space defined by Z^* from (1.23). As before, calculating the Euclidean norm for data in principal component space is equivalent to the Mahalanobis distance in the original data space.

Similarly to the first phase, weights for each robust distance are calculated according to (1.26) and setting M^2 equal to the $\chi_{p^*}^2$ 25th quantile and c^2 equal to the $\chi_{p^*}^2$ 99th quantile. Call the weights calculated this way, w_{2i} , $i = 1, \dots, n$.

Finally, the algorithm combine weights from these two steps to calculate final weights w_i , $i = 1, \dots, n$, according to

$$w_i = \frac{(w_{1i} + s)(w_{2i} + s)}{(1 + s)^2}, \quad (1.28)$$

where typically the scaling constant $s = 0.25$. Outliers are then classified as points that have weight $w_i < 0.25$.

These phases can be summarized in the following scheme:

Phase 1 Detection of location outliers.

- Step 1 Robustly sphere the data according to (1.20) using the coordinate-wise median and the median absolute deviation (MAD).
Calculate the sample covariance matrix of the transformed data X^* .
- Step 2 Compute a principal component decomposition of the semi-robust covariance matrix from Step 1, and retain only those p^* eigenvectors whose eigenvalues contribute to at least 99% of the total variance. Robustly sphere the transformed data as in (1.23).
- Step 3 Compute the robust kurtosis weights for each component as in (1.24), and hence weighted norms for the sphere data from Step 2. Since the data have been scaled by the MAD, these Euclidean norms in principal component space are equivalent to robust Mahalanobis distances. Transform these distances according to (1.25).
- Step 4 Determine weights $w - 1i$ for each robust distance according the translate biweight in (1.26), with M equal to the $33\frac{1}{3}$ rd quantile of the distances $\{d_1, \dots, d_n\}$ and $c = med(d_1, \dots, d_n) + 2.5MAD(d_1, \dots, d_n)$.

Phase 2 Detection of scatter outliers.

- Step 5 Use the same semi-robust principal component decomposition calculated in Step 2 and compute the (unweighted) Euclidean norms of the data in principal component space. Transform according to (1.25) to yield a set of distances for use in Step 6.
- Step 6 Determine weights w_{2i} for each robust distance according to the translated biweight in (1.26) with c^2 equal to the $X_{p^*}^2$ 99th quantile and M^2 equal to the $X_{p^*}^2$ 25th quantile.

Combining Phase 1 and Phase 2: Use the weights from Step 4 and 6 to determine final weights for all observations according to (1.28).

Model-Based Clustering, Classification and Density Estimation

([Fra+12],[FR02]). The Cluster Analysis refers to the partitioning of a groups of individuals into groups following some criteria involving the variables measured on them.

In our case, we aim to classify a population of particles using the measures of a set physical and morphological variables. It would improve the degree of interpretation in the physical phenomenon and to better describe

the population. Note that we do not know, a priori, the number of groups existing in the population.

We use a Model-Based Clustering in order to classify a population of particles. We can see each component probability distribution in a finite mixture model as a cluster. Thus, problems that are usually associated with clustering can be addressed as a statistical model choice problem. For example, The problem of determining the number of clusters can be see as the problem of compare models that differ in number of component distributions. Outliers can be modeled also as a special component distribution (or distributions) representing the atypical data.

We use the library Mclust ([Fra+12]) in the R Language ([R C17]) in order to perform the cluster analysis using the mixture model-based with the presence of outliers.

The strategy arose from two methods based on multivariate normal mixture models with covariances parameterized by eigenvalue decomposition. These methods are hierarchical agglomeration based on the classification likelihood and the EM algorithm for maximum likelihood estimation of multivariate mixture models. The two approaches are complementary; model-based hierarchical agglomeration tends to produce reasonably good partitions even when started without any information about the groupings, whereas initialization is critical in expectation-maximization (EM) because the likelihood surface tends to have multiple modes, although EM typically produces improved partitions when started from reasonable ones.

By initializing EM with partitions from model-based hierarchical agglomeration and using approximate Bayes factors with the Bayesian Information Criterion (BIC) approximation to determine the number of groups present in the data. The method proposed allows to select the parameterization of the model as well as the number of clusters simultaneously using BIC.

Mixture Models Lets denote \mathbf{y} a data matrix formed by independent multivariate observations of observations $\mathbf{y}_1, \dots, \mathbf{y}_n$, the likelihood for a mixture model with C components is

$$\mathcal{L}(\theta_1, \dots, \theta_C; \tau_1, \dots, \tau_C | \mathbf{y}) = \prod_{i=1}^n \sum_{k=1}^C \tau_k f_k(\mathbf{y}_i | \theta_k) \quad (1.29)$$

where f_k and θ_k are the density and parameters of the k th component in the mixture and τ_k is the probability that an observation belongs to the k th component with $\tau_k \geq 0$ and $\sum_{k=1}^C \tau_k = 1$. Most commonly, f_k is the multivariate normal (Gaussian) density ϕ_k parameterized by its mean μ_k

and covariance matrix Σ_k ,

$$\phi_k(\mathbf{y}_i | \mu_k, \Sigma_k) = \frac{\exp\{-\frac{1}{2}(\mathbf{y}_i - \mu_k)^T \Sigma_k^{-1}(\mathbf{y}_i - \mu_k)\}}{\sqrt{\det(2\pi\Sigma_k)}}. \quad (1.30)$$

Data generated by mixtures of multivariate normal densities are characterized by groups or clusters centered at the means μ_k . Geometric features (shape, volume, orientation) of the clusters are determined by the covariances Σ_k .

The covariance matrix Σ_k can be used to characterize several possible mixture models for the clusters. We mention the mixture model implemented in the Mclust library:

- $\Sigma_k = \lambda I$: all clusters are spherical of the same size,
- $\Sigma_k = \Sigma$ constant across clusters: all the clusters have the same geometry which is not necessarily spherical,
- $\Sigma_k = \lambda_k I$: clusters are spherical but have different volume,
- $\Sigma_k = \lambda_k A_k$: all the covariances are diagonal but their size, shapes and orientation are allowed to vary,
- Σ_k unrestricted: each cluster may have a different geometry.

Cluster Analysis. The strategy presented by the authors is based on mixture models. A parameterization for the covariance matrices Σ_k through eigenvalues decomposition is used as basis for a class of models that is sufficiently flexible to accommodate data with widely varying characteristics. The parameterization is

$$\Sigma_k = \lambda_k D_k A_k D_k^T \quad (1.31)$$

where D_k is the orthogonal matrix of eigenvectors, A_k is a diagonal matrix whose elements are proportional to the eigenvalues, and λ_k is an associated constant of proportionality.

The strategy comprises three core elements: initialization via model-based hierarchical agglomerative clustering, maximum likelihood estimation via EM algorithm, and selection of the model and the number of clusters using approximate Bayes factors with the BIC approximation.

Model-based hierarchical agglomerative clustering is an approach to computing an approximate maximum for the *classification likelihood*,

$$\mathcal{L}_{CL}(\theta_1, \dots, \theta_G; \ell_1, \dots, \ell_n | \mathbf{y}) = \prod_{i=1}^n f_{\ell_i}(\mathbf{y}_i | \theta_{\ell_i}) \quad (1.32)$$

where the ℓ_i are labels indicating a unique classification of each observation, $\ell_i = k$ if \mathbf{y}_i belongs to the k th component. In the mixture likelihood (1.29), each component is weighted by the probability that an observation belongs to that component. The presence of the class labels in the classification likelihood (1.32) introduces a combinatorial aspect that makes exact maximization impractical.

Model-based hierarchical agglomerative clustering proceeds by successively merging pairs of clusters corresponding to the greatest increase in the classification likelihood (1.32) among all possible pairs. In the absence of any information about groupings, the procedure starts by treating each observation as a singleton cluster.

The *EM algorithm* (Dempster, Laird, and Rubin 1977) is a general approach to maximum likelihood estimation for problems in which the data can be viewed as consisting of n multivariate observations \mathbf{x}_i recoverable from $(\mathbf{y}_i, \mathbf{z}_i)$, in which \mathbf{y}_i is observed and \mathbf{z}_i is unobserved. If the \mathbf{x}_i are independent and identically distributed (iid) according to a probability distribution f with parameter θ , then the *complete-data likelihood* is

$$\mathcal{L}_C(\mathbf{x}_i | \theta) = \prod_{i=1}^n f(\mathbf{x}_i | \theta). \quad (1.33)$$

Further, if the probability that a particular variable is unobserved depends only on the observed data \mathbf{y} and not on \mathbf{z} , then the observed-data likelihood, $\mathcal{L}_O(\mathbf{y} | \theta)$, can be obtained by integrating \mathbf{z} out of the complete-data likelihood,

$$\mathcal{L}_O(\mathbf{y} | \theta) = \int \mathcal{L}_C(\mathbf{x} | \theta) d\mathbf{z}. \quad (1.34)$$

The EM algorithm alternates between two steps, an "E step", in which the conditional expectation of the complete-data log-likelihood given the observed data and the current parameter estimates is computed, and a "M step", in which parameters that maximize the expected log-likelihood from the "E step" are determined. Under regularity conditions, EM can be shown to converge to a local maximum of the observed-data likelihood. The unobserved portion of the data may involve values that are missing due to nonresponse and/or quantities that are introduced to reformulate the problem for EM.

In EM for mixture models, the "complete data" are considered to be $\mathbf{x}_i = (\mathbf{y}_i, \mathbf{z}_i)$, where $\mathbf{z}_i = (z_{i1}, \dots, z_{iG})$ is the unobserved portion of the data, with

$$z_{ik} = \begin{cases} 1, & \text{if } \mathbf{x}_i \text{ belongs to group } k \\ 0, & \text{otherwise.} \end{cases} \quad (1.35)$$

Assuming that each \mathbf{z}_i is iid according to a multinomial distribution of one draw from G categories with probabilities τ_1, \dots, τ_G , and that the density of an observation \mathbf{y}_i given \mathbf{z}_i is given by $\prod_{k=1}^G f_k(\mathbf{y}_i | \theta_k)^{z_{ik}}$, the resulting complete-data log-likelihood is

$$l(\theta_k, \tau_k, z_{ik} | \mathbf{x}) = \sum_{i=1}^n \sum_{k=1}^G z_{ik} \log [\tau_k f_k(\mathbf{y}_i | \theta_k)]. \quad (1.36)$$

The E step of the EM algorithm for mixture models is given by

$$\hat{z}_{ik} \leftarrow \frac{\hat{\tau}_k f_k(\mathbf{y}_i | \hat{\theta}_k)}{\sum_{j=1}^G \hat{\tau}_j f_j(\mathbf{y}_i | \hat{\theta}_j)}, \quad (1.37)$$

while the M step involves maximizing (1.36) in terms of τ_k and θ_k with z_{ik} fixed at the values computed in the E step, \hat{z}_{ik} . The value z_{ik}^* of \hat{z}_{ik} at a maximum of (1.29) is the estimated conditional probability that observation i belongs to group k .

For multivariate normal mixtures, the E step is given by (1.37) with f_k replaced by ϕ as defined in (1.30), regardless of the parameterization. For the M step, estimates of the means and probabilities have simple closed-form expressions involving the data and \hat{z}_{ik} from the E step,

$$\hat{\tau}_k \leftarrow \frac{n_k}{n}; \quad \hat{\mu}_k \leftarrow \frac{\sum_{i=1}^n \hat{z}_{ik} \mathbf{y}_i}{n_k}; \quad n_k \equiv \sum_{i=1}^n \hat{z}_{ik}. \quad (1.38)$$

Computation of the covariance estimate $\hat{\Sigma}_k$ depends on its parameterization.

The approach proposed by the authors to the problem of model selection in clustering is based on Bayesian model selection via Bayes factors and posterior model probabilities. The basic idea is that if several models M_1, \dots, M_K , are considered, with posterior probabilities $p(M_k)$, $k = 1, \dots, K$ (often taken equal), then, by Bayes's theorem, the posterior probability of model M_k given data D is proportional to the probability of the data given model M_k , times the model's prior probability, namely

$$p(M_k | D) \propto p(D | M_k) p(M_k). \quad (1.39)$$

When there are unknown parameters, by the law of total probability, $p(D | M_k)$ is obtained by integrating over parameters, that is

$$p(D | M_k) = \int p(D | \theta_k, M_k) p(\theta_k | M_k) d\theta_k, \quad (1.40)$$

where $p(\theta_k | M_k)$ is the prior distribution of θ_k , the parameter vector of model M_k . The quantity $p(D | M_k)$ is known as the integrated likelihood of model k . In hierarchical agglomeration, each stage of merging corresponds to a unique number of clusters and an unique partition of the data. A given partition can be transformed into indicator variables (1.35), which can then be used as conditional probabilities in a M step of EM for parameter estimation, initializing an EM algorithm. This, combined with Bayes factors as approximated by BIC for model selection, yields a comprehensive clustering strategy:

- Determine a maximum number of clusters, M , and a set of mixture models to consider.
- Perform hierarchical agglomeration to approximately maximize the classification likelihood for each model, and obtain the corresponding classification for up to M groups.
- Apply the EM algorithms for each model and each number of clusters $2, \dots, M$, starting with the classification from hierarchical agglomeration. Compute BIC for the one-cluster case for each model and for the mixture model with the optimal parameters from EM for $2, \dots, M$ clusters.

Strong evidence for a model and an associated number of clusters is taken to correspond to a decisive maximum of the BIC.

Linear Discriminant Analysis [Ren02] Discriminant analysis is used in situations where the individuals are classified a priori in groups or clusters. There are two major objective in separation of groups:

- *Description of groups separation.* in which linear functions of the variables (discriminant functions) are used to describe or elucidate the differences between two or more groups. Descriptive discriminant analysis aims to identify the relative contributions of the explicative variables to separation of the groups. This tool helps to find a subset of the original variables that separates the groups almost as well as the original set. Also the interest is to find an optimal plane on which the points can be projected in order to best illustrate the configuration of the groups.
- *Prediction of allocation of observations to groups.* The classification functions are employed to assign an individual unit to one of the groups finding the group to which the individual most likely belongs.

For k groups, with n_i observations in the i th group, the observation vector y_{ij} are transformed to obtain $z_{ij} = a'y_{ij}$, for $i = 1, 2, \dots, k$; $j = 1, 2, \dots, n_i$, and find the means $\bar{z}_i = a'\bar{y}_i$, where $\bar{y}_i = \sum_{j=1}^{n_i} y_{ij}/n_i$.

We seek a vector a that maximally separates $\bar{z}_1, \bar{z}_2, \dots, \bar{z}_k$. To express separation among $\bar{z}_1, \bar{z}_2, \dots, \bar{z}_k$, we use the H and E matrix and we write

$$\lambda = \frac{a'Ha}{a'Ea},$$

which can also be expressed as

$$\lambda = \frac{SSH(z)}{SSE(z)},$$

where $SSH(z)$ and $SSE(z)$ are the between and within sums of squares for z , and where

$$H = n \sum_{i=1}^k (\bar{y}_i - \hat{y}_{..}) (\bar{y}_i - \hat{y}_{..})',$$

$$E = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i) (y_{ij} - \bar{y}_i)'$$

and

$$y_{i.} = \sum_{j=1}^{n_i} y_{ij}$$

$$y_{..} = \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}.$$

We can write this expression in the form

$$a'Ha = \lambda a'Ea \Leftrightarrow a'(Ha - \lambda Ea) = 0.$$

We search for the value of a that results in maximum λ (the solution $a' = 0'$ is not permissible because it gives $\lambda = \frac{0}{0}$). Other solutions are found from

$$Ha - \lambda Ea = 0 \Leftrightarrow (E^{-1}H - \lambda I) a = 0. \quad (1.41)$$

The solutions of (1.41) are the eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_s$ and associated to the eigenvectors a_1, a_2, \dots, a_s of $E^{-1}H$, and we consider them to be ranked $\lambda_1 > \lambda_2 > \dots > \lambda_s$. The number of (nonzero) eigenvalues s is the rank of H , which can be found as the smaller of $k - 1$ and p . Thus, the largest eigenvalue λ_1 is the maximum value of $\lambda = a'Ha/a'Ea$, and the coefficient vector

that produces the maximum is the corresponding eigenvector a_1 . Hence, the discriminant function that maximally separates the means is $z_1 = a'_1 y$, z_1 represents the dimension or direction that maximally separates the means.

Then, we obtain s discriminant functions $z_i = a'_i y$, for $i = 1, 2, \dots, s$, which show the dimensions or directions of differences among $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_k$. This discriminant functions are uncorrelated, but they are not orthogonal because $E^{-1}H$ is nor symmetric.

The relative importance of each discriminant function z_i is given by

$$\frac{\lambda_i}{\sum_{j=1}^s \lambda_j}. \quad (1.42)$$

Standardized discriminant functions. The contributions of the explicative variables y 's to separation of several groups can be examined from the coefficients of the standardized discriminant functions. If we denote the r th coefficient in the m th discriminant function by a_{mr} , $m = 1, 2, \dots, s$; $r = 1, 2, \dots, p$, then the standardized form is

$$a_{mr}^* = s_r a_{mr}, \quad (1.43)$$

where s_r is the within-group standard deviation obtained from the diagonal of $S_{pl} = E/v_E$.

Analysis PCA The initial distribution of flocs was characterized by $Rg \in [10, 200]$ and $C_I \in [0.3, 0.95]$. Variety of size and shape are not of the same nature.

It is possible to obtain empirical data from the flocculation phenomenon using techniques like Granulometry for laser diffraction in Jar test ([Vli14]). This processes usually involves the treatment of floc's images, measuring a set of internals characteristics that can be classified into:

- Size describing measures,
- Shape or morphology describing measures,
- Localization measures.

Among this measures, we have:

- **Size describing measures**
 - Volume distribution of Circle Equivalent Diameter (CED): diameter of the sphere having the same image of diffraction.

- Mean radius of gyration (Rg): it is the mean square of distances from the components of the floc to the center.
- Area (A), Perimeter (P): the area (perimeter) of the disk having the same area of the object's image.
- The size can also be characterized for the biggest dimension (the maximum distance between two pixels belonging to the object's image).

- **Shape or morphology describing measures**

- Fractal dimension: the notion of fractal, based on the invariance for scale change (repetition of a motif) (Df): the exponent relating the mass of the floc with the diameter apparent.

Mass Df : $N = Kg \left(\frac{Rg}{r_0} \right)^{Df}$, where N : number of primary particles contained (similar to the mass of the aggregate adimensionalized for the mass of one primary particle), Rg : radius of gyration, r_0 : radius of the primary particles, Kg : form factor. $1 \leq Df \leq 3$.

The flocculation process can result in a great quantity of aggregates and they can be very different, with a large variability. Usually the granulometry software takes several hours to analyse one image. Moreover, the description of flocculation processes involve the time evolution of the flocs present in the system, then several images are taken (and analyzed) over experimental time ([Mor]).

Depending on Shape and Size describing measures, the flocs have different kind of chemical properties ([Mor]). Because of this, it is important to understand the behavior of these measures as individually as in the multivariate context.

The main objective of measuring all these characteristics is to know the number distribution of flocs in function of its internal variables. Nevertheless, the minimal set of these variables is required in order to facilitate the interpretation. The classical analysis implies the study of the number distribution in function of one measure, usually a size describing one although recent studies pursued describing number distribution in function of some size and shape measures ([WMpt], [Vli14]).

1.4.6 Bentonite Experimental Data

We also found that the properties of size are highly positively correlated, and having a very similar distribution each other. For the shape measures, a similar behavior was found, being the less correlated convexity and circularity.

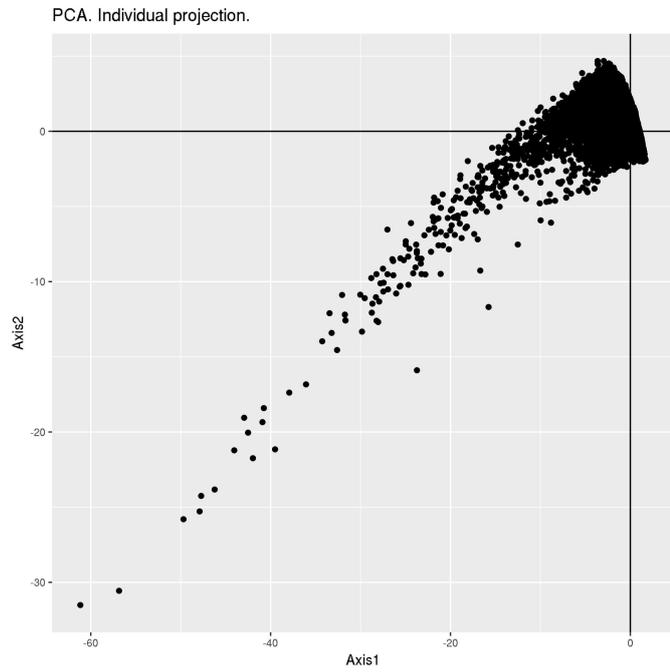
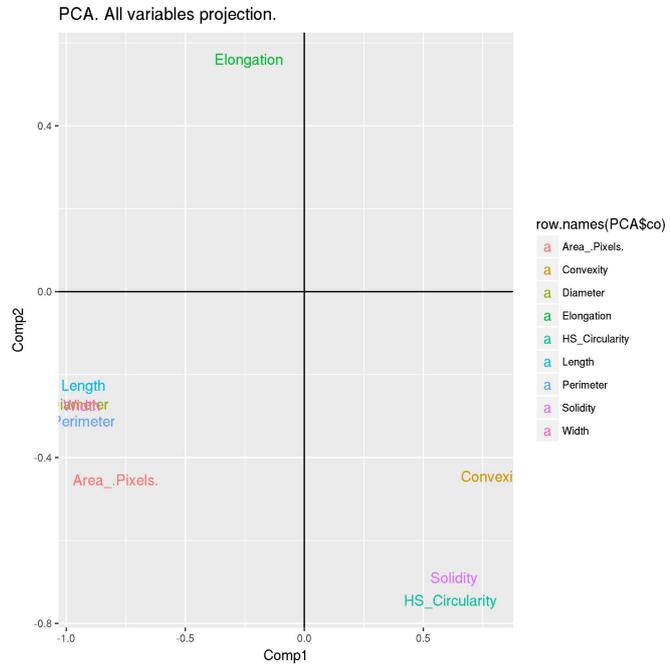


Figure 1.18: CPA variables

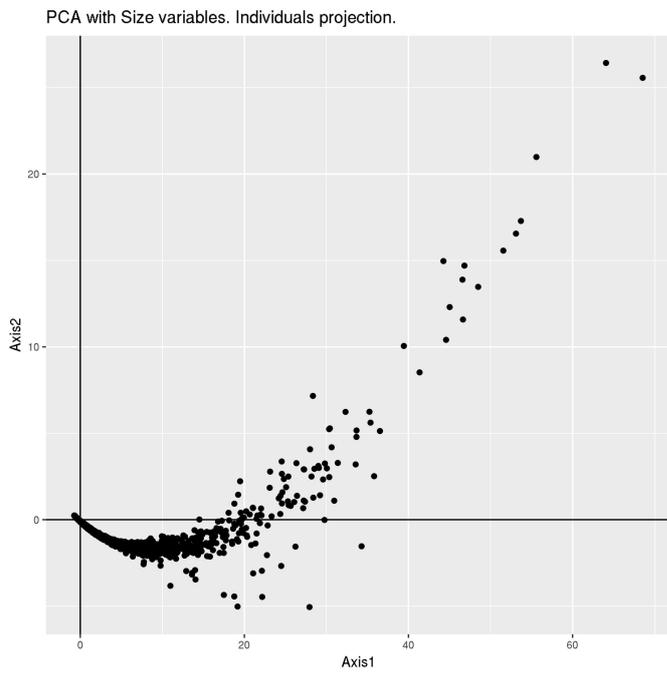
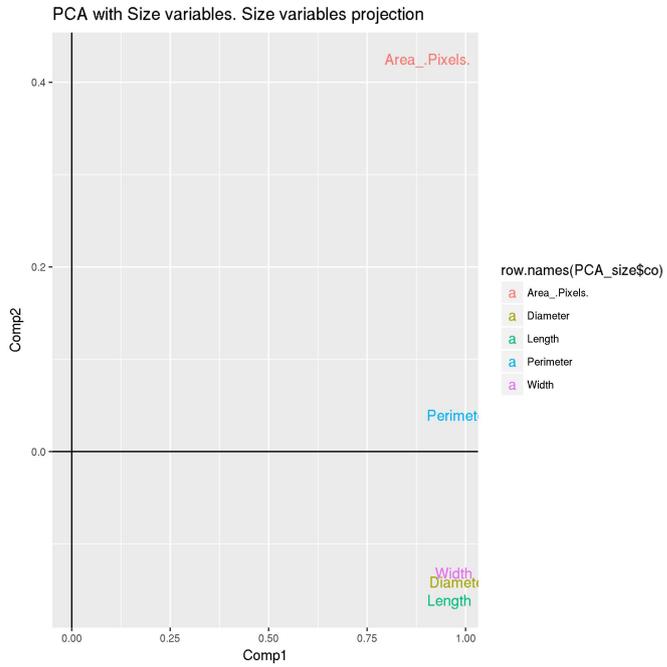


Figure 1.19: CPA variables

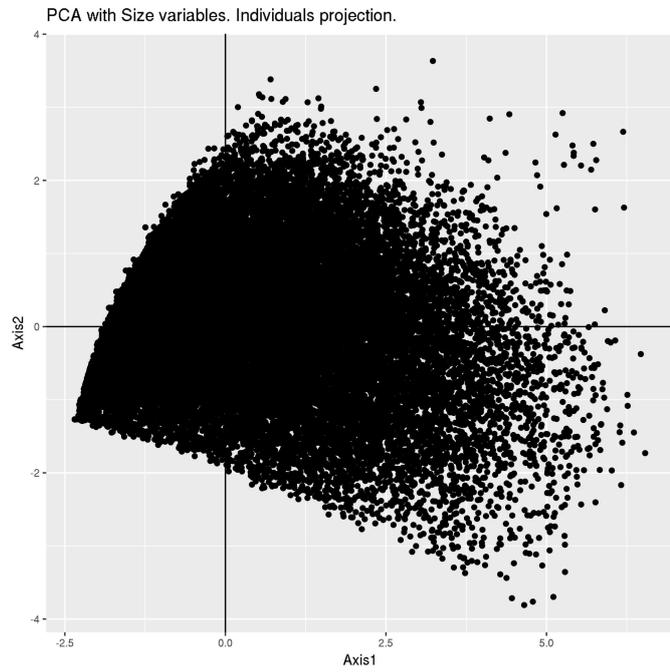
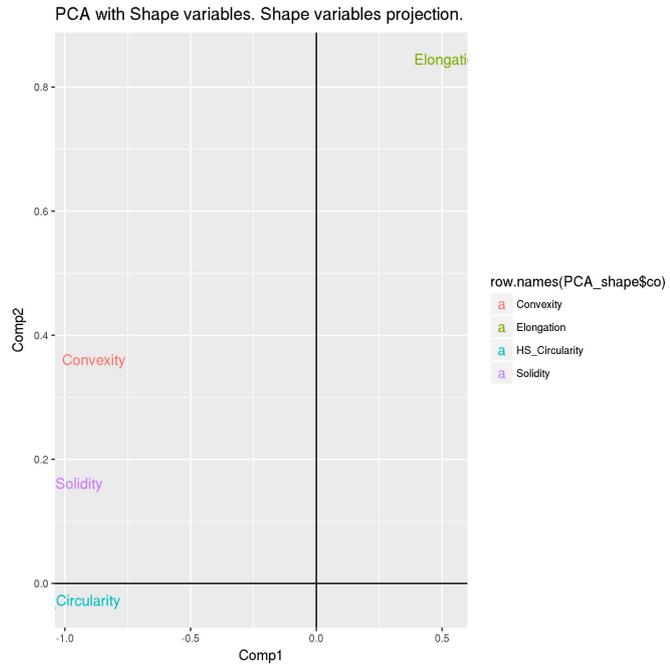
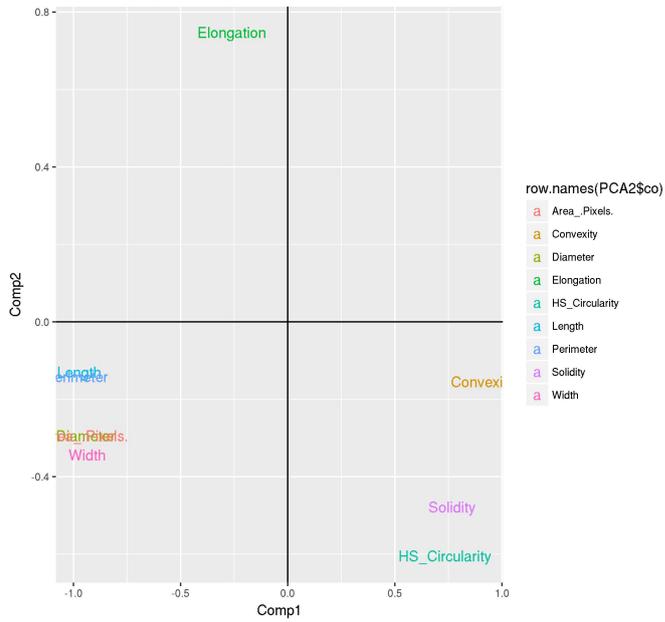


Figure 1.20: CPA variables

PCA with all variables and Size variables in logarithmic transformation.
Variables projection.



PCA with all variables and Size variables in logarithmic transformation.
Individuals projection.

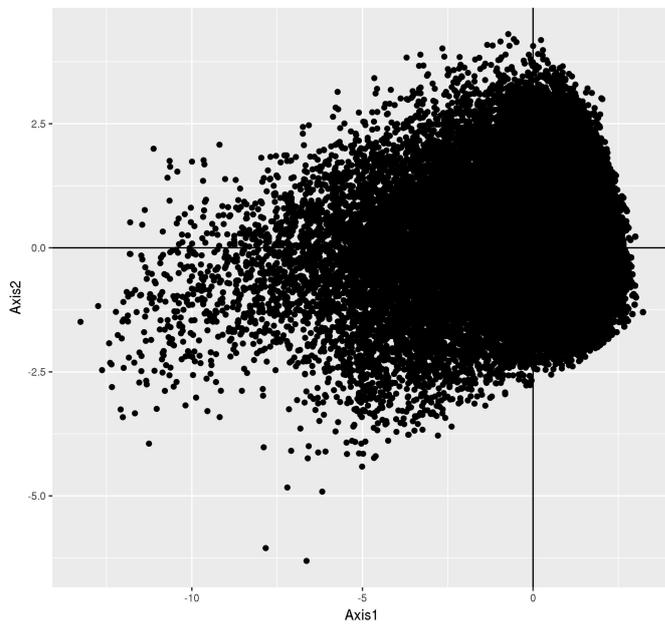
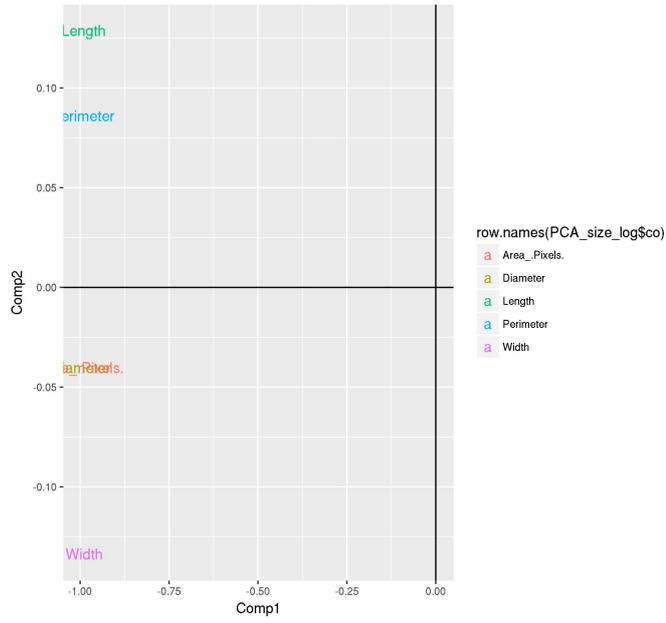


Figure 1.21: CPA variables

PCA with Size variables in logarithmic transformation.
Variables projection.



PCA with Size variables in logarithmic transformation.
Individuals projection.

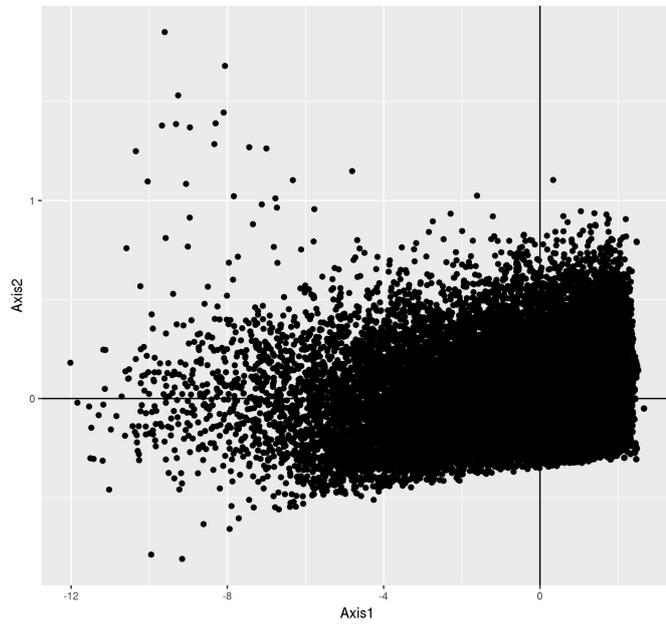


Figure 1.22: CPA variables

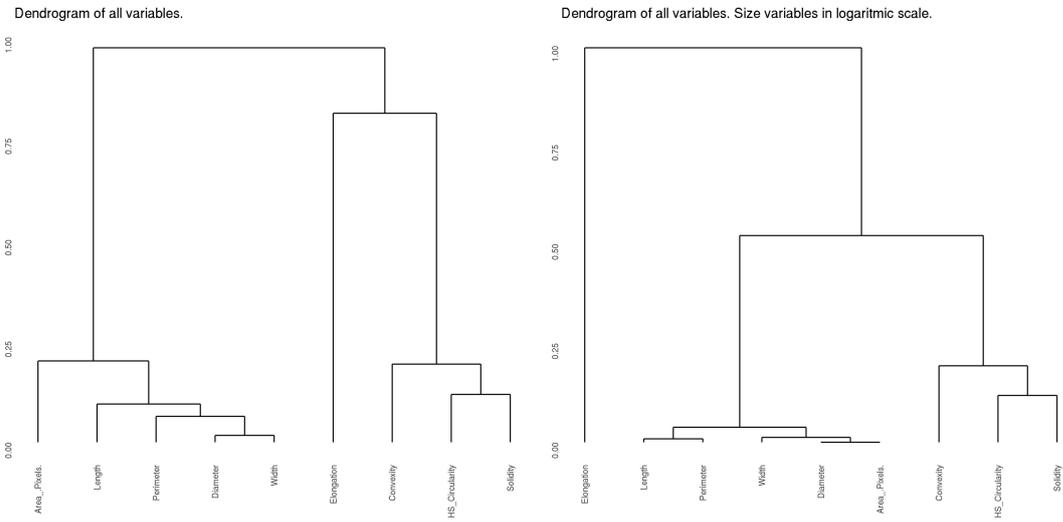


Figure 1.23: Dendrogram for all the variables

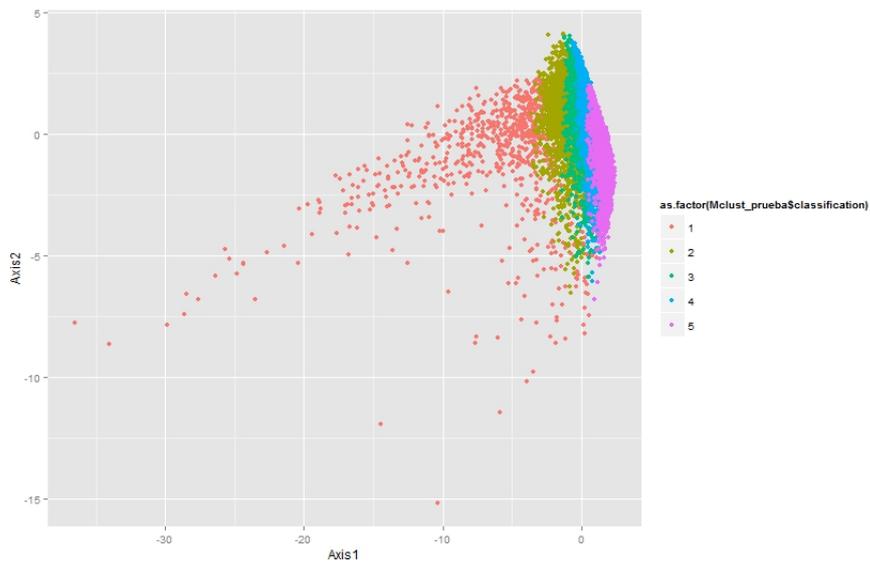


Figure 1.24: PCA. Individuals

In order to have a better representation of the results, we made the same comparisons using a logarithmic transformation only for the Size variables, because of the large range of variation in Size variables. We repeat the CPA analysis and we find the same trend in the relation between the variables. In General, the first Principal Component represent the Size variables, with the Perimeter and the Diameter as the more important ones. The second Principal Component can be interpreted as related to Shape variables with the Circularity and the Elongation as the more important ones.

Further, an PCA was performed only taking account one group of variables at time. Thus, for the Size variables using the logarithmic transformation, we have that the first component explains almost 97% of the total variability and, the more important variable was the Diameter. In the other hand, for the Shape variables, the most important variable was the Circularity. Then, we can summarize the behavior of the Size variables using the Diameter and the behavior of the Shape variables using the Circularity.

In the PCA graphic for the individuals, we can remark a possible division of the original population of flocs into sub-populations, having different behavior for Size and Shape variables. Also, it is possible to find outliers individuals. We use an outlier multivariate detection technique in order to mark some individuals that have a very different behavior from the others. For that, we use the algorithm of Filzmoser, Marona and Werner. Then we use an hierarchical cluster analysis in order to characterize the sub-populations. We use a Normal Mixture Modeling cluster technique, that compare a set of models in order to find the better fit with the population and estimates the number of groups present in the population. The procedure exclude the individuals marked as outliers.

We find using the Bayesian Information Criteria the presence of five groups or sub-populations. We can describe those sub-populations using the Diameter for summarizing the Size variables behavior and the Circularity for summarizing the Shape variables.

Morphological analysis of Bentonite's flocs for different hydrodynamic conditions

Experimental empirical data was obtained in ([Vli14]). They analyzed the behavior of one micron (μm) bentonite primary particles in Jar test. They used several mixing speed in order to get different conditions for aggregation and breakage processes. The speed of mixing used were 30, 50, 70 and 90 revolutions per minute (rpm).

Bentonite in a kind of clay of type smectite, montmorillonite. The clay particles are the plates consisting of a stack of sheets separated by an inter-

lacing space.

The experiments originating the empirical data implement the flocculation of Bentonite provided by CECA Chemicals. The mass concentration of the suspension is 30 mgL^{-1} . The suspension is diluted as much as allowing to take images.

The bentonite's mass needed for the experiments was putted in the water at least 24 hours before the experiments. This suspension had high intensity agitation by 45 *min*, before introduction into the reactor. Therefore, the suspension is constituted of primary agglomerates, compound of several elementary particles.

The initial population of this primary agglomerates shows that they have high variation about the Circle Equivalent Diameter (CED), with some values larger than $200 \mu\text{m}$. The most part of particles had size about some tens of micrometers and there are a small proportion of agglomerates with size larger than $200 \mu\text{m}$. The water used was demineralized water, the Ph is about 4.5 ± 0.1 and temperature between 20 and 25 °C. The experiments were performed in a Taylor-Couette Reactor.

In the experiments, they used a mechanism of aggregation by charges neutralization, using aluminum sulfate $Al_2(SO_4)_3$ as coagulant, which is used commonly in water treatment. It allows to produce flocs weak enough to interact with the hydrodynamic. The concentration of aluminum sulphate in the reactor is $3.5 \times 10^{-5} \text{ molL}^{-1}$.

1.4.7 Conclusions

A continuation, the most important conclusions obtained from the descriptive data analysis are mentioned. The conclusions are divided in:

- Conclusions from the descriptive analysis of the data.
- Conclusions from the exploratory multivariate analysis of the data.

Conclusions of the descriptive analysis of the data.

The analysis was done using a natural classification for the variables in "size" variables and "shape" variables. Then, the description searched to interpret the data in that terms.

Initial population of Bentonite flocks.

Populations of Bentonite flocks under different hydrodynamics conditions.

Populations of Bentonite flocks under different hydrodynamics conditions thought time.

Chapter 2

Numerical methods for recovering the moments of the number density function

2.1 Introduction

Population Balances for simultaneous coagulation and breakage are employed in describing many systems including aerosols, powders and polymers, and many unit operations including reactors, crystallizers, and size reduction (or enlargement) equipment. Solving this kind of equations is very important for Computational Fluid Dynamic (CFD) simulations or with process flow-sheet simulations ([WVF05],[Soo+07b]). The classical methods direct calculating the particle-size distribution evolution consume large computational resources. Often, this techniques include scale size discretization and solving the equation on each interval or the use of Monte Carlo methods, that makes the implementation in such simulations nonviable.

Conventional moments models are computationally less demanding, but are restricted to those systems for which the set of moments equations are closed. Furthermore, unless the moments themselves are the targets of the model, the method also poses the inversion problem. Recent approaches allows to relate the CFD simulation with techniques involving the moments of the distribution. One of those approaches is the Quadrature Method of Moments (QMOM), first proposed by McGraw (1997)[McG97], where the closure problem is solved by using a quadrature approximation involving some finite support set (abscissas) and weights.

2.2 Formulation of the model

The PBE is an equation in the foregoing number density and may be regarded as representing a *number balance* on particles of a particular state. The equation is often coupled with conservation equations for entities in the particles' environmental (or continuous) phase [Ram00].

The population balance equation basically accounts for various ways in which particles of a specific state can either form in or disappear from the system. When particle states are continuous, then processes, which cause their smooth variation in time, must contribute to the rates of formation and disappearance of specific particle types. There are several ways in which number of particles of a particular type can change is by processes that create new particles ("birth" processes) and destroy existing particles ("death" processes).

Birth of new particles can occur due to breakage or splitting processes, aggregation processes and so on. Breakage and aggregation processes also contribute to death processes, for a particle type that either breaks (into other particles) or aggregates with another particle no longer exists as such following the event [Mar+03].

We are interested in the analysis of this kind of populations in liquid-solid suspensions containing particles, where the evolution of the system is reflected in the time evolution of the number distribution of the particles about a size characteristic (longitude, area, volume, for example). In this kind of systems, the mass in the suspension is not modified. The aggregations and breakages models describe the probabilities of the particles having some properties are added or broken ([Van00],[Soo+07a]).

The **Balance Population Equation (PBE)** including aggregation and breakage is defined in terms the number distribution in function of only the volume of the particle as main property by [MVF03]

In the precedent chapter, we presented in section 1.2, proposition 3.2 the **Balance Population Equation (PBE)** (equation 3.2) including aggregation and breakage is defined in terms the number distribution in function of only the size of the particle as main property. We show that equation as

reference

$$\begin{aligned}
\frac{\partial n(L, t)}{\partial t} &= B^a(L; t) - D^a(L; t) + B^b(L; t) - D^b(L; t) \\
&= \frac{L^2}{2} \int_0^L \frac{\beta\left(\frac{(L^3 - \lambda^3)^{1/3}}{L}, \lambda\right)}{(L^3 - \lambda^3)^{2/3}} n\left(\frac{(L^3 - \lambda^3)^{1/3}}{L}, t\right) n(\lambda, t) d\lambda \\
&\quad - n(L, t) \int_0^\infty \beta(L, \lambda) n(\lambda, t) d\lambda \\
&\quad + \int_L^\infty a(\lambda) b(L | \lambda) n(L, t) d\lambda \\
&\quad - a(L) n(L, t),
\end{aligned}$$

Further, in section 1.3, we presented the PBE describing the time evolution of the standard moments of the number density. This equation was deduced from equation (3.2) as we have shown in proposition 1.3 (equation (1.11)). We show the equation here also as reference

$$\begin{aligned}
\frac{\partial m_k(t)}{\partial t} &= \frac{1}{2} \int_0^L n(\lambda; t) \int_0^\infty \beta(u, \lambda) n(u; t) (u^3 + \lambda^3)^{k/3} du d\lambda \\
&\quad - \int_0^\infty L^k n(L; t) \int_0^\infty \beta(L, \lambda) n(\lambda; t) d\lambda dL \\
&\quad + \int_0^\infty L^k \int_0^\infty a(\lambda) b(L/\lambda) n(\lambda; t) d\lambda dL \\
&\quad - \int_0^\infty a(\lambda) n(L; t) L^k dL
\end{aligned}$$

The methods that we are about to present in this chapter use the PBE describing the time evolution of the standard moments of the number density. The objective is to recover the time evolution of a finite set of standard moments of the number density from an initial set of standard moments values and also knowing the information from the aggregation and breakage kernels and parameters.

In order to apply some methods for recovering the time evolution of the standard moments of the number distribution, we are going to obtain a system of equations in differences from the Population Balance Equation in 3.2. This system of equation is expressed in function of a finite number of some abscissas points and weights representing the number distribution.

Starting from a initial estimation of those abscissas points and weights, then the system of equation in differences can be solve using some numerical differentiation scheme ([Mar+03],[Lag07]).

Proposition 2.1. *The Population Balance Equation can be used in order to obtain the following system of equation in differences*

$$\begin{aligned} \frac{\partial m_k}{\partial t} = & \frac{1}{2} \sum_i W_i \sum_j W_j (L_i^3 + L_j^3)^{k/3} \beta_{ij} - \sum_i L_i^k W_i \sum_j W_j \beta_{ij} \\ & + \sum_i a_i b_i^{-(k)} W_i - \sum_i a_i L_i^k W_i \end{aligned} \quad (2.1)$$

where L_i are a finite number of abscissas points and W_i a finite number of weights representing the number density distribution.

Proof. From proposition 1.3, we know that the PBE in terms of the standard moments of the number density function is given by equation (1.11)), this is

$$\begin{aligned} \frac{\partial m_k(t)}{\partial t} = & \frac{1}{2} \int_0^L n(\lambda; t) \int_0^\infty \beta(u, \lambda) n(u; t) (u^3 + \lambda^3)^{k/3} du d\lambda \\ & - \int_0^\infty L^k n(L; t) \int_0^\infty \beta(L, \lambda) n(\lambda; t) d\lambda dL \\ & + \int_0^\infty L^k \int_0^\infty a(\lambda) b(L/\lambda) n(\lambda; t) d\lambda dL \\ & - \int_0^\infty a(\lambda) n(L; t) L^k dL \end{aligned}$$

Now, we introduce the quadrature approximation

$$m_k = \int_0^\infty n(L) L^k dL \approx \sum_{i=1}^N W_i L_i^k \quad (2.2)$$

W_i , L_i and m_k depends on t

using this approximation into the equation (1.18) we get

$$\frac{\partial m_k(t)}{\partial t} = \overline{B}_k^a - \overline{D}_k^a + \overline{B}_k^b - \overline{D}_k^b \quad (2.3)$$

if we use L_i for λ and L_j for u , we have

$$\begin{aligned} \overline{B}_k^a(t) = & \frac{1}{2} \int_0^L n(\lambda; t) \int_0^\infty \beta(u, \lambda) n(u; t) (u^3 + \lambda^3)^{k/3} du d\lambda \\ \overline{B}_k^a(t) \approx & \frac{1}{2} \sum_i W_i \sum_j W_j (L_i^3 + L_j^3)^{k/3} \beta(L_j, L_i) \end{aligned} \quad (2.4)$$

$$\begin{aligned}\overline{D}_k^a(t) &= \int_0^\infty L^k n(L;t) \int_0^\infty \beta(L,\lambda) n(\lambda;t) d\lambda dL \\ \overline{D}_k^a(t) &\approx \sum_i L_i^k W_i \sum_j W_j \beta(L_i, L_j)\end{aligned}\tag{2.5}$$

$$\begin{aligned}\overline{B}_k^b(t) &= \int_0^\infty L^k \int_0^\infty a(\lambda) b(L/\lambda) n(\lambda;t) d\lambda dL \\ &= \int_0^\infty a(\lambda) \left[\int_0^\infty L^k b(L/\lambda) dL \right] n(\lambda;t) dL\end{aligned}\tag{2.6}$$

$$\begin{aligned}\overline{B}_k^b(t) &\approx \sum_i a(L_i) \left[\int_0^\infty L^k b(L/L_i) dL \right] W_i \\ \overline{D}_k^b(t) &= \int_0^\infty a(\lambda) n(L;t) L^k dL \\ \overline{D}_k^b(t) &\approx \sum_i a(L_i) L_i^k W_i\end{aligned}\tag{2.7}$$

Now, if we denote

- $\beta_{ij} = \beta(L_i, L_j)$
- $a_i = a(L_i)$
- $b_i^{-(k)} = \int_0^\infty L^k b(L/L_i) dL$

we get

$$\begin{aligned}\frac{\partial m_k}{\partial t} &\approx \frac{1}{2} \sum_i W_i \sum_j W_j (L_i^3 + L_j^3)^{k/3} \beta_{ij} - \sum_i L_i^k W_i \sum_j W_j \beta_{ij} \\ &\quad + \sum_i a_i b_i^{-(k)} W_i - \sum_i a_i L_i^k W_i\end{aligned}\tag{2.8}$$

□

2.3 Basis Pursuit

In the famous article [CDS98] it is treated the characterization of a signal s viewed as the decomposition

$$s = \sum_{\gamma \in \Gamma} \alpha_\gamma \phi_\gamma\tag{2.9}$$

or as the approximate decomposition

$$s = \sum_{i=1}^m \alpha_{\gamma_i} \phi_{\gamma_i} + R^{(m)} \quad (2.10)$$

where $R^{(m)}$ is a residual, and the ϕ_{γ} are waveforms in a dictionary $D = (\phi_{\gamma})_{\gamma \in \Gamma}$. Specifically:

- **Overcomplete representation.** In an overcomplete representation, the signal is represented as a superposition of waveforms, as $s = (s_t : 0 \leq t < n)$ where s is a discrete signal of length n
- **Dictionaries and atoms.** A dictionary is a collection of discrete time signals of length n as $D = (\phi_{\gamma})_{\gamma \in \Gamma}$ where the atoms ϕ_{γ} are discrete-time signals of length n parametrized by the vector γ .

The decomposition of the signal can be seen as

$$s = \Phi \alpha \quad (2.11)$$

where Φ is a $n \times p$ matrix containing p waveforms as columns $\alpha = (\alpha_r)$ is the vector of coefficients. When the dictionary furnishes a basis, Φ is a $n \times n$ non-singular matrix and we have unique representation $\alpha = \Phi^{-1}s$. When the atoms are, in addition, mutually orthogonal, then $\Phi^{-1} = \Phi^T$.

Given a dictionary of waveforms, one can distinguish *analysis* from *synthesis*:

- *Synthesis*; is the operation of building up a signal by superposing atoms. It involves a $n \times p$ matrix and $s = \Phi \alpha$.
- *Analysis*; involves the operation of associating with each signal a vector of coefficients attached to atoms, it involves a $p \times n$ matrix and $\tilde{\alpha} = \Phi^T s$.

In the overcomplete case, we are interested in $p \gg n$ and Φ is not invertible. The goals of adaptive representation are:

- **Sparsity:** To find the sparsiest possible representation, which means to find the representation of the signal with the fewest significant coefficients.
- **Superresolution:** We should obtain a resolution of sparse objects that is much higher resolution than that possible with traditional approach.
- **Speed:** A representation in order $O(n)$ or $O(n \log(n))$ time

Basis Pursuit: Basis Pursuit (BP) finds signal representations in overcomplete dictionaries by convex optimization; it obtains the decomposition that minimizes the ℓ^1 norm of the coefficients occurring in the representation.

BP can be used with noisy data by solving an optimization problem trading a quadratic misfit measure with an ℓ^1 norm of coefficients.

We assume that the dictionary is overcomplete, so that there are in general many representations $s = \sum_r \alpha_r \phi_r$. BP solves the problem:

$$\min \|\alpha\|_1 \quad \text{subject to } \Phi\alpha = s \quad (2.12)$$

BP is connected with linear programming, the problem treated in BP can be written as a linear program and solved using interior-point methods or primal-dual log-barrier methods.

Linear program. The Linear Program (LP) in so-called standard form is a constrained optimization problem defined in terms of a variable $x \in \mathbb{R}^m$ by

$$\min c^T x \quad \text{subject to } Ax = b, x \geq 0, \quad (2.13)$$

where $c^T x$ is the objective function, $Ax = b$ is a collection of equality constraints, and $x \geq 0$ is a set of bounds. The main question is which variables should be zero. The BP problem can be equivalently reformulated as a linear program in the standard form by making the following translations: $m \Leftrightarrow 2p$, $A \Leftrightarrow (\Phi, -\Phi)$, $b = s$, $c \Leftrightarrow (1, 1)$, $x \Leftrightarrow (u, v)$, and $\alpha = u - v$, with $(1, 1)$ a $1 \times 2p$ vector.

In the solution of the LP problem, suppose A is a $n \times m$ matrix with $n > m$ and suppose an optimal solution exists. It is known that a solution exists in which at most n of the entries in the optimal x are nonzero (in the generic case, the solution is called non-degenerated, and there are exactly n nonzeros). The nonzero coefficients are associated with n columns of A and these columns make up a basis of \mathbb{R}^n . Once the basis is identified, the solution is uniquely dictated by the basis. Thus, finding a solution to the LP problem is identical to find the optimal basis.

Then, we have from the LP results, the following decomposition

$$s = \sum_{i=1}^n \alpha_{\gamma_i}^* \phi_{\gamma_i} \quad (2.14)$$

the waveforms (ϕ_{γ_i}) are linearly independent but not necessarily orthogonal. The collection γ_i is not known in advance but depends on the problem data (the signal s). The selection of waveforms is therefore adaptive.

The algorithm BP-interior. The collection of feasible points $\{x : Ax = b, x \geq 0\}$ is a convex polyhedron in \mathbb{R}^m or a "simplex". The simplex

method works by walking around the boundary of this simplex, jumping from one vertex (extreme point) of the polyhedron to an adjacent vertex at which the objective function is "better".

Interior point methods instead starts from a point $x^{(0)}$ well inside the interior of the simplex ($x^{(0)} \gg 0$) and go "through the interior" of the simplex. Since the solution of a linear program is always at an extreme point of the simplex, as interior-point method converges, the current iterate $x^{(k)}$ approaches the boundary. Then, one may abandon the basic interior-point iteration and invoke a "crossover" procedure that uses simplex interactions to find the optimizing extreme point.

Translating this LP algorithm into BP terminology, one starts from a solution to the overcomplete representation problem $\Phi\alpha^{(0)} = s$, with $\alpha^{(0)} > 0$. One iteratively modifies the coefficients, maintaining feasibility $\Phi\alpha^{(k)} = s$ and apply a transformation that effectively sparsifies the vector $\alpha^{(k)}$. At some iteration, the vector has $\leq n$ significantly nonzero entries, and those correspond to the atoms appearing in the final solution. One forces all the other coefficients to zero and "jumps" to the decomposition in terms of the $\leq n$ selected atoms. More general interior-point algorithms starts with $\alpha^{(0)} > 0$ but do not require the feasibility $\Phi\alpha^{(k)} = s$ throughout, they achieve feasibility eventually.

Basis Pursuit denoising. Basis pursuit denoising refers to the solution of

$$\min \frac{1}{2} \| y - \Phi\alpha \|^2 + \lambda \| \alpha \|_1 \quad (2.15)$$

the solution $\alpha^{(\lambda)}$ is a function of the parameter λ . It yields a decomposition into signal-plus-residual

$$Y = s^{(\lambda)} + r^{(\lambda)} \quad (2.16)$$

where $s^{(\lambda)} = \Phi\alpha^{(\lambda)}$. The size of the residual is controlled by λ . As $\lambda \rightarrow 0$, the residual goes to zero and the solution behaves exactly like BP applied to y . As $\lambda \rightarrow \infty$, the residual gets large, we have $r^{(\lambda)} \rightarrow y$ and $s^{(\lambda)} \rightarrow 0$. (2.3) is equivalent to the following perturbed linear program:

$$\begin{aligned} \min_{x,p} \quad & c^T x + \frac{1}{2} \| p \|^2 \quad \text{subject to} \quad Ax + \delta p = b, \\ & x \geq 0, \\ & \delta = 1 \end{aligned} \quad (2.17)$$

where $A = (\Phi, -\Phi)$, $b = y$, $c = \lambda(1, 1)$, $x = (u, v)$, $\alpha = u - v$. Perturbed BP is really quadratic programming, but it retains a structure similar to BP. Hence, we can have a similar classification of algorithms into BPND-simplex and BPND-interior-point.

2.4 Exact reconstruction using Generalized Minimal Extrapolation

- They show that measures with finite support on the real line are the unique solution to an algorithm, named **Generalized Minimal Extrapolation (GME)**, involving only a finite number of generalized moments.
- GME shares related geometric properties with the **Basis Pursuit (BP)** [CDS98] approach.
- They also extend some standard results of compressed sensing (the dual polynomial, the nullspace property) to the **signed measure** framework
- They express exact reconstruction in terms of a **simple interpolation problem**.
- They prove that **every nonnegative measure**, supported by a set containing s points, can be exactly recovered from only $2s + 1$ generalized moments.
- The last result leads to a new construction of deterministic sensing matrices for compressed sensing.
- In this paper, they show that the exact reconstruction of a **signed measure** is still possible when one only knows the values of a finite number of a finite number of **non adaptive** linear measurements.
- Surprisingly, GME appears to uncover exact reconstruction results related to basis pursuit.
- **More precisely**, consider a *signed discrete measure* σ on a set $I := [-1, 1]$. Consider the *Jordan decomposition* $\sigma = \sigma^+ - \sigma^-$, and denote by S^+ (resp. S^-) the support of σ^+ (resp. σ^-). Let us define the *Jordan support* of the measure σ as the pair $J = (S^+, S^-)$. Assume further that $S := S^+ \cup S^-$ is **finite** and has cardinality s . Moreover, suppose that J belongs to a family Ψ of pairs of subsets of I . We call Ψ a *Jordan support family*. **The measure σ can be written as**

$$\sigma = \sum_{i=1}^s \sigma_i \delta_{x_i},$$

where $S = \{x_1, \dots, x_s\}$, $\sigma_1, \dots, \sigma_s$ are nonzero real numbers, and δ_x denotes the Dirac measure at point x .

- Let $F = \{u_0, u_1, \dots, u_n\}$ be **any** family of continuous functions on \bar{I} , where \bar{I} is the closure of I . Let μ be a **signed measure** on I . The k -th generalized moment of μ is defined by

$$c_k(\mu) = \int_I u_k d\mu$$

for all the indices $k = 0, 1, \dots, n$.

- **The main issue:** The reconstruction of the *target measure* σ from the observation of $K_n = (c_0(\sigma), \dots, c_n(\sigma))$. Assume that the support S and the weights σ_i of the target measure are **unknown**. We want to recover σ uniquely from K_n (*does an algorithm fitting $K_k(\sigma)$ among all the signed measures of I recover the measure σ ?*).

Generalized minimal extrapolation is the process of reconstructing a target measure σ from the observation $K_n(\sigma) = (c_0(\sigma), \dots, c_n(\sigma))$ of its first $n + 1$ generalized moments $c_k(\sigma)$ by finding a solution to the problem

$$\sigma^* \in \arg \min_{\mu \in \mathcal{M}} \|\mu\|_{TV} \quad \text{subject to} \quad K_n(\mu) = K_n(\sigma) \quad (2.18)$$

where the supremum is taken over all partition Π of I into a finite number of disjoint measurable sets.

By analogy with basis pursuit,

- on one hand, basis pursuit minimizes the ℓ_1 -norm subject to linear constraints,
- on the other hand, generalized minimal extrapolation naturally substitutes the TV -norm for the ℓ_1 -norm,
- for the case of Fourier coefficients, (GME) is simply Beurling Minimal Extrapolation.

Basis Pursuit: It is the process of reconstructing a target vector $x_0 \in \mathbb{R}^p$ from the observation $b = Ax_0$ by finding a sparse solution x^* to an underdetermined system of equations

$$x^* \in \arg \min_{y \in \mathbb{R}^p} \|y\|_1 \quad \text{subject to} \quad Ay = Ax_0 \quad (2.19)$$

where $A \in \mathbb{R}^{n \times p}$ is the design matrix.

GME looks for a minimizer among all the signed measures on I . Nevertheless, the target measure σ is assumed to be of extrema Jordan type.

By analogy with compressed sensing: if σ is of extrema Jordan type, then σ is a point of contact between the ball of radius $\|\sigma\|_{TV}$ and the affine space $\{\mu \in \mathcal{M}, K_n(\mu) = K_n(\sigma)\}$, where n is greater than a bound depending only on the structure of the Jordan support of σ .

Definition 2.2. Extrema Jordan Type A signed measure μ is of **extrema Jordan type** (with respect to a family $F = \{u_0, u_1, \dots, u_n\}$) **if and only if** its Jordan decomposition $\mu = \mu^+ - \mu^-$ satisfies

$$\text{Supp}(\mu^+) \subset E_P^+ \text{ and } \text{Supp}(\mu^-) \subset E_P^-,$$

where $\text{Supp}(\nu)$ is defined as the support of the measure ν , and

- P denotes any linear combination of elements of F ,
- P is not constant and $\|P\|_\infty \leq 1$,
- E_P^+ (resp. E_P^-) is the set of all points x_i such that $P(x_i) = 1$ (resp. $P(x_i) = -1$).
- In the paper, they give exact reconstruction results for the following three kinds of *extrema Jordan type* measures:
 - **Nonnegative measures:** Assume that F is an homogeneous M-system. Then, *any nonnegative measure* σ is the **unique** solution to GME given the observation $K_n(\sigma)$, where n is not less than twice the size of the support of *sigma*.
 - **Generalized Chebyshev measures:** Assume that F is a M-system. *let σ be a signed measure having Jordan support included in $(E_{\mathfrak{S}_k}^+, E_{\mathfrak{S}_k}^-)$* , for some $1 \leq k \leq n$, where \mathfrak{S}_k denotes the k -th generalized Chebyshev polynomial. **Then** σ is the **unique** solution to GME given the observation $K_n(\sigma)$.
 - **Δ -interpolation:** Considering $F_P^n = \{1, x, x^2, \dots, x^n\}$, the standard family, *GME exactly recovers any Δ -spaced out type measure* σ from the observation $K_n(\sigma)$, where n is greater than a bound depending only on Δ .
- Let σ be an extrema Jordan type measure. Then σ is a point of contact between the ball of radius $\|\sigma\|_{TV}$ and the affine space $\{\mu \in M, K_n(\mu) = K_n(\sigma)\}$, where n is greater than a bound depending only on the structure of the Jordan support of σ .

Generalized dual polynomial

- The existence of a generalized dual polynomial is a sufficient condition for the exact reconstruction of a *signed measure with finite support*.

Theorem 2.3. GME, [DG11] Let \mathcal{F} be an homogeneous M -system on I . Consider a non-negative measure σ with finite support included in I . Then, the measure σ is the unique solution to generalized minimal extrapolation given observation $K_n(\sigma)$, where n is not less than twice the size of support of σ .

Every measure with support size s depends on $2s$ parameters (s for its support and s for its weights). Surprisingly, this information can be recovered from only $2s + 1$ of its generalized moments.

Lemma 1. The generalized dual polynomials Let n be a positive integer. Let $S = \{x_1, \dots, x_s\} \subset I$ and $(\varepsilon_1, \dots, \varepsilon_s) \in \{\pm 1\}^s$ **If** there exist a linear combination $P = \sum_{k=0}^n a_k u_k$ such that

1. the generalized Vandermonde system

$$\begin{pmatrix} u_0(x_1) & u_0(x_2) & \cdots & u_0(x_s) \\ u_1(x_1) & u_1(x_2) & \cdots & u_1(x_s) \\ \vdots & \vdots & & \vdots \\ u_n(x_1) & u_n(x_2) & \cdots & u_n(x_s) \end{pmatrix}$$

has full column rank,

2. $P(x_i) = \varepsilon_i, \forall i = 1, \dots, s,$
3. $|P(x)| < 1, \forall x \in [-1, 1] \setminus S$

Then every measure $\sigma = \sum_{i=1}^s \sigma_i \delta_{x_i}$, such that $\text{sgn}(\sigma_i) = \varepsilon_i$ is the **unique** solution of GME given the observation $K_n(\sigma)$.

- The linear combination P considered above is called a **generalized dual polynomial**.

2.4.1 Reconstruction of a cone

- GME recovers exactly all measures σ of which support is included in $S = \{x_1, \dots, x_s\}$ and such that $\text{sgn}(\sigma_i) = \varepsilon_i$ for all nonzero σ_i **if** the generalized interpolation problem defined in the lemma above has solution.

- Let us denote this set by $C(x_1, \varepsilon_1, \dots, x_s, \varepsilon_s)$. It is exactly the *cone* defined by

$$C(x_1, \varepsilon_1, \dots, x_s, \varepsilon_s) = \left\{ \sum_{i=1}^s \mu_i \delta_{x_i} \mid \forall \mu_i \neq 0, \operatorname{sgn}(\mu_i) = \varepsilon_i \right\}$$

Thus, the existence of P implies the exact reconstruction of **all** measures in this cone.

- The cone $C(x_1, \varepsilon_1, \dots, x_s, \varepsilon_s)$ is the conic span of an $(s-1)$ -dimensional face of the TV -unit ball, that is

$$F(x_1, \varepsilon_1, \dots, x_s, \varepsilon_s) = \left\{ \sum_{i=1}^s \varepsilon_i \lambda_i \delta_{x_i} \mid \forall i, \lambda_i \geq 0 \text{ and } \sum_{i=1}^s \lambda_i = 1 \right\}.$$

Furthermore, the affine space $\{\mu, K_n(\mu) = K_n(\sigma)\}$ is tangent to the TV -unit ball at any point $\sigma \in F(x_1, \varepsilon_1, \dots, x_s, \varepsilon_s)$, as we can see in the following remark.

- **Remark.** From the convex optimization point of view, the **dual certificates** and the generalized dual polynomials are deeply related: the existence of an generalized dual polynomial P implies that, for all $\sigma \in F(x_1, \varepsilon_1, \dots, x_s, \varepsilon_s)$, a subgradient Φ_P of the TV -norm at the point σ is perpendicular to the set of the feasible points, that is

$$\{\mu, K_n(\mu) = K_n(\sigma)\} \subset \ker(\Phi_P)$$

where \ker denotes the nullspace.

- The condition 2 and 3 in the Lemma 1 ensure that the solutions to GME belong to the cone $C(x_1, \varepsilon_1, \dots, x_s, \varepsilon_s)$, whereas condition 1 gives uniqueness.

2.5 Exact reconstruction of the nonnegatives measures

- They show that if the underlying family $F = \{u_0, u_1, \dots, u_n\}$ is an *homogeneous M -system*, **then** GME recovers exactly each **finitely supported nonnegative measure** μ from the observation of a surprisingly few generalized moments.

2.5.1 Markov systems

Definition 1. T-systems of order k . [SS66] Denote by $\{u_0, u_1, \dots, u_k\}$ a set of continuous real (or complex) functions on \bar{I} . This set is a T-system of degree k if and only if every generalized polynomial

$$P = \sum_{l=0}^k a_l u_l$$

has at most k zeros in I , where $(a_0, a_1, \dots, a_k) \neq (0, 0, \dots, 0)$.

- A finite combination of elements of a T-system is called a **generalized polynomial**.
- This definition is equivalent to each of the two following conditions
 1. For all x_0, x_1, \dots, x_k distinct elements of I and all y_0, y_1, \dots, y_k real (or complex) numbers, there exists a unique generalized polynomial P such that $P(x_i) = y_i$, for all $i = 0, 1, 2, \dots, k$.
 2. For all x_0, x_1, \dots, x_k distinct elements of I , generalized Vandermonde system

$$\begin{pmatrix} u_0(x_1) & u_0(x_2) & \cdots & u_0(x_k) \\ u_1(x_1) & u_1(x_2) & \cdots & u_1(x_k) \\ \vdots & \vdots & & \vdots \\ u_k(x_1) & u_k(x_2) & \cdots & u_k(x_k) \end{pmatrix}$$

has full rank.

Definition 2.4. M-system. We say that the family $F = \{u_0, u_1, \dots, u_n\}$ is an M-system if and only if it is a T-system of degree k for all $0 \leq k \leq n$

Definition 2.5. We say that the family $F = \{u_0, u_1, \dots, u_n\}$ is an homogeneous M-system if and only if it is an M-system and u_0 is a constant function.

- Using homogeneous M-systems, they show that one can exactly recover all non-negative measures from a few generalized moments.

Theorem 2.6. *Let F be an homogeneous M-system on I . Consider a non-negative measure σ with finite support included in I . Then the measure σ is the **unique** solution to GME given the observation $K_n(\sigma)$, where n is not less than twice the size of the support of σ .*

Non-negative interpolation An important property of M-systems is the existence of a non-negative generalized polynomial that vanishes exactly at a prescribed set of points $\{t_1, \dots, t_m\}$, where $t_i \in I$ for all $i = 1, \dots, m$. Indeed, define the index as

$$\text{Index}(t_1, \dots, t_m) = \sum_{i=1}^m \chi(t_i) \quad (2.20)$$

where

$$\chi(t) = \begin{cases} 2 & \text{if } t \in \overset{\circ}{I} \\ 1 & \text{otherwise} \end{cases} \quad (2.21)$$

and where $\overset{\circ}{I}$ denotes the interior of I .

Lemma 2.7. *Non-negative generalized polynomial Consider*

Deterministic matrices for compressed sensing. In the following, p denotes the number of predictors, or from a signal processing view point, the length of the signal.

- *Deterministic Design:* For

$$O_{p,s \rightarrow \infty} \left(s \log \left(\frac{p}{s} \right) \right) \quad (2.22)$$

there exists a deterministic matrix $A \in \mathbb{R}^{n \times p}$ such that basis pursuit recovers all s -sparse vectors from the observation Ax_0 .

- *Random design:* If

$$n \geq Cs \log \left(\frac{p}{s} \right) \quad (2.23)$$

where $C > 0$ is a universal constant, there exists (with high probability) a random matrix $A \in \mathbb{R}^{n \times p}$ such that basis pursuit recovers all s -sparse vectors from the observation Ax_0 .

Considering non-negatives sparse vectors, it is possible to drop the bound on n to $2s + 1$.

Theorem 2.8. *Deterministic design matrix. Let n, p , and s be integers such that $s \leq \min(n/2, p)$. Let A be generalized Vandermonde system defined by*

$$A = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ u_1(t_1) & u_1(t_2) & \cdots & u_1(t_p) \\ \vdots & \vdots & & \vdots \\ u_n(t_1) & u_n(t_2) & \cdots & u_n(t_p) \end{pmatrix} \quad (2.24)$$

then, basis pursuit exactly recovers all non-negative s -sparse vectors $x_0 \in \mathbb{R}^p$ from the observation Ax_0 .

The program basis pursuit can be recast as a linear program. Then, we use an interior point method to solve basis pursuit.

Lets us denote $K : t \mapsto (1, u_1(t), \dots, u_n(t))$. The columns of A are the values of this map at points t_1, \dots, t_p .

The nullspace property for measures. We aim at deriving a sufficient condition for exact reconstruction of signed measures. Note that the solutions to program (GME) depend only on the first $n + 1$ elements of \mathcal{F} and on the target measure σ . We investigate the condition that the family \mathcal{F} must satisfy to ensure exact reconstruction (nullspace property).

The nullspace property for generalized minimal extrapolation. Consider the linear map $K_n : \mu \mapsto (c_0(\mu), \dots, c_n(\mu))$ from μ to \mathbb{R}^{n+1} that defines the generalized moment morphism.

Its nullspace $\ker(K_n)$ is a linear subspace of \mathcal{M} . The Lebesgue decomposition theorem is the precious tool used to define the nullspace property.

Definition 2.9. Nullspace property with respect to a Jordan family Y We say that the generalized moment morphism K_n satisfies the nullspace property with respect to a Jordan support family Y if and only if it satisfies the following property. For all nonzero measure μ in the nullspace of K_n , and for all $(S^+, S^-) \in Y$

(nullspace property)

$$\| \mu_S \| < \| \mu_{S^c} \| \tag{2.25}$$

(weak nullspace property)

$$\| \mu_S \| \leq \| \mu_{S^c} \| \tag{2.26}$$

where $S = S^+ \cup S^-$

2.6 BLASSO

Super-resolution

- **The Super-resolution phenomenon** is the ability to recover the information beyond the physical limitations[EC12].
- This paper offers quantitative detection guarantees from noisy observations.
- The authors provide a tractable algorithm (**BLASSO**) and quantitative estimates of a train of complex valued spikes from very few noisy observations.

- Their analysis involves an estimate of the magnitude of the noise perturbation in the signal domain using the **Rice method**. In particular, they derive explicit bounds for tuning parameter appearing in BLASSO.

General model and notation

Theorem 2.10. *General model and notation*

Let \mathbb{T} be a compact set homeomorphic to either the interval $[0, 1]$ or the unit circle \mathbb{S}^1 . Let Δ be a complex measure on \mathbb{T} with discrete support of (unknown) size s . In particular, Δ has polar decomposition:

$$\Delta = \sum_{k=1}^s \Delta_k \exp(i\theta_k) \delta_{T_k}, \quad (2.27)$$

where $\Delta_k > 0$, $\theta_k \in \mathbb{R}$, $T_k \in \mathbb{T}$ for $k = 1, \dots, s$ and δ_x denotes the Dirac measure at point x .

Let m be a positive integer and $F = \{\phi_0, \phi_1, \dots, \phi_m\}$ be a family of complex continuous functions on \mathbb{T} . Define the k -th generalized moment of a complex measure μ on \mathbb{T} as:

$$c_k(\mu) = \int_{\mathbb{T}} \phi_k d\mu, \quad (2.28)$$

for all indices $k = 0, 1, \dots, m$. Assume that we observe $y = (y_k)_{k=0}^m$ defined as:

$$y = \int_{\mathbb{T}} \Phi d\mu + \varepsilon, \quad (2.29)$$

where $\Phi = (\phi_0, \phi_1, \dots, \phi_m)$. We aim at reconstructing the complex measure Δ from the $m + 1$ measurements given by y .

Definition 2.11. Beurling LASSO (BLASSO) [ADG14] Denote by \mathcal{M} the set of finite complex measures on \mathbb{T} and by $\|\cdot\|_{TV}$ the total variation norm. We recall that for all $\mu \in \mathcal{M}$,

$$\|\mu\|_{TV} = \sup_{\pi \in \Pi} \sum_{E \in \pi} |\mu(E)| \quad (2.30)$$

by analogy with LASSO, Beurling LASSO (BLASSO) is the process of reconstructing a discrete measure Δ from the samples y by finding a solution to

$$\hat{\Delta} \in \arg \min_{\mu \in \mathcal{M}} \frac{1}{2} \left\| \int_{\mathbb{T}} \Phi d\mu - y \right\|_2^2 + \lambda \|\mu\|_{TV} \quad (2.31)$$

where λ is a tuning parameter.

Fendrel dual program: The usual convex analysis shows that BLASSO can be viewed as a Fendrel dual program. As a matter of fact, any solution to BLASSO can be faithfully computed from a companion program that builds a dual certificate of $\hat{\Delta}$.

Definition 2.12. The problem

$$\min_{a \in \mathbb{C}^{m+1}} \frac{\|a - y\|_2^2}{2} + \mathbb{I}_{\{a \in \mathbb{C}; \|\sum_{k=0}^m a_k \psi_k\|_\infty \leq \lambda\}}(a) \quad (2.32)$$

has its Fendrel dual with the same minimizer as BLASSO. Here, the indicator $\mathbb{I}_E(v)$ of a set $E \subset \mathbb{C}$ is defined by $\mathbb{I}_E(v) = 0$ if $v \in E$ and $\mathbb{I}_E(v) = +\infty$ otherwise.

Using the predual problem (2.32), it is possible to derive optimality conditions for BLASSO. Hence, we mention that all solution to BLASSO is SM.

Proposition 2.13. [BP 10] *The optimization problem (BLASSO) admits at least a solution. Moreover, all solution $\hat{\Delta}$ is SM and it has a dual certificate $\hat{P} = \sum_{k=0}^m \hat{a}_k \psi_k$ where*

$$\forall k \in \{0, \dots, m\}, \quad \hat{a}_k = \frac{c_k (\hat{\Delta} - y_k)}{\lambda} \quad (2.33)$$

Remark: We have an explicit formulation of a dual certificate \hat{P} of $\hat{\Delta}$ using 2.33. Moreover, all solution to BLASSO is discrete, SM and satisfies

$$\{x \in \mathbb{T}, |\hat{\Delta}(\{x\})| > 0\} \subseteq \{x \in \mathbb{T}, |\hat{P}(x)| = 1\}. \quad (2.34)$$

In other words, the support of $\hat{\Delta}$ is included in the set of the points for which $|\hat{P}|$ is maximal.

On the algorithmic side, the program (2.32) allows us to compute a dual certificate of a solution $\hat{\Delta}$ to (BLASSO). As a matter of fact, it takes the form:

$$\hat{a} \in \arg \min_{a \in \mathbb{C}^{m+1}} \left\| a - \frac{y}{\lambda} \right\|_2^2 \quad \text{subject to} \quad \left\| \sum_{k=0}^m a_k \psi_k \right\|_\infty \leq 1 \quad (2.35)$$

Once the support is estimated accurately, a solution to (BLASSO) can be found by solving a well-posed linear problem.

2.7 The Quadrature Method of Moments

In (Gordon, 1968) [Gor68], it was developed a methodology for calculating a Gaussian quadrature whose weight function is an arbitrary distribution function whose support belongs to $[0, \infty)$. This quadrature was applied by (McGraw, 1997) [McG97], to develop the *quadrature method of moments* (QMOM) for numerically solving the population balance equation [Lag07].

The Gordon quadrature of k points can be derived for any given continuous or discrete distribution for which the $2k$ first standard moments can be calculated, using the *Product-Difference Algorithm (PDA)*.

The calculation of a k -points Gordon quadrature implies in the determination of k weights (W_i) and k abscissas (L_i). This $2k$ variables can be found using the PDA from the low order moments. The PDA is based on minimizing the error tasked when we replace the integral in equation (2.8) with its quadrature approximation ([Mar+03],[Lag07]).

2.7.1 The Product Difference Algorithm

The PDA consists in two principal steps. First, a upper triangular matrix P of size $(2k + 1) \times (2k + 1)$. Then, the second step consists in the calculation of the vector α which is used to compute a tridiagonal symmetric Jacobi matrix J of size $k \times k$. The weights and abscissas are computed using the characteristic values of this matrix and the characteristic vectors of J ([Lag07], [Gor68]).

Construction of the matrix P .

We start with the construction of the first column of P

$$P_{i,1} = \delta_{i1}, \quad i = 1, \dots, 2N + 1; \quad \delta_{i1}: \text{Kronecker delta} \quad (2.36)$$

then, the second column is

$$P_{i,2} = (-1)^{i-1}, \quad i = 1, \dots, 2N + 1 \quad (2.37)$$

Since the final weight can be corrected multiplying by the real value of m_0 , the computations are made assuming a normalized distribution (exp. $m_0 = 1$). The rest of the components are found using the PD algorithm

$$P_{i,j} = P_{1,j-1}P_{i+1,j-2} - P_{1,j-2}P_{i+1,j-1} \quad (2.38)$$

where $\begin{cases} j = 3, \dots, 2N + 1 \\ i = 1, \dots, 2N + 2 - j \end{cases}$

Calculation of the weights and abscissas.

The continuous fraction coefficients (α_i) are computed assigning $\alpha_1 = 0$ and the rest of them can be computed recursively by the following relation

$$\alpha_i = \frac{P_{1,i+1}}{P_{1,i}P_{1,i-1}}, \text{ for } i = 2, \dots, 2N \quad (2.39)$$

A symmetric tridiagonal matrix can be get from the α_i 's using

$$\begin{aligned} a_i &= \alpha_{2i} + \alpha_{2i-1}, \text{ for } i = 1, \dots, N \\ b_i &= \sqrt{\alpha_{2i+1}\alpha_{2i}}, \text{ for } i = 1, \dots, N - 1 \end{aligned} \quad (2.40)$$

where a_i and b_i are the diagonal and codiagonal of the Jacobi matrix respectively. The weights and the abscissas are determined finding the eigenvalues and eigenvectors of the Jacobi matrix.

$$\begin{aligned} L_i &= \text{eigenvalues} \\ W_i &= m_0 v_{i,1}^2 \\ \text{where } v_{i,1} &= \text{first component of the } i\text{-th eigenvector} \end{aligned} \quad (2.41)$$

2.8 Examples of Implementation of QMOM

2.8.1 Theoretical moments of the case of Silva 2010

In ([Sil+10],[Sco67],[PA98]) the authors studied particular cases where an theoretical solution to the Population Balance Equation exists. In order to investigate the behavior of methods for recovering the time evolution of a finite number of standard moments, we compute the standard moments of the distribution given in (Silva 2010) as theoretical solution for the PBE.

The authors studied three cases where they could find an analytic solution. The first one modelize systems where the aggregation and breakage processes were equally important. The second case modelize systems where the aggregation processes are significantly more important than breakage processes and finally the third case modelize systems where the breakage processes are more important than aggregation. We have then

Proposition 2.14. *Theoretical moments Following [Sil+10],[Vli14], we have a system with the following initial conditions: $\mu_0(0) = 1$, $\mu_1 = 1$, and $C = 1$. The initial distributions are*

$$\begin{aligned} n_1(L, 0) &= 3L^2 e^{-L^3} \\ n_2(L, 0) &= 12L^5 e^{-2L^3} \end{aligned} \quad (2.42)$$

For the initial conditions, the initial moments are respectively

$$\begin{aligned} m_k(t)_1 &= \Gamma\left(\frac{k}{3} + 1\right), \\ m_k(t)_2 &= \frac{1}{2^{\frac{k}{3}}}\Gamma\left(\frac{k}{3} + 2\right). \end{aligned} \quad (2.43)$$

- Case $\Phi(\infty) = 1$ (aggregation and breakage equally important).
In this case, the theoretical solution is

$$n^a(L, t) = \sum_{i=1}^2 \frac{K_1(t) + P_i(t) K_2(t)}{L_2(t) + 4P_i(t)} 3L^2 e^{P_i(t)L^3}, \quad t > 0 \quad (2.44)$$

and the standard moments can be computed as

$$\mu_k(t) = \Gamma\left(\frac{k}{3} + 1\right) \sum_{i=1}^2 \frac{1}{(-P_i(t))^{k/3+1}} \frac{K_1(t) + P_i(t) K_2(t)}{L_2(t) + 4P_i(t)} \quad (2.45)$$

- Case $\Phi(\infty) \neq 1$ (aggregation more important than breakage if $\Phi(\infty) > 1$, and breakage more important than aggregation if $\Phi(\infty) < 1$).
In this case, the theoretical solution is

$$n^a(L, t) = 3L^2 [\Phi(t)]^2 e^{-\Phi(t)L^3}, \quad t > 0, \quad (2.46)$$

and the standard moments can be computed as

$$\mu_k(t) = \frac{\Gamma(k/3 + 1)}{[\Phi(t)]^{k/3-1}} \quad (2.47)$$

Proof. We use the same ideas as ([Vli+11], [MF05]) in order to express the number distribution function in terms of the size of the particle the population bilan equation

- Case $\Phi(\infty) = 1$. We have that the standard moments of the theoretical solution is defined as

$$\begin{aligned} \mu_k(t) &= \int_0^\infty L^k n^a(L, t) dL \\ &= \int_0^\infty L^k \left[\sum_{i=1}^2 \frac{K_1(t) + P_i(t) K_2(t)}{L_2(t) + 4P_i(t)} 3L^2 e^{P_i(t)L^3} \right] dL \end{aligned} \quad (2.48)$$

if we call

$$\Lambda_i(t) = \frac{K_1(t) + P_i(t) K_2(t)}{L_2(t) + 4P_i(t)} \quad (2.49)$$

which is a function that depends only on time t . Then, we can write this integral as

$$\mu_k(t) = \Lambda_1(t) \int_0^\infty L^k 3L^2 e^{P_1(t)L^3} dL + \Lambda_2(t) \int_0^\infty L^k 3L^2 e^{P_2(t)L^3} dL. \quad (2.50)$$

The two integrals involved can be evaluated using the Gamma function as

$$\mu_k(t) = \Lambda_1(t) \frac{\Gamma(k/3 + 1)}{[-P_1(t)]^{k/3+1}} + \Lambda_2(t) \frac{\Gamma(k/3 + 1)}{[-P_2(t)]^{k/3+1}}, \quad (2.51)$$

and finally we obtain

$$\mu_k(t) = \Gamma(k/3 + 1) \sum_{i=1}^2 \frac{3}{(-P_i(t))^{k/3+1}} \frac{K_1(t) + P_i(t) K_2(t)}{L_2(t) + 4P_i(t)} \quad (2.52)$$

- Case $\Phi(\infty) \neq 1$. In a similar way we have

$$\begin{aligned} \mu_k(t) &= \int_0^\infty L^k n^a(L, t) dL \\ &= \int_0^\infty 3L^{k+2} [\Phi(t)]^2 e^{-\Phi(t)L^3} dL. \end{aligned} \quad (2.53)$$

The integral involved can be evaluated using the Gamma function and finally we obtain

$$\mu_k(t) = \frac{\Gamma(k/3 + 1)}{[\Phi(t)]^{k/3+1}}, \quad (2.54)$$

□

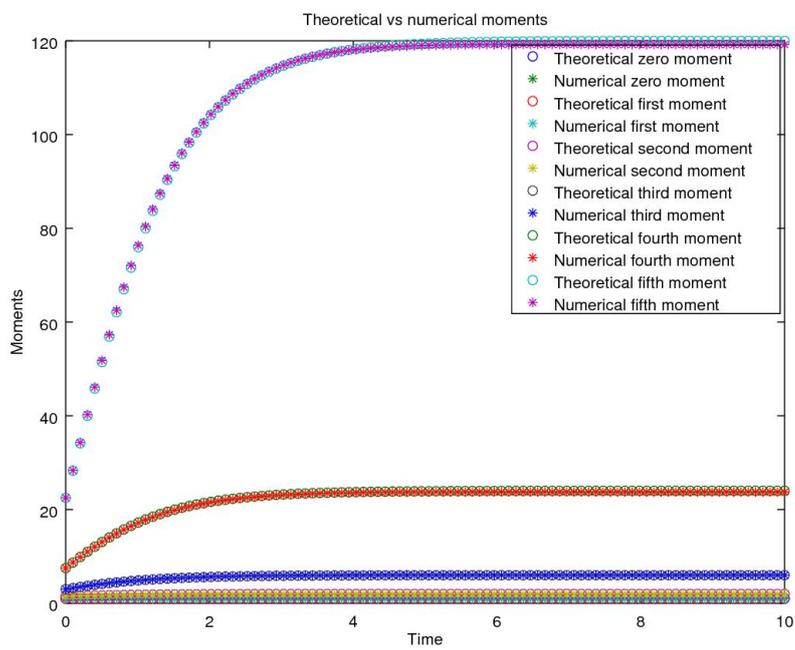


Figure 2.1: First 6 standard moments of the case of (Silva, 2011) (strength line) and the estimation with QMOM (dotted line)

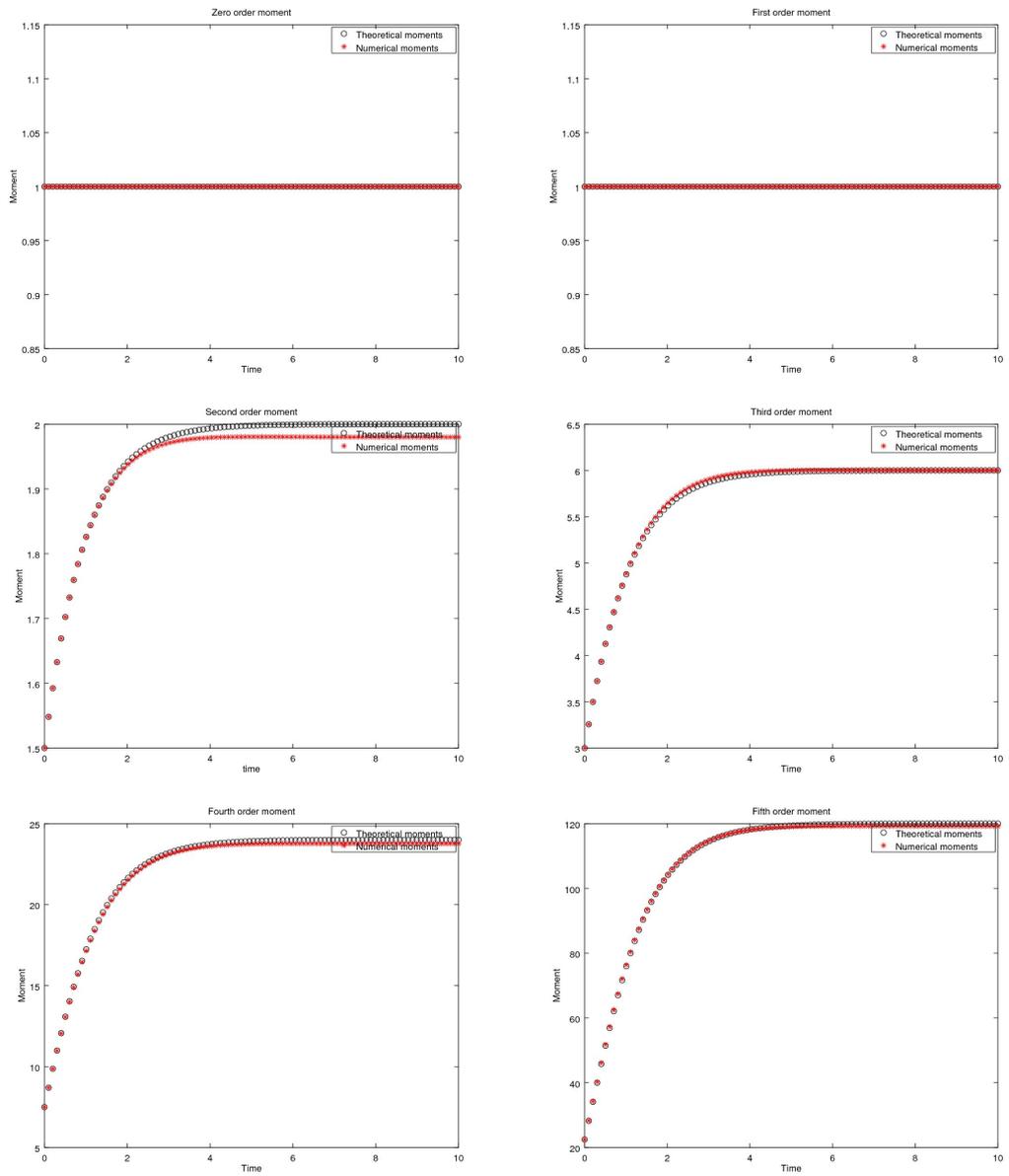


Figure 2.2: First 6 standard moments of the case of (Silva, 2011) (strength line) and the estimation with QMOM (dotted line)

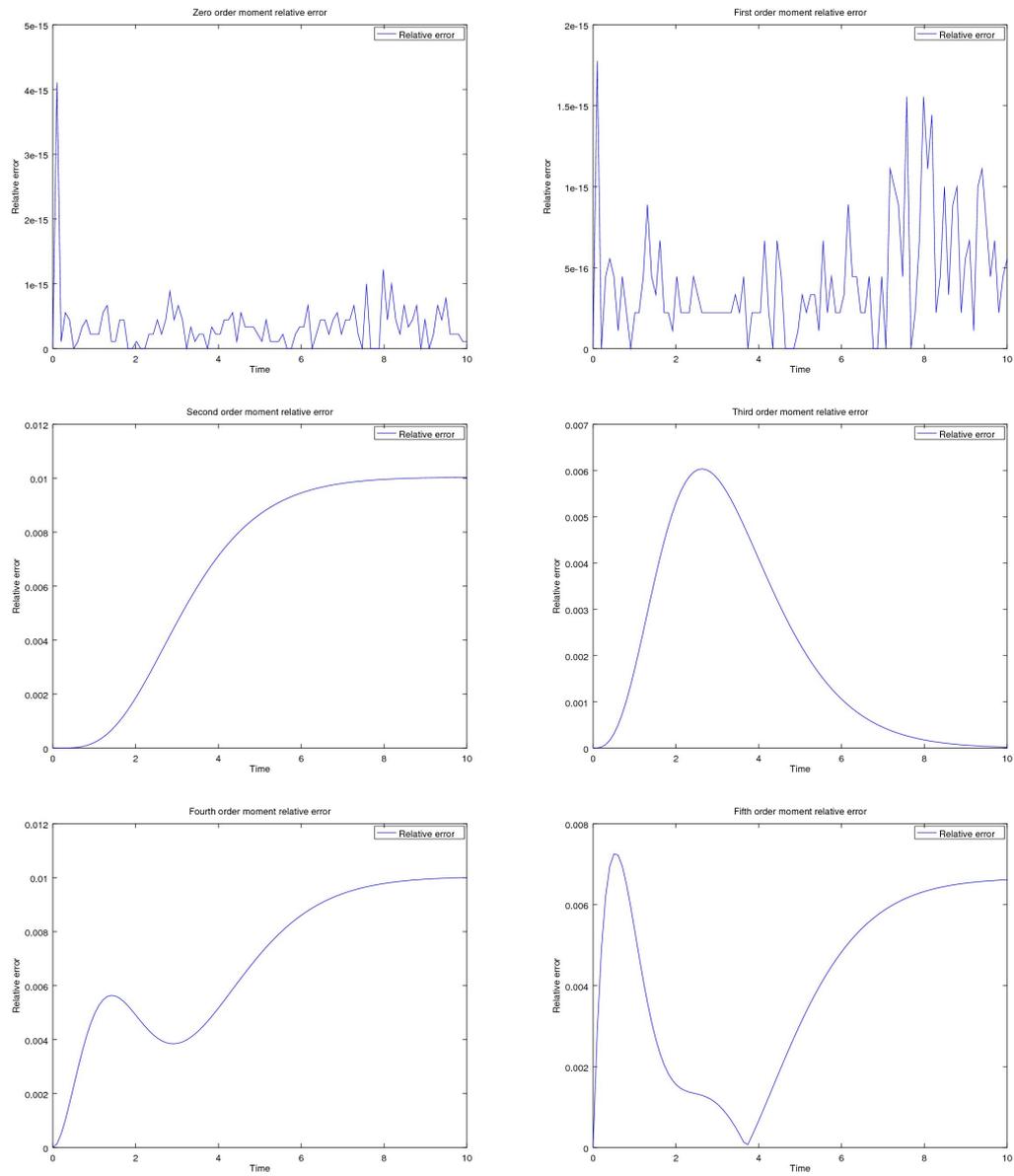


Figure 2.3: First 6 standard moments of the case of (Silva, 2011) (strength line) and the estimation with QMOM (dotted line)

Chapter 3

Numerical resolution for the Population Balance Equation

3.1 Introduction and theoretical considerations

Particles are encountered in an innumerable variety of systems. The particles are either naturally presented in these systems or engineered into them. In either case, the particles often significantly affect the behavior of such systems. In many other situations, systems are associated with processes in which particles are formed either as a main product or as a by-product. We will refer to systems containing particles as dispersed phase systems or particulate systems regardless of the precise role of the particles in them.

Population balances are essential to scientists and engineers of wide varying of disciplines. They are of interest to physicist (astrophysicist, high-energy, geophysicists, meteorologists) and chemists (colloidal chemists, statistical mechanicians). Biophysicists concerned with populations of cells of various kinds, food scientists dealing with preparations of emulsions or sterilization of food all have an indispensable need for population balances.

Among engineers, population balance concepts are of importance to aeronautical, chemical, civil (environmental), mechanical, and materials engineers. Chemical engineers have put population balances to the most diverse use. Applications have covered a wide range of dispersed phase systems, such as solid-liquid dispersions (although with incidental emphasis on crystallization systems), and gas-liquid, gas-solid, and liquid-liquid dispersions [Soo+07a].

Although most of the foregoing applications are known, it is significant to cite more modern applications such as the preparation of ceramic mixtures and fine particles (nanoparticles) for a variety of applications, in which

population balances play a critical role in the analysis, design and control of such processes. For example, the manufacture of superconducting ceramic mixtures requires very tight specifications on their composition on a fine scale of mixing. The knowledge of the time evolution of the number particle distribution is also important in many applications like in computational fluid dynamic simulation [Mar+03].

Analysis of particulate systems seeks to synthesize the behavior of the population of particles and its environment from the behavior of single particles in their local environments. In the application of population balances, one is more interested in the distribution of particles populations and their effect on the system behavior. The system of interest is that they contain particles which are continually being created and destroyed by processes such as particles breakage and agglomeration [Van00].

Fundamental to the formulation of population balances is the assumption that there exist a number density of particles at every point of the particle state space. The number of particles in any region of the state space is obtained by integrating the number density over the region desired.

The theoretical solution of population balance equation means to solve the integral-partial differential equation involved. Except for some cases, it is not straightforward to obtain such solution. The efforts for solving this equation comprise numerical integration and the discretization of the range of the considered property (like volume or length), solving the equation numerically in some intervals. Also we can find methods solving the equation for the evolution of a set of moments of the number distribution. The application of methods finding the time evolution of moments often involves to work with no closed equations.

We propose a discretization scheme in order to find a numerical approximation to the solution of population balance equations involving only aggregation and breakage processes acting as particle's modifiers. This numerical solution is compared to some known theoretical solution for simple aggregation and breakage kernels. Also, we consider a more general case of aggregation and breakage kernels. We compare the performance of the numerical approximation to an empirical estimation found in (?) from computed from experimental data.

3.1.1 The framework of Population Balance

We are concerned with systems consisting of particles dispersed in an environmental phase which we shall refer to as the continuous phase.

The particles may interact between themselves as well as with the continuous phase. Such behavior may vary from particle to particle depending

upon a number of "properties" that may be associated with the particle. Continuous variables may be encountered more frequently in population balance analysis. The external coordinates denoting the position vector of (the centroid of) a particle describing continuous motion through space represent continuous variables.

The temporal evolution of the particulate system, we shall regard time as varying continuously and inquire into the rate of change of the particle state variables. It is convenient to deal with continuous variables in this regard. A fundamental assumption here is that the rate of change of state of any particle is a function only of the state of the particle in question and the local continuous phase variables. Thus we exclude the possibility of direct interactions between particles, although indirect interaction between particles via the continuous phase is indeed accounted for because of the dependence of particle behavior on the "local" continuous phase variables. In order to enable such a local characterization of the continuous phase variables, it is necessary to assume that the particles are considerably smaller than the length scale in which the continuous phase quantities vary. The continuous phase variables may be assumed to satisfy the usual transport equations with due regard to interaction with the particulate phase. Thus, such transport equations will be coupled with the population balance equation.

Particle state vector We are concerned with particle phase variables that are continuous. In general, the choice of the particle state is determined by the variable needed to specify:

- The rate of change of those of direct interest to the application, and
- The birth and death processes.

The particle state may generally be characterized by a finite dimensional vector.

- External coordinate $\mathbf{r} \equiv (r_1, r_2, r_3)$ denote the position (of the centroid) of the particle.
- Internal coordinates $\mathbf{x} \equiv (x_1, \dots, x_d)$ representing d different quantities associated with the particle.

The particle state vector (\mathbf{x}, \mathbf{r}) accounts for both internal and external coordinates. We shall let Ω_x represent the domain of internal coordinates and Ω_r be the domain of external coordinates, which is the set of points in physical space in which the particles are present. These domains may be bounded or may have infinite boundaries.

The particle population may be regarded as being randomly distributed in the particle state space, which include both external or internal coordinates.

Our concern will be about large populations, which will display relatively deterministic behavior because the random behavior of individual particles will be averaged out.

The continuous phase vector. The continuous phase variables may be collated into a finite c -dimensional vector field. The continuous phase variables affect the behavior of each particle.

We define a continuous phase vector. $Y(r, t) = [Y_1(r, t), \dots, Y_c(r, t)]$ which is clearly a function only of the external coordinate r and time t .

The evolution of this field in space and time is governed by the laws of transport and interaction with the particles.

In some applications, a continuous phase balance may not be necessary because interactions between the population and the continuous phase may not bring about any (or a substantial enough) change in the continuous phase. In such case, analysis of the population involves only the population balance equation.

The number density function. We postulate that there exist an average number density function defined on the particle state space,

$$E[f_1(x, r, t)] \equiv n(x, r, t)$$

with $x \in \Omega_x$ and $r \in \Omega_r$, where $E[f_1(x, r, t)]$ denote the expectation or the average of the actual number density $f_1(x, r, t)$, while $n(x, r, t)$ denotes the average number density. This definition implies that the average number of particles in the infinitesimal volume $dV_x dV_r$ (in the particle state space) about the particle state (x, r) is $n(x, r, t) dV_x dV_r$. However, we will refer to particles in volume $dV_x dV_r$ about the particle state (x, r) .

The average number density $n(x, r, t)$ is assumed to be sufficiently smooth to allow differentiation with respect to any of its arguments as many times as may become necessary.

The (average) number density allows one to calculate the (average) number of particles in any region of particle state space. Thus, the (average) total number of particles in the entire system is given by

$$\int_{\Omega_x} dV_x \int_{\Omega_r} dV_r n(x, r, t)$$

where dV_x and dV_r are infinitesimal volume measures in the spaces of internal and external coordinates respectively.

The local (average) number density in physical space, i. e. the (average) total number of particles per unit volume of physical space, denoted $N(r, t)$ is given by

$$N(r, t) = \int_{\Omega_x} dV_x n(x, r, t).$$

Other densities such as volume or mass density may also be defined for the particle population. Thus, if $v(x)$ is the volume of the particle of internal state x , then the volume density may be defined as $v(x) f_1(x, r, t)$.

The volume fraction density $\phi(x, r, t)$ of a particle state is defined by

$$\phi(x, r, t) = \frac{1}{\Phi(r, t)} v(x) n(x, r, t)$$

where

$$\Phi(r, t) = \int_{\Omega_x} dV_x v(x) n(x, r, t)$$

the denominator above represents the total volume fraction of all particle. Similarly, mass fractions can also be defined. For the case of scalar interval state using volume, the volume fraction density of particles of volume v becomes

$$\phi(v, r, t) = \frac{vn(v, r, t)}{\Phi(r, t)}$$

where

$$\Phi(r, t) = \int_0^\infty vn(v, r, t) dv.$$

In contrast with number density, volume or mass denote the amount of dispersed phase material.

The rate of change of particle state vector We observe earlier that particle state might vary in time. We are concerned with smooth changes in particle state describable by some vector field defined over the particle state space both internal and external coordinates.

While changes of external coordinates refers to motion through physical space, that of internal coordinates refers to motion through an abstract property space (for example size).

We had collectively referred to thm as convective processes for the reasons that they might be likened to physical motion.

It will be convenient to define "velocity" $\dot{R}(x, r, Y, t)$ for internal coordinates and $\dot{R}(x, r, Y, t)$ for external coordinates. These functions are assumed to be as smooth as necessary.

The velocity just defined may be random processes in space and time. Thus, $n(x, r, t) \dot{R}(x, r, Y, t)$ represents the particle flux through physical space and $n(x, r, t) \dot{X}(x, r, Y, t)$ is the particle flux through internal coordinate space.

The Population balance equation. Consider a population of particles distributed according to their size x which we shall take to be the mass of the particle and allow it to vary between 0 and ∞ .

The particles are uniformly distributed in space so that the number density is independent of external coordinates. Further, we assume for the present that the environment does not play any explicit role in particle behavior.

We let $\dot{X}(x, t)$ be the growth rate of the particle size x and let $n(x, t)$ denote the number density. All functions involved are assumed to be sufficiently smooth. Thus, we have the population balance equation

$$\frac{\partial n(x, t)}{\partial t} + \frac{\partial \dot{X}(x, t) n(x, t)}{\partial x} = 0.$$

In the above derivation, we did not take in account the birth and death of particles. To assess the rates of these contributions detailed modeling of breakage and aggregation processes will be needed. Let $h(x, t) dx$ the net rate of generation of particles in the size range x to $x + dx$, where the identity of $h(x, t)$ would depend on the models of breakage and aggregation. In this case, the population balance equation becomes

$$\frac{\partial n(x, t)}{\partial t} + \frac{\partial \dot{X}(x, t) n(x, t)}{\partial x} = h(x, t)$$

The preceding equation must be supplemented with initial and boundary conditions. The initial condition must clearly stipulate the distribution of particles in the particle state space.

The Population Balance Equation (PBE) is an equation that describes the evolution of one population of particles in colloidal systems. Changes in this kind of population are due to aggregation or breakage processes that can be seen as processes of birth and death. The evolution of the population is characterized by the particle size. The formulation of PBE is traditionally made in terms of the particle's volume as size property.

This four terms at the right side of the equation are the corresponding processes of birth and death due to aggregation or breakage.

This equation is expressed like

Definition 3.1. Population Balance Equation The equation governing the evolution in time of the number distribution of a population of colloidal particles is known as Population Balance Equation, and it is defined as

$$\begin{aligned}
\frac{\partial \eta(v; t)}{\partial t} &= B_a(v; t) - D_a(v; t) + B_b(v; t) - D_b(v; t) \\
&= \frac{1}{2} \int_0^v \phi(v - \epsilon, \epsilon) \eta(v - \epsilon; t) \eta(\epsilon; t) d\epsilon \\
&\quad - \eta(v; t) \int_0^\infty \phi(v, \epsilon) \eta(\epsilon; t) d\epsilon \\
&\quad + \int_v^\infty \psi(\epsilon) \rho(v/\epsilon) \eta(\epsilon; t) d\epsilon \\
&\quad - \psi(v) \eta(v; t),
\end{aligned} \tag{3.1}$$

where

- $B_a(v; t) = \frac{1}{2} \int_0^v \phi(v - \epsilon, \epsilon) \eta(v - \epsilon; t) \eta(\epsilon; t) d\epsilon$: birth rate of particles with volume v by aggregation of little particles,
- $D_a(v; t) = \eta(v; t) \int_0^\infty \phi(v, \epsilon) \eta(\epsilon; t) d\epsilon$: death rate of particles with volume v by aggregation with other particles,
- $B_b(v; t) = \int_v^\infty \psi(\epsilon) \rho(v/\epsilon) \eta(\epsilon; t) d\epsilon$: birth rate of particles with volume v by breakage of big particles,
- $D_b(v; t) = \psi(v) \eta(v; t)$: death rate of particles with volume v by breakage into little particles.

and where

- $\eta(v; t)$: number density function using volume as coordinate,
- $\phi(v, \epsilon)$: aggregation kernel using volume as coordinate,
- $\psi(v)$: breakage kernel using volume as coordinate,
- $\rho(v/\epsilon)$: distribution function of fragments.

In some applications, it is interesting to express the PBE in terms of the Length of Diameter particle instead of the volume. Because of this, we are going to see how we can transform the PBE using the particle's volume as distribution variable to the PBE using the particle's size as distribution variable. We can see that this formulation is

Proposition 3.2. *The PBE in terms of size particle coordinate. The PBE can be formulate in terms of the length coordinate like*

$$\begin{aligned}
\frac{\partial n(L, t)}{\partial t} &= B^a(L; t) - D^a(L; t) + B^b(L; t) - D^b(L; t) \\
&= \frac{L^2}{2} \int_0^L \frac{\beta\left((L^3 - \lambda^3)^{1/3}, \lambda\right)}{(L^3 - \lambda^3)^{2/3}} n\left((L^3 - \lambda^3)^{1/3}, t\right) n(\lambda, t) d\lambda \\
&\quad - n(L, t) \int_0^\infty \beta(L, \lambda) n(\lambda, t) d\lambda \\
&\quad + \int_L^\infty a(\lambda) b(L | \lambda) n(L, t) d\lambda \\
&\quad - a(L) n(L, t),
\end{aligned} \tag{3.2}$$

where where

- $B^a(L; t) = \frac{L^2}{2} \int_0^L \frac{\beta\left((L^3 - \lambda^3)^{1/3}, \lambda\right)}{(L^3 - \lambda^3)^{2/3}} n\left((L^3 - \lambda^3)^{1/3}, t\right) n(\lambda, t) d\lambda$: birth rate of particles with length L by aggregation of little particles,
- $D^a(L; t) = n(L, t) \int_0^\infty \beta(L, \lambda) n(\lambda, t) d\lambda$: death rate of particles with length L by aggregation with other particles,
- $B^b(v; t) = \int_L^\infty a(\lambda) b(L | \lambda) n(L, t) d\lambda$: birth rate of particles with length L by breakage of big particles,
- $D^b(v; t) = a(L) n(L, t)$: death rate of particles with length L by breakage into little particles.

and where

- $n(L; t)$: number density function using length as coordinate,
- $\beta(L, \lambda)$: aggregation kernel using length as coordinate,
- $a(L)$: breakage kernel using length as coordinate,
- $b(L/\lambda)$: distribution function of fragments.

Except for some theoretical cases is not possible to find analytically the solution of the Population Balance Equation. In order to find an approximation to this solution, we are going to propose a discretization scheme.

Let L be the particle's Size. This coordinate has its natural variation range in $]0, +\infty[$. We start building a discretization scheme for the coordinate

L . Lets denote \underline{L} a conveniently small value for L and \overline{L} a conveniently large value of L . That is, $L \in [\underline{L}, \overline{L}]$. We choose a discretization scheme, building a geometrical grid of $N + 1$ values. Lets denote s the multiplicative step for the geometrical grid, where with $s = \left(\frac{L_N}{L_0}\right)^{\frac{1}{N}}$. Lets denote $L_0 = \underline{L}$, $L_N = \overline{L}$ and for $i = 1, \dots, N - 1$, $L_i = s^i L_0$.

The PBE is an integro-differential equation of the number distribution. We are going to use the extended trapezoidal rule in the geometrical grid in order to numerically solve the integral parts involved and the Euler's step foward method for the derivative part.

3.1.2 Trapezoidal rule for regular and geometric grid

Now, we want to compare the performance of the trapezoidal rule implemented with a regular grid to the same rule using a geometric grid instead. We are going to use the particular Gamma distribution described before at the beginning because of its 'tailed' behavior and then we are going to use others distributions like the uniform distribution.

We use the theoretical value of the integral of the Gamma function for comparing the performance between the two estimation methods.

The Extended Trapezoidal rule

We want to approximate the value of the integral of certain function between the points a and b with the integral of a linear function between the same rang, that is

$$\int_a^b f(x) dx \approx \int_a^b [\beta x + \alpha] dx.$$

Definition 3.3. The trapezoidal rule[Nak92] The integral of certain function in $[a, b]$ can be approximated by

$$\int_a^b f(x) dx = \left[\frac{b-a}{2}\right] (f(a) + f(b)) + o(1) \quad (3.3)$$

Now, if we want to use $N+1$ points to approximate the theoretical integral value, we use the trapezoidal rule N times, from which we have the Extended

Trapezoidal rule. That is

$$\begin{aligned}
\int_a^b f(x) dx &= \int_{x_0}^{x_1} f(x) dx + \cdots + \int_{x_{N-1}}^{x_N} f(x) dx \\
&= \left[\frac{x_1 - x_0}{2} \right] (f(x_0) + f(x_1)) + \cdots \\
&\quad + \left[\frac{x_N - x_{N-1}}{2} \right] (f(x_{N-1}) + f(x_N)) + o(1) \\
&= \sum_{i=1}^N \left[\frac{x_i - x_{i-1}}{2} \right] (f(x_{i-1}) + f(x_i)) + o(1)
\end{aligned}$$

then, if we use a uniform grid $x_{i+1} = x_i + h$ for $i = 0, \dots, N$, and $h = \frac{x_N - x_0}{N}$ we get

$$\begin{aligned}
\int_a^b f(x) dx &= \sum_{i=1}^N \left[\frac{x_i - x_{i-1}}{2} \right] (f(x_{i-1}) + f(x_i)) + o(1) \\
&= \frac{1}{2} \sum_{i=1}^N [x_i - x_{i-1}] (f(x_{i-1}) + f(x_i)) + o(1) \\
&= \frac{1}{2} \sum_{i=1}^N h (f(x_{i-1}) + f(x_i)) + o(1) \\
&= \frac{h}{2} \left[2 \sum_{i=1}^{N-1} f(x_i) + f(x_0) + f(x_N) \right] + o(1).
\end{aligned}$$

If we use a geometric grid $x_{i+1} = s x_i = s^{i+1} x_0$ for $i = 0, \dots, N$, and with $s = \left(\frac{x_N}{x_0} \right)^{1/N}$ we get

$$\begin{aligned}
\int_a^b f(x) dx &= \sum_{i=1}^N \left[\frac{x_i - x_{i-1}}{2} \right] (f(x_{i-1}) + f(x_i)) + o(1) \\
&= \frac{1}{2} \sum_{i=1}^N [x_i - x_{i-1}] (f(x_{i-1}) + f(x_i)) + o(1) \\
&= \frac{x_0}{2} \sum_{i=1}^N [s^i - s^{i-1}] (f(x_{i-1}) + f(x_i)) + o(1) \\
&= \frac{x_0}{2} \left\{ \sum_{i=1}^{N-1} (s^{i+1} - s^{i-1}) f(x_i) + (s^1 - s^0) f(x_0) \right. \\
&\quad \left. + (s^N - s^{N-1}) f(x_N) \right\} + o(1).
\end{aligned}$$

Definition 3.4. The extended trapezoidal rule The integral of certain function in $[a, b]$ can be approximated by

$$\int_a^b f(x) dx = \sum_{i=1}^N \left[\frac{x_i - x_{i-1}}{2} \right] (f(x_{i-1}) + f(x_i)) + o(1) \quad (3.4)$$

If we use a uniform grid $x_{i+1} = x_i + h$ for $i = 0, \dots, N$, and $h = \frac{x_N - x_0}{N}$ we get

$$\int_a^b f(x) dx = \frac{h}{2} \left[2 \sum_{i=1}^{N-1} f(x_i) + f(x_0) + f(x_N) \right] + o(1) \quad (3.5)$$

And if we use a geometric grid $x_{i+1} = s x_i = s^{i+1} x_0$ for $i = 0, \dots, N$, and with $s = \left(\frac{x_N}{x_0} \right)^{1/N}$ we get

$$\int_a^b f(x) dx = \frac{x_0}{2} \left\{ \sum_{i=1}^{N-1} (s^{i+1} - s^{i-1}) f(x_i) + (s^1 - s^0) f(x_0) + (s^N - s^{N-1}) f(x_N) \right\} + o(1) \quad (3.6)$$

The interest of solving the PBE is to know the time evolution of the number distribution $n(L, t)$. Lets denote n_t the vector containing the values of this distribution at the set of the chosen coordinate points in the grid at instant of time t . Then, n_t is a $N + 1$ column vector, that is

$$n_k = \begin{bmatrix} n(L_0, k) \\ n(L_1, k) \\ \vdots \\ n(L_N, k) \end{bmatrix}_{N+1 \times 1}.$$

Proposition 3.5. Discretization scheme for the Population Balance Equation. The solution of the PBE can be computed numerically by an iterative procedure like

$$\begin{aligned} n_{k+1} &= n_k + \varepsilon [B_k^a - D_k^a + B_k^b - D_k^b] + o(1) \\ n_{k+1} &= n_k + \varepsilon \left[n_t \circ [(\bar{B} \circ \bar{S}) n_t] - n_t \circ (\beta S n_t) + (\Omega \circ \bar{S}) A n_t - A n_t \right] + o(1) \end{aligned} \quad (3.7)$$

having $A, B, S, \bar{B}, \bar{S}, \bar{S}$ and Ω as appropriate $N + 1 \times N + 1$ matrices:

- n_k is the number distribution vector at time instant k ,

- A is the matrix containing the values of the breakage kernel,
- B and \bar{B} are the matrices containing the values of the aggregation kernel,
- S , \bar{S} , and $\bar{\bar{S}}$ are matrices involving the geometric multiplicative step s , and
- Ω is a matrix containing the values of the fragment distribution function of daughter particles.

where \circ denotes the Hadamar product and where

- B_k^a is the vector representing the process of birth due to aggregation,
- D_k^a is the vector representing the process of death due to aggregation,
- B_k^b is the vector representing the process of birth due to breakage, and
- D_k^b is the vector representing the process of death due to breakage.

Proof. We may write the previous equation for each coordinate point L_i for $i = 0, \dots, N$ as

$$\begin{aligned}
\frac{\partial n(L_i, t)}{\partial t} &= B^a(L_i, t) - D^a(L_i, t) + B^b(L_i, t) - D^b(L_i, t) \\
&= \frac{L_i^2}{2} \int_0^{L_i} \frac{\beta\left(\left(L_i^3 - \lambda^3\right)^{1/3}, \lambda\right)}{\left(L_i^3 - \lambda^3\right)^{2/3}} n\left(\left(L_i^3 - \lambda^3\right)^{1/3}, t\right) n(\lambda, t) d\lambda \\
&\quad - n(L_i, t) \int_0^\infty \beta(L_i, \lambda) n(\lambda, t) d\lambda \\
&\quad + \int_{L_i}^\infty a(\lambda) b(L_i|\lambda) n(\lambda, t) d\lambda \\
&\quad - a(L_i) n(L_i, t).
\end{aligned}$$

We are going to deal with each term in the right side of the PBE as they involve integrals. We will use a numerical integration scheme on a large enough integration interval $[\underline{L}, \bar{L}]$. This will be performed by using the extended trapezoidal rule (assuming that the integrated function is regular in the whole integration domain). Then, we are going to find a numerical expression for each term.

Proposition 3.6. Numerical integration of the birth due to aggregation term. The term representing the process of birth due to aggregation $B^a(L_i, t)$ can be expressed as

$$B^a(L, t) = B_t^a$$

represented by an vector evaluated in each one of the points of the grid in the discretization scheme, and it can be estimated as

$$B_k^a = n_t \circ [(\bar{B} \circ \bar{S}) n_t] + o(1)$$

where

$$B_t^a = \begin{bmatrix} B^a(L_0, t) \\ \vdots \\ B^a(L_N, t) \end{bmatrix}_{N+1 \times 1}$$

$$\bar{\beta} = \begin{bmatrix} \bar{\beta}_0 \\ \vdots \\ \bar{\beta}_N \end{bmatrix}_{N+1 \times N+1}$$

$$\bar{S} = (\bar{S})_{i,j} = \begin{cases} 1 & \text{if } i = j \\ \frac{1}{2}L_0(S^1 - S^0), & \text{if } j = 0, \\ \frac{1}{2}L_0(S^{i-1} - S^{i-1}), & \text{if } j = i - 1, \\ \frac{1}{2}L_0(S^{j+1} - S^{j-1}), & \text{if } j = 1, \dots, i - 2 \\ 0 & \text{if } i < j \end{cases}$$

Proposition 3.7. Numerical integration of the death due to aggregation term. The term representing the process of death due to aggregation $D^a(L, t)$ can be expressed as

$$D^a(L, t) = D_t^a \tag{3.8}$$

represented by an vector evaluated in each one of the points of the grid in the discretization scheme, and it can be estimated as

$$D_t^a = n_t \circ (\beta S n_t) + o(1).$$

where

$$D_t^a = \begin{bmatrix} D^a(L_0, t) \\ \vdots \\ D^a(L_N, t) \end{bmatrix}_{N+1 \times 1}$$

$$\beta = \begin{bmatrix} \beta(L_0, L_0) & \cdots & \beta(L_0, L_N) \\ \vdots & & \vdots \\ \beta(L_N, L_0) & \cdots & \beta(L_N, L_N) \end{bmatrix}_{N+1 \times N+1}$$

and

$$S_{N+1 \times N+1} = (S)_{i,j} = \begin{cases} \text{if } i = j & \begin{cases} \frac{L_0}{2} s^1 - s^0, & \text{if } j = 0 \\ \frac{L_0}{2} s^N - s^{N-1}, & \text{if } j = N \\ \frac{L_0}{2} s^{j+1} - s^{j-1}, & \text{if } j = 1, \dots, N-1 \end{cases} \\ \text{else} & 0 \end{cases}$$

Proposition 3.8. Numerical integration of the birth due to breakage term. The term representing the process of birth due to breakage $B^b(L, t)$ can be expressed as

$$B^b(L, t) = B_t^b$$

represented by an vector evaluated in each one of the points of the grid in the discretization scheme, and it can be estimated as

$$B_t^b = \left(b \circ \bar{S} \right) a n_t + o(1).$$

where

$$B_t^b = \begin{bmatrix} B^b(L_0, t) \\ \vdots \\ B^b(L_N, t) \end{bmatrix} \quad a = \begin{bmatrix} a(L_0) & 0 & \cdots & 0 \\ 0 & a(L_1) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a(L_N) \end{bmatrix}_{N+1 \times N+1}$$

$$b = \begin{bmatrix} b(L_0 | L_0) & \cdots & b(L_0 | L_N) \\ \vdots & & \vdots \\ b(L_N | L_0) & \cdots & b(L_N | L_N) \end{bmatrix}_{N+1 \times N+1}$$

where

$$b(L_i | L_j) = \begin{cases} \frac{3L_i^2}{L_j^3}, & \text{if } i \leq j; \\ 0, & \text{if } i > j. \end{cases}$$

and

$$\bar{S}_{N+1 \times N+1} = \left(\bar{S} \right)_{i,j} = \begin{cases} \text{if } i \leq j & \begin{cases} \frac{L_0}{2} (s^{i+1} - s^i), & \text{if } j = i \\ \frac{L_0}{2} (s^{j+1} - s^{j-1}), & \text{if } j = i+1, \dots, N-1 \\ \frac{L_0}{2} (s^N - s^{N-1}), & \text{if } j = N \end{cases} \\ \text{if } i > j & 0 \end{cases}$$

Proposition 3.9. *Numerical integration of the death due to breakage term.* The term representing the process of birth due to breakage $D^b(L, t)$ can be expressed as

$$D^b(L, t) = D_t^b$$

represented by an vector evaluated in each one of the points of the grid in the discretization scheme, and it can be estimated as

$$D_t^b = an_t + o(1).$$

where

$$D_t^b = \begin{bmatrix} D^b(L_0, t) \\ \vdots \\ D^b(L_N, t) \end{bmatrix}.$$

and a is the same as in proposition 3.9.

Proof of proposition 3.6. From the definition 3.2 we have the expression of the term for each coordinate point L_i as

$$B^a(L_i; t) = \frac{L_i^2}{2} \int_0^{L_i} \frac{\beta\left((L_i^3 - \lambda^3)^{1/3}, \lambda\right)}{(L_i^3 - \lambda^3)^{2/3}} n\left((L_i^3 - \lambda^3)^{1/3}, t\right) n(\lambda, t) d\lambda.$$

We notice that there is a singularity in a neighborhood of L_i . In order to calculate numerically the value of this integral, we can split this integral into two parts and use the extended trapezoidal rule in the geometrical grid

$$\begin{aligned} B^a(L_i, t) &= \frac{L_i^2}{2} \int_0^{L_i} \frac{\beta\left((L_i^3 - \lambda^3)^{1/3}, \lambda\right)}{(L_i^3 - \lambda^3)^{2/3}} n\left((L_i^3 - \lambda^3)^{1/3}, t\right) n(\lambda, t) d\lambda = \\ &\frac{L_i^2}{2} \int_0^{L_{i-1}} \frac{\beta\left((L_i^3 - \lambda^3)^{1/3}, \lambda\right)}{(L_i^3 - \lambda^3)^{2/3}} n\left((L_i^3 - \lambda^3)^{1/3}, t\right) n(\lambda, t) d\lambda + \\ &\frac{L_i^2}{2} \int_{L_{i-1}}^{L_i} \frac{\beta\left((L_i^3 - \lambda^3)^{1/3}, \lambda\right)}{(L_i^3 - \lambda^3)^{2/3}} n\left((L_i^3 - \lambda^3)^{1/3}, t\right) n(\lambda, t) d\lambda. \end{aligned}$$

We use the extended trapezoidal rule for the numeric integration of the first

term so that the first integral is evaluated using the approximation

$$\begin{aligned}
& \frac{L_i^2}{2} \int_0^{L_{i-1}} \frac{\beta\left((L_i^3 - \lambda^3)^{1/3}, \lambda\right)}{(L_i^3 - \lambda^3)^{2/3}} n\left((L_i^3 - \lambda^3)^{1/3}, t\right) n(\lambda, t) d\lambda = \\
& \frac{1}{2} L_0 \left\{ \sum_{j=1}^{i-2} (s^{j+1} - s^{j-1}) \frac{L_i^2}{2} \frac{\beta\left((L_i^3 - L_j^3)^{1/3}, L_j\right)}{(L_i^3 - L_j^3)^{2/3}} n\left((L_i^3 - L_j^3)^{1/3}, t\right) n(L_j, t) \right. \\
& + (s^1 - s^0) \frac{L_i^2}{2} \frac{\beta\left((L_i^3 - L_0^3)^{1/3}, L_0\right)}{(L_i^3 - L_0^3)^{2/3}} n\left((L_i^3 - L_0^3)^{1/3}, t\right) n(L_0, t) \\
& + (s^{i-1} - s^{i-2}) \frac{L_i^2}{2} \frac{\beta\left((L_i^3 - L_{i-1}^3)^{1/3}, L_{i-1}\right)}{(L_i^3 - L_{i-1}^3)^{2/3}} n\left((L_i^3 - L_{i-1}^3)^{1/3}, t\right) n(L_{i-1}, t) \left. \right\} \\
& + o(1).
\end{aligned}$$

Proposition 3.10. Improper integral The improper integral in 3.9 can be numerically approximated like

$$\int_{L_{i-1}}^{L_i} \frac{1}{(L_i^3 - \lambda^3)^{2/3}} d\lambda = \frac{1}{3} \int_x^1 \frac{du}{u^{2/3} (L_i^3 - u)^{2/3}} \quad (3.9)$$

which can be calculated using an incomplete beta function $Beta(1/3, 1/3)$, where $x = \left(\frac{L_{i-1}}{L_i}\right)^3$.

Proof of the proposition 3.10. In order to compute the value of the integral

$$\int_{L_{i-1}}^{L_i} \frac{d\lambda}{(L_i^3 - \lambda^3)^{2/3}} \quad (3.10)$$

If we propose the transformation $u = \lambda^3$, then $du = 3\lambda^2 d\lambda$, $u \xrightarrow{\lambda \rightarrow L_{i-1}} L_{i-1}^3$, and $u \xrightarrow{\lambda \rightarrow L_i} L_i^3$, then we can write the integral as

$$\int_{L_{i-1}}^{L_i} \frac{d\lambda}{(L_i^3 - \lambda^3)^{2/3}} = \int_{L_{i-1}^3}^{L_i^3} \frac{du}{3u^{2/3} (L_i^3 - u)^{2/3}} \quad (3.11)$$

□

□

Proof. Proof of the proposition 3.7 From the definition 3.2 we have the expression of the term for each coordinate point L_i as

$$D^a(L, t) = n(L, t) \int_0^\infty \beta(L, \lambda) n(\lambda, t) d\lambda$$

$$D^a(L_i, t) = n(L_i, t) \int_0^\infty \beta(L_i, \lambda) n(\lambda, t) d\lambda$$

Again, applying the extended trapezoidal approximation rule for this integral in $\lambda = L_j$ for $j = 0, 1, \dots, N$, we have

$$D_t^a = \begin{bmatrix} D^a(L_0, t) \\ \vdots \\ D^a(L_N, t) \end{bmatrix}_{N+1 \times 1} \quad \beta = \begin{bmatrix} \beta(L_0, L_0) & \cdots & \beta(L_0, L_N) \\ \vdots & & \vdots \\ \beta(L_N, L_0) & \cdots & \beta(L_N, L_N) \end{bmatrix}_{N+1 \times N+1}$$

and

$$S_{N+1 \times N+1} = (S)_{i,j} = \begin{cases} \text{if } i = j \\ \quad \begin{cases} \frac{L_0}{2} s^1 - s^0, & \text{if } j = 0 \\ \frac{L_0}{2} s^N - s^{N-1}, & \text{if } j = N \\ \frac{L_0}{2} s^{j+1} - s^{j-1}, & \text{if } j = 1, \dots, N-1 \end{cases} \\ \text{else} \\ 0 \end{cases}$$

then, we can express the later equation like

$$D_t^a = n_t \circ (\beta S n_t) + o(1).$$

□

Proof. Proof of proposition 3.8 From the definition 3.2 we have the expression of the term for each coordinate point L_i as

$$B^b(L, t) = \int_L^\infty a(\lambda) b(L | \lambda) n(\lambda, t) d\lambda.$$

Continuing with the discretization scheme, we work now on the breakage birth term. We begin with

$$B^b(L_i, t) = \int_{L_i}^\infty a(\lambda) b(L_i | \lambda) n(\lambda, t) d\lambda.$$

and, applying the extended trapezoidal rule for solve this integral numerically, we get

$$B_t^b = \begin{bmatrix} B^b(L_0, t) \\ \vdots \\ B^b(L_N, t) \end{bmatrix} \quad a = \begin{bmatrix} a(L_0) & 0 & \cdots & 0 \\ 0 & a(L_1) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a(L_N) \end{bmatrix}_{N+1 \times N+1}$$

$$b = \begin{bmatrix} b(L_0 | L_0) & \cdots & b(L_0 | L_N) \\ \vdots & & \vdots \\ b(L_N | L_0) & \cdots & b(L_N | L_N) \end{bmatrix}_{N+1 \times N+1}$$

where

$$b(L_i | L_j) = \begin{cases} \frac{3L_i^2}{L_j^3}, & \text{if } i \leq j; \\ 0, & \text{if } i > j. \end{cases}$$

and

$$\bar{\bar{S}}_{N+1 \times N+1} = \left(\bar{\bar{S}} \right)_{i,j} = \begin{cases} \text{if } i \leq j \\ \begin{cases} \frac{L_0}{2} (s^{i+1} - s^i), & \text{if } j = i \\ \frac{L_0}{2} (s^{j+1} - s^{j-1}), & \text{if } j = i + 1, \dots, N - 1 \\ \frac{L_0}{2} (s^N - s^{N-1}), & \text{if } j = N \end{cases} \\ \text{if } i > j \\ 0 \end{cases}$$

then, we can write

$$B_t^b = \left(b \circ \bar{\bar{S}} \right) a n_t + o(1).$$

□

Proof. Proof of the proposition 3.9 From the definition 3.2 we have the expression of the term for each coordinate point L_i as

$$D^b(L, t) = a(L, t) n(L_i, t)$$

Using the same discretization scheme, we can express the term for death due to breakage like

$$D^b(L_i, t) = a(L_i, t) n(L_i, t).$$

for $i = 0, 1, \dots, N$.

Using matricial notation, let's denote

$$D_t^b = \begin{bmatrix} D^b(L_0, t) \\ \vdots \\ D^b(L_N, t) \end{bmatrix}.$$

Finally, we can express the last equation like

$$D_t^b = an_t + o(1).$$

□

Now, in order to find an approximation for

$$\frac{\partial n(L, t)}{\partial t}$$

we can write an approximation for the PBE as

$$\frac{\partial n_t}{\partial t} = B_t^a + D_t^a + B_t^b - D_t^b + o(1)$$

$$\frac{\partial n_t}{\partial t} = n_t \circ [(\bar{\beta} \circ \bar{S}) n_t] - n_t \circ (\beta S n_t) + (b \circ \bar{S}) an_t - an_t + o(1).$$

Now, we are going to use the forward one step Euler's method from t_k to $t_{k+1} = t_k + \varepsilon$ for handling the numerical derivative

$$\frac{\partial n_{t_k}}{\partial t} \approx \frac{n_{t_{k+1}} - n_{t_k}}{\varepsilon} + o(1)$$

hence, we can write the state equation like

$$n_{t_{k+1}} = n_{t_k} + \varepsilon \left[\frac{\partial n_{t_k}}{\partial t} \right] + o(1)$$

$$n_{t_{k+1}} = n_{t_k} + \varepsilon [B_{t_k}^a + D_{t_k}^a + B_{t_k}^b - D_{t_k}^b] + o(1)$$

$$n_{t_{k+1}} = n_{t_k} + \varepsilon [n_t \circ [(\bar{\beta} \circ \bar{S}) n_t] - n_t \circ (\beta S n_t) + (b \circ \bar{S}) an_t - an_t] + o(1).$$

□

3.1.3 Simulations using the Discretized PBE

We are going to evaluate the discretized PBE presented before, using the results from [Sil+10] in which they find the analytical solution for a PBE in the case of the number density distribution for the volume of the particles. We use this result for comparing the performance of the analytical and numerical solution.

The PBE is the conservation equation for the mean number density distribution function of particles, $n(v, t)$, whose dimensions depend on the particle properties, v , considered as distribution variables [Silva_2011].

The theoretical distribution

In [Sil+10], they studied the analytic solution given by Patil and Andrews (1998) [PA98] for a special case where the total number of particles is constant. Also, they use the solution given by McCoy and Madras (2003) for a more general case, where the number of particles is not constant, but using a different initial condition. The distributions are originally given using the particle's volume as distribution variable. We are going to use the same solutions but using particle's size as distribution variable. The initial equation is

$$\begin{aligned} \frac{\partial n'(v; t)}{\partial t} = & \frac{1}{2} \int_0^v \beta'(v - \epsilon, \epsilon) n'(v - \epsilon; t) n'(\epsilon; t) d\epsilon \\ & - n'(v; t) \int_0^\infty \beta'(v, \epsilon) n'(\epsilon; t) d\epsilon \\ & + 2 \int_v^\infty a'(\epsilon) b'(v/\epsilon) n'(\epsilon; t) d\epsilon \\ & - a'(v) n'(v; t). \end{aligned} \quad (3.12)$$

This PBE is subjected to the following initial conditions:

$$n'(v; 0) = \mu_0(0) \left(\frac{\mu_0(0)}{\mu_1(0)} \right) e^{-\left(\frac{\mu_0(0)}{\mu_1(0)}\right)v} \quad (3.13)$$

or

$$n'(v; 0) = \mu_0(0) \left[2 \frac{\mu_0(0)}{\mu_1(0)} \right]^2 v e^{-2\frac{\mu_0(0)}{\mu_1(0)}v} \quad (3.14)$$

where $\mu_0(0)$ and $\mu_1(0)$ are the initial zero and first-order moments. Due to mass conservation, μ_1 is constant for the considered problems. The aggregation and breakage kernels are:

- $\beta'(v, \epsilon) = C$, with C constant,
- $a'(v) = Sv$, with S constant,
- $b'(v | \epsilon) = \frac{2}{\epsilon}$.

For these choice of aggregation and breakage kernels, if the initial distribution is normalized ($\mu_0(0) = 1$), and $\mu_1 = 1$, and using $C = 1$ and $S = \frac{[\Phi(\infty)]^2}{2}$, and $\Phi(\infty) = 1$, the number of particles and the density distribution are constants. The analytical solution given by Patil and Andrews

(1998), assuming no variation in the total number of particles (and using the constants defined before), is

$$n'_a(v; t) = \sum_{i=1}^2 \frac{K_1(t) + p_i(t) K_2(t)}{L_2(t) + 4p_i(t)} e^{p_i(t)v} \text{ for } \forall t > 0; \quad (3.15)$$

where

$$\begin{aligned} K_1(t) &= 7 + t + e^{-t} \\ K_2(t) &= 2 - 2e^{-t} \\ L_2(t) &= 9 + t - e^{-t} \\ p_{1,2}(t) &= \frac{1}{4} (e^{-t} - t - 9) \pm \sqrt{d(t)} \\ d(t) &= t^2 + (10 - 2e^{-t})t + 25 - 26e^{-t} + e^{-2t}. \end{aligned} \quad (3.16)$$

McCoy and Madras (2003) treated the general case where the total number of particles is not constant. Thus, $\Phi(\infty)$ can assume arbitrary values, that represents systems with predominant breakage ($\Phi(\infty) > 1$) or aggregation ($\Phi(\infty) < 1$), they find the following solution when the first initial condition is used:

$$n'_a(v; t) = [\Phi(t)]^2 e^{-\Phi(t)v} \quad (3.17)$$

where

$$\Phi(t) = \Phi(\infty) \left[\frac{1 + \Phi(\infty) \tanh(\Phi(\infty)t/2)}{\Phi(\infty) + \tanh(\Phi(\infty)t/2)} \right]. \quad (3.18)$$

Using the particles size like distribution variable, we have

$$\begin{aligned} \frac{\partial n(L, t)}{\partial t} &= \frac{L^2}{2} \int_0^L \frac{\beta\left(\left(L^3 - \lambda^3\right)^{1/3}, \lambda\right)}{\left(L^3 - \lambda^3\right)^{2/3}} n\left(\left(L^3 - \lambda^3\right)^{1/3}, t\right) n(\lambda, t) d\lambda \\ &\quad - n(L, t) \int_0^\infty \beta(L, \lambda) n(\lambda, t) d\lambda \\ &\quad + 2 \int_L^\infty a(\lambda) b(L | \lambda) n(L, t) d\lambda \\ &\quad - a(L) n(L, t), \end{aligned} \quad (3.19)$$

with the following initial conditions

$$n(L, 0) = \mu_0(0) \left(\frac{\mu_0(0)}{\mu_1(0)} \right) (3L^2) e^{-\left(\frac{\mu_0(0)}{\mu_1(0)}\right)L^3} \quad (3.20)$$

or

$$n(L, 0) = \mu_0(0) \left[2 \frac{\mu_0(0)}{\mu_1(0)} \right]^2 (3L^5) e^{-2\frac{\mu_0(0)}{\mu_1(0)}L^3} \quad (3.21)$$

and with the following aggregation and breakage kernels

- $\beta(L, \lambda) = C$ with C constant,
- $a(L) = SL^3$ with S constant,
- $b(L | \lambda) = \frac{6L^2}{\lambda^3}$,

where $\mu_0(0)$ and $\mu_1(0)$ are the initial zero and first-order moments of $n'(v; 0)$. Then, the analytical solution given by Patil and Andrews (1998) is expressed by

$$n^a(L, t) = 3L^2 \left[\sum_{i=1}^2 \frac{K_1(t) + p_i(t) K_2(t)}{L_2(t) + 4p_i(t)} e^{p_i(t)L^3} \right] \text{ for } \forall t > 0; \quad (3.22)$$

and, the solution given by McCoy and Madras (2003) is expressed by

$$n^a(L, t) = 3L^2 [\Phi(t)]^2 e^{-\Phi(t)L^3}, \quad (3.23)$$

using the constants defined before.

Chapter 4

Parameter Estimation via Extended Kalman Filter and Least Squares

4.1 Introduction

In the previous chapter, the population balance equation (PBE), the integro-differential equation governing the time evolution of the number density distribution of a particles population has been solved numerically using a discretization scheme. However, the implementation of such a solution needs the knowledge of the kernels involved in the modelization and that of its parameters.

Recent investigation in colloidal particles systems are centered in the identification of those theoretical kernels for aggregation and breakage processes and also in the estimation of its parameters. The kernel modelization is done using hydrodynamic theoretical considerations. Nevertheless, the parameters involved have to be estimated. Some of those researches are based on data obtained from experimental controlled environments using a population of primary particles of size and volume known [Vli14]. The resulting data are expressed in terms of the volume distribution in time. From this data, experts are able to give an empirical estimation of the parameter vector.

The time number distribution function of particles can be used to estimate the parameter vector. Although this distribution is not directly measurable, the experimental data expressed as the volume distribution function contains information about the number distribution so it can be seen as a function of the number distribution. Now, the aggregation and breakage kernels being identified and having at hand some initial estimation of the parameters, we

propose a procedure to recover the number distribution from the volume distribution using the Extended Kalman Filtering.

For this, our procedure use the discretized scheme of the PBE as state equation. This equation involve the chosen kernels and the initial estimation of the parameters. We then consider the volume distribution as a measure equation involving the number distribution to be recovered. The Extended Kalman Filter algorithm predicts the number distribution at a discretized set of time instants.

Furthermore, in the same framework, we can also use the observation to estimate the unknown parameters of the PBE. Using the discretized scheme of the PBE as theoretical model, we propose a procedure to estimate the parameters iteratively. The procedure involves the misfit between the number distribution obtained from the discretized model and the predicted by the Extended Kalman Filter. We propose a Least Squares Estimator. This estimator is computed iteratively at each time. Notice that as the discretized model is not lineal, we have to use a linear approximation model. Moreover, the computation of the parameter vector estimation is done via ridge regression.

4.2 Estimation of the Number Distribution using Extended Kalman Filter

4.2.1 The Extended Kalman Filter

The Extended Kalman Filter (EKF) is a suboptimal solution for the filtering problem in the case where the state or the measure equations contains non-linear functions. It is included in the Analytic Approximation Filters because to be implemented it needs a linearization of the nonlinear functions of the state dynamic or measure models.[RAG04] [RS85]

Definition 4.1. Extended Kalman Filter The EKF is derived for nonlinear systems with additive noise, that is, for $k \in \mathbb{N}$

$$\begin{aligned}x_k &= f_{k-1}(x_{k-1}) + v_{k-1} \\z_k &= h_k(x_k) + w_k\end{aligned}\tag{4.1}$$

where f_{k-1} is the non-linear function linking x_k the state at time k with x_{k-1} the state at time $k-1$. h_k is the non-linear function linking z_k the measure at time k and the state at time k . v_{k-1} and w_k are random sequences, mutually independent, zero mean white Gaussian noise with covariances Q_{k-1} and R_k respectively.

The nonlinear functions in (4.1) are approximated by the first term in their Taylor series expansion [DV10]. The EKF is based on the assumption that the local linearization may be a sufficient description of nonlinearity. The posterior pdf $p(x_k | Z_k)$ is approximated by a Gaussian density and the following relationships are assumed to hold

$$\begin{aligned} p(x_{k-1} | Z_{k-1}) &= N(x_{k-1}; \hat{x}_{k-1|k-1}, P_{k-1|k-1}) \\ p(x_k | Z_{k-1}) &= N(x_k; \hat{x}_{k|k-1}, P_{k|k-1}) \\ p(x_k | Z_k) &= N(x_k; \hat{x}_{k|k}, P_{k|k}) \end{aligned} \quad (4.2)$$

where $p(x_{k-1} | Z_{k-1}) = N(x; m, P)$ is a Gaussian density with argument x , mean m and covariance P .

The mean and the covariance of the underlying Gaussian density are computed recursively as follows

$$\begin{aligned} \hat{x}_{k|k-1} &= f_{k-1}(\hat{x}_{k-1|k-1}) \\ P_{k|k-1} &= Q_{k-1} + \hat{F}_{k-1} P_{k-1|k-1} \hat{F}_{k-1}^T \\ \hat{x}_{k|k} &= \hat{x}_{k|k-1} + K_k (z_k - h_k(\hat{x}_{k|k-1})) \\ P_{k|k} &= P_{k|k-1} - K_k S_k K_k^T \end{aligned} \quad (4.3)$$

where

$$\begin{aligned} S_k &= \hat{H}_k P_{k|k-1} \hat{H}_k^T + R_k \\ K_k &= P_{k|k-1} \hat{H}_k^T S_k^{-1} \end{aligned} \quad (4.4)$$

and \hat{F}_{k-1} and \hat{H}_k are the local linearization of nonlinear functions f_{k-1} and h_k respectively. They are defined as Jacobians evaluated at $\hat{x}_{k-1|k-1}$ and $\hat{x}_{k|k-1}$ respectively, that is:

$$\begin{aligned} \hat{F}_{k-1} &= [\nabla_{x_{k-1}} f_{k-1}^T(x_{k-1})]^T |_{x_{k-1}=\hat{x}_{k-1|k-1}} \\ \hat{H}_k &= [\nabla_{x_k} h_k^T(x_k)]^T |_{x_k=\hat{x}_{k|k-1}} \end{aligned} \quad (4.5)$$

where

$$\nabla_{x_k} = \left[\frac{\partial}{\partial x_k[1]} \quad \cdots \quad \frac{\partial}{\partial x_k[n_x]} \right]^T \quad (4.6)$$

with $x_k[i]$; $i = 1, \dots, n_x$, being the i -th component of vector x_k . An element of, say, \hat{H}_k is then given by:

$$\hat{H}_k[i, j] = \frac{\partial h_k[i]}{\partial x_k[j]} |_{x_k=\hat{x}_{k|k-1}} \quad (4.7)$$

where $h_k[i]$ denotes the i -th component of vector $h_k(x_k)$.

For the number distribution function, we are going to use the discretized scheme of the PBE as state equation as state equation. Thus, the state vector is going to be represented by n_k the vector containing the values of this distribution at the set of the chosen coordinate points in the grid at instant of time k . Then, n_k is a $N + 1$ column vector, that is

$$n_k = \begin{bmatrix} n(L_0, k) \\ n(L_1, k) \\ \vdots \\ n(L_N, k) \end{bmatrix}_{N+1 \times 1},$$

and, the state equation is

$$\begin{aligned} n_k = & n_{k-1} + \varepsilon \{ n_{k-1} \circ [(\bar{B} \circ \bar{S}) n_{k-1}] - n_{k-1} \circ (\beta S n_{k-1}) \\ & + (\Omega \circ \bar{S}) A n_{k-1} - A n_{k-1} \} + v_{k-1} \end{aligned} \quad (4.8)$$

where v_{k-1} is a vector residues. The measure equation can be obtained as

$$v_i = \frac{v_{i,t}}{V_T} = \int_{y_i}^{y_{i+1}} \frac{\varphi}{V_T} L^3 n(L, t) dL;$$

where

$$V_T = \int_0^\infty \varphi L^3 n(L, t) dL;$$

using the discretization scheme, and denoting by $\varphi(L) = \varphi L^3$, $\varphi = \frac{\pi}{6} (2R_0)^3$, and R_0 is the primary particle's ratio, we get

$$\begin{aligned} v_{i,t} = & \frac{L_0}{2} \frac{1}{V_T} \left\{ \sum_{j=1}^{N-1} (s^{j+1} - s^{j-1}) \varphi_i(L_j) n(L_j, t) + (s^1 - s^0) \varphi_i(L_0) n(L_0, t) \right. \\ & \left. + (s^N - s^{N-1}) \varphi_i(L_N) n(L_N, t) \right\} + o(1), \end{aligned}$$

where $\varphi_i(L) = \frac{\pi}{6} (2R_0)^3 L^3 \mathbf{1}_{[y_i, y_{i+1})}(L)$ for $i = 1, \dots, d$. If we denote

$$v_t = \begin{bmatrix} v_{1,t} \\ \vdots \\ v_{d,t} \end{bmatrix}_{d \times 1}, \quad \varphi' = \begin{bmatrix} \varphi_1(L_0) & \cdots & \varphi_1(L_0) \\ \vdots & & \vdots \\ \varphi_d(L_0) & \cdots & \varphi_d(L_0) \end{bmatrix}_{d \times N+1}$$

then we get

$$v_t = \varphi' S n_t,$$

and, the constant V_T included in the model is known.

4.3 Least Squares Estimators for PBE parameters

4.3.1 General Information

Let's denote r the number of experiment's replicates, k the observation time, where $0 \leq k \leq t_{max}$, $z_{k,r}$ the volume fraction measured for the k -th time and for the r -th sample. Then, the available data can be represented as

$$\begin{bmatrix} z_{1,1} & z_{1,2} & \cdots & z_{1,r} \\ z_{2,1} & z_{2,2} & \cdots & z_{2,r} \\ \vdots & \vdots & & \vdots \\ z_{k,1} & z_{k,2} & \cdots & z_{k,r} \\ \vdots & \vdots & & \vdots \\ z_{t_{max},1} & z_{t_{max},2} & \cdots & z_{t_{max},r} \end{bmatrix}. \quad (4.9)$$

We are interested in the number distribution, and we apply a **Kalman Filtering** for recover this distribution from the measured matrix described above

$$\begin{aligned} x_k &= f_\theta(x_{k-1}) + v_{k-1} \\ z_k &= h(x_k) + w_k \end{aligned}$$

where the parameter vector $\theta \in \mathbb{R}^p$, and the probabilistic conditions for the **Extended Kalman Filter**

From this, we have the recovered objective distribution, denoted by

$$\begin{bmatrix} \hat{x}_{1,1} & \hat{x}_{1,2} & \cdots & \hat{x}_{1,r} \\ \hat{x}_{2,1} & \hat{x}_{2,2} & \cdots & \hat{x}_{2,r} \\ \vdots & \vdots & & \vdots \\ \hat{x}_{k,1} & \hat{x}_{k,2} & \cdots & \hat{x}_{k,r} \\ \vdots & \vdots & & \vdots \\ \hat{x}_{t_{max},1} & \hat{x}_{t_{max},2} & \cdots & \hat{x}_{t_{max},r} \end{bmatrix}$$

where

- $\hat{x}_{k,r} \in \mathbb{R}^{N+1}$: the number distribution filtered from the r -th sample in the k -th time.
- **Assumptions:**
 - Neither the f_θ and h functions nor the θ vector of parameters change in time.

– The error's covariance matrices Q_{k-1} and R_k can change in time.

The estimation problem can be stated as follows: We denote

- $\hat{x}_{k,r}$: the filtered state from the real measured data in k -th time and r -th sample.
- $\tilde{x}_{k,r}$: simulated state by the discretized PBE model in k -th time and r -th sample.

we want to find an estimator for θ using the **least Squares Method**, by minimizing the following **objective function**

$$\begin{aligned}\Gamma_k &= \sum_{l=1}^r [\hat{x}_{k,l} - \tilde{x}_{k,l}]^T [\hat{x}_{k,l} - \tilde{x}_{k,l}] \\ \Gamma_k &= \sum_{l=1}^r [\hat{x}_{k,l} - \tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta)]^T [\hat{x}_{k,l} - \tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta)]\end{aligned}$$

for $k = 1, \dots, t_{max}$, where $\tilde{x}_{k,l} = \tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta) = f_\theta(\hat{x}_{k-1,l})$

$$\begin{aligned}\tilde{x}_{k,l} &= \hat{x}_{k-1,l} + \varepsilon \{ \hat{x}_{k-1,l} \circ (\bar{\beta} S \hat{x}_{k-1,l}) - \hat{x}_{k-1,l} \circ (\beta S \hat{x}_{k-1,l}) \\ &\quad - ba S \hat{x}_{k-1,l} + a \hat{x}_{k-1,l} \}\end{aligned}$$

where the symbol \circ denote de Hadamar matricial product, and we have the measure equation which links the measure with the aimed state,

$$z_k = h(\hat{x}_k) = \varphi S \hat{x}_k.$$

In this case, we consider the following aggregation frequency and efficacy kernel and rupture frequency kernel, like

$$\begin{aligned}\beta(L, \lambda) &= \alpha_m \exp \left[-x \left[1 - \left(\frac{\max \left[\left(\frac{L}{2} \right)^{Df}, \left(\frac{\lambda}{2} \right)^{Df} \right]}{\min \left[\left(\frac{L}{2} \right)^{Df}, \left(\frac{\lambda}{2} \right)^{Df} \right]} \right) \right] \right] \times \\ &\quad \frac{G}{6} \frac{(L + \lambda)^3}{\left[\left(\frac{L}{2} \right)^{Df} \left(\frac{\lambda}{2} \right)^{Df} \right]^y}\end{aligned}$$

$$a(L) = \alpha_b G^{b_0} \left(\frac{L}{2} \right)^c,$$

this kernels are considered fixed and not time variant, and the vector of unknown parameters is $\theta = [\alpha_m, x, y, \alpha_b, b_0, c]^T \in \mathbb{R}^6$.

The matrices involved in the above equations are, after considering the fixed kernels, expressed like

$$\bar{\beta} = \begin{bmatrix} \bar{\beta}(L_0, L_0) & \cdots & \bar{\beta}(L_0, L_N) \\ \vdots & & \vdots \\ \bar{\beta}(L_N, L_0) & \cdots & \bar{\beta}(L_N, L_N) \end{bmatrix}_{N+1 \times N+1}$$

where

$$\bar{\beta}(L_i, L_j) = \begin{cases} \alpha_m e_1^{i,j}(x) m_1^{i,j}(y) l_1^{i,j}, & \text{if } i > j; \\ \alpha_m e_2^{i,j}(x) m_2^{i,j}(y) l_2^{i,j}, & \text{if } i = j; \\ 0, & \text{if } i < j. \end{cases}$$

where the functions are

$$e_1^{i,j}(x) = \exp \left[-x \left[1 - \frac{\left(\max \left[\left(\frac{(L_i^3 - L_j^3)^{1/3}}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]}{\left(\min \left[\left(\frac{(L_i^3 - L_j^3)^{1/3}}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]} \right) \right]^2 \right],$$

$$m_1^{i,j}(y) = \frac{G}{6} \frac{\left((L_i^3 - L_j^3)^{1/3} + L_j \right)^3}{\left[\left(\frac{(L_i^3 - L_j^3)^{1/3}}{2} \right)^{Df} \left(\frac{L_j}{2} \right)^{Df} \right]^y},$$

$$l_1^{i,j} = \frac{L_i^2}{2(L_i^3 - L_j^3)^{2/3}},$$

$$e_2^{i,j}(x) = \exp \left[-x \left[1 - \frac{\left(\max \left[\left(\frac{\left(L_i^3 - \left(\frac{L_{i-1} + L_i}{2} \right)^3 \right)^{1/3}}{2} \right)^{Df}, \left(\frac{L_{i-1} + L_i}{2} \right)^{Df} \right]}{\left(\min \left[\left(\frac{\left(L_i^3 - \left(\frac{L_{i-1} + L_i}{2} \right)^3 \right)^{1/3}}{2} \right)^{Df}, \left(\frac{L_{i-1} + L_i}{2} \right)^{Df} \right]} \right) \right]^2 \right],$$

$$m_2^{i,j}(y) = \frac{G}{6} \frac{\left(\left(L_i^3 - \left(\frac{L_{i-1} + L_i}{2} \right)^3 \right)^{1/3} + \frac{L_{i-1} + L_i}{2} \right)^3}{\left[\left(\left(L_i^3 - \left(\frac{L_{i-1} + L_i}{2} \right)^3 \right)^{1/3} \right)^{Df} \left(\frac{L_{i-1} + L_i}{2} \right)^{Df} \right]^y},$$

$$l_2^{i,j} = \frac{L_i^2 (L_i - L_{i-1})^{1/3}}{2(3L_i^2)^{2/3}}.$$

If we define the following matrices

$$\bar{e}(x)_{N+1 \times N+1} = (\bar{e}(x))_{i,j} = \begin{cases} e_1^{i,j}(x), & \text{if } i > j; \\ e_2^{i,j}(x), & \text{if } i = j; \\ 0, & \text{if } i < j, \end{cases}$$

$$\bar{m}(x)_{N+1 \times N+1} = (\bar{m}(x))_{i,j} = \begin{cases} m_1^{i,j}(y), & \text{if } i > j; \\ m_2^{i,j}(y), & \text{if } i = j; \\ 0, & \text{if } i < j, \end{cases}$$

$$\bar{D}_{N+1 \times N+1} = (\bar{D})_{i,j} = \begin{cases} l_1^{i,j}, & \text{if } i > j; \\ l_2^{i,j}, & \text{if } i = j; \\ 0, & \text{if } i < j. \end{cases}$$

we can rewrite the $\bar{\beta}$ matrix above in terms of this matrices like

$$\bar{\beta} = \alpha_m (\bar{e}(x) \circ \bar{m}(y) \circ \bar{D}).$$

For the β and a matrices, we can write a similar developement. This is

$$\beta_{N+1 \times N+1} = (\beta)_{i,j} = \beta(L_i, L_j)$$

where

$$\beta(L_i, L_j) = \alpha_m e_3^{i,j}(x) m_3^{i,j}(y),$$

and

$$e_3^{i,j}(x) = \exp \left[-x \left[1 - \left(\frac{\max \left[\left(\frac{L_i}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]}{\min \left[\left(\frac{L_i}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]} \right) \right] \right]$$

$$m_3^{i,j}(y) = \frac{G}{6} \frac{(L + \lambda)^3}{\left[\left(\frac{L}{2} \right)^{Df} \left(\frac{\lambda}{2} \right)^{Df} \right]^y}.$$

If we call

$$e(x)_{N+1 \times N+1} = (e(x))_{i,j} = e_3^{i,j}(x),$$

$$m(y)_{N+1 \times N+1} = (m(y))_{i,j} = m_3^{i,j}(y)$$

we have

$$\beta = \alpha_m (e(x) \circ m(y)).$$

In the same way, we have

$$a_{N+1 \times N+1} = (a)_{i,j} = \begin{cases} a(L_i), & \text{if } i = j; \\ 0, & \text{if } i \neq j \end{cases}$$

where

$$a(L_i) = \alpha_b G^{b_0} T^i(c),$$

and

$$T^i(c) = \left(\frac{L_i}{2}\right)^c.$$

If we call

$$T(c)_{N+1 \times N+1} = (T(c))_{i,j} = \begin{cases} T^i(c), & \text{if } i = j; \\ 0, & \text{if } i \neq j \end{cases}$$

the we can rewrite the a matrix like

$$a = \alpha_b G^{b_0} T(c).$$

The rest of the involved matrices are

$$b_{N+1 \times N+1} = (b)_{i,j} = b(L_i | L_j),$$

where

$$b(L_i | L_j) = \begin{cases} 2Df\left(\frac{L_i^{Df-1}}{L_j^{Df}}\right), & \text{if } i \leq j; \\ 0, & \text{if } i > j \end{cases}$$

and

$$S_{N+1 \times N+1} = (S)_{i,j} = \begin{cases} s-1, & \text{if } i = j = 0; \\ \frac{L_0}{2} s^{i+1} - s^{i-1}, & \text{if } i = j = 1, \dots, N-1; \\ s^N - s^{N-1}, & \text{if } i = j = N+1; \\ 0, & \text{if } i \neq j. \end{cases}$$

Considering the expressions above we can write the state equation like

$$\begin{aligned} \tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta) = & \hat{x}_{k-1,l} + \varepsilon \{ \hat{x}_{k-1,l} \circ [\alpha_m(\bar{e}(x) \circ \bar{m}(y) \circ \bar{D}) S \hat{x}_{k-1,l}] \\ & - \hat{x}_{k-1,l} \circ [\alpha_m(e(x) \circ m(y)) S \hat{x}_{k-1,l}] \\ & + b(\alpha_b G^{b_0} T(c)) S \hat{x}_{k-1,l} - (\alpha_b G^{b_0} T(c)) \hat{x}_{k-1,l} \} \end{aligned}$$

The Least Squares Estimator is defined like the point in the parametrical space where the objective function Γ_k is minimal for each $k = 1, \dots, t_{max}$. In

this case, the function $\tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta)$ is not linear, we propose the following strategy for computing this minimum. We start with a linearization for approximate the function $\tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta)$ like

$$\tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta) \approx C_{k,l} + W_{k,l}^T [\theta - \hat{\theta}_{k-1}]$$

where $\hat{\theta}_{k-1}$ is a previous approximation, and

$$C_{k,l} = \tilde{x}_{k,l}(\hat{x}_{k-1,l}; \hat{\theta}_{k-1})$$

$$W_{k,l} = \begin{bmatrix} \frac{\partial \tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial \tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta)}{\partial \theta_p} \end{bmatrix}^T \hat{\theta}_{k-1}$$

We can express the last equations like

$$[\hat{x}_{k,l} - \tilde{x}_{k,l}] \approx \left[\hat{x}_{k,l} - \left(C_{k,l} + W_{k,l}^T [\theta - \hat{\theta}_{k-1}] \right) \right]$$

$$\approx \left[\hat{x}_{k,l} - C_{k,l} + W_{k,l}^T \hat{\theta}_{k-1} - W_{k,l}^T \theta \right]$$

if we call

$$Z_{k,l} = \hat{x}_{k,l} - C_{k,l} + W_{k,l}^T \hat{\theta}_{k-1}$$

then we can write

$$[\hat{x}_{k,l} - \tilde{x}_{k,l}] \approx [Z_{k,l} - W_{k,l}^T \theta].$$

We can now use an approximation objective function like

$$\Gamma_k^* = \sum_{l=1}^r [Z_{k,l} - W_{k,l}^T \theta]^T [Z_{k,l} - W_{k,l}^T \theta]$$

and to search for the $\hat{\theta}$ that minimize this function, for each $k = 1, \dots, t_{max}$.

We can rewrite this problem as an linear least squares estimation like

$$\begin{bmatrix} Z_{k,1} \\ \vdots \\ Z_{k,r} \end{bmatrix} = \begin{bmatrix} W_{k,1}^T \\ \vdots \\ W_{k,r}^T \end{bmatrix} \theta + \epsilon$$

We can rewrite the last equation like

$$\bar{Z}_k = \bar{W}_k \theta + \epsilon$$

with the usual solution

$$\hat{\theta}_k = (\bar{W}_k^T \bar{W}_k)^{-1} \bar{W}_k^T \bar{Z}_k.$$

The expression for the $W_{k,l}$ considering the structure of the $\tilde{x}_{k,l}(\hat{x}_{k-1,l}; \theta)$ is

$$W_{k,l} = \left[\frac{\partial \tilde{x}_{k,l}}{\partial \alpha_m}, \frac{\partial \tilde{x}_{k,l}}{\partial x}, \frac{\partial \tilde{x}_{k,l}}{\partial y}, \frac{\partial \tilde{x}_{k,l}}{\partial \alpha_b}, \frac{\partial \tilde{x}_{k,l}}{\partial b_0}, \frac{\partial \tilde{x}_{k,l}}{\partial c} \right]^T,$$

where

$$\begin{aligned} \frac{\partial \tilde{x}_{k,l}}{\partial \alpha_m} &= \varepsilon \{ \hat{x}_{k-1} \circ [(\bar{e}(x) \circ \bar{m}(y) \circ \bar{D}) S \hat{x}_{k-1}] - \hat{x}_{k-1} \circ [(e(x) \circ m(y)) S \hat{x}_{k-1}] \}, \\ \frac{\partial \tilde{x}_{k,l}}{\partial x} &= \varepsilon \{ \hat{x}_{k-1} \circ \left[\alpha_m \left(\frac{\partial \bar{e}(x)}{\partial x} \circ \bar{m}(y) \circ \bar{D} \right) S \hat{x}_{k-1} \right] - \hat{x}_{k-1} \circ \left[\alpha_m \left(\frac{\partial e(x)}{\partial x} \circ m(y) \right) S \hat{x}_{k-1} \right] \}, \\ \frac{\partial \tilde{x}_{k,l}}{\partial y} &= \varepsilon \{ \hat{x}_{k-1} \circ \left[\alpha_m \left(\bar{e}(x) \circ \frac{\partial \bar{m}(y)}{\partial y} \circ \bar{D} \right) S \hat{x}_{k-1} \right] - \hat{x}_{k-1} \circ \left[\alpha_m \left(e(x) \circ \frac{\partial m(y)}{\partial y} \right) S \hat{x}_{k-1} \right] \}, \\ \frac{\partial \tilde{x}_{k,l}}{\partial \alpha_b} &= \varepsilon \{ b (G^{b_0} T(c)) S \hat{x}_{k-1} - (G^{b_0} T(c)) \hat{x}_{k-1} \}, \\ \frac{\partial \tilde{x}_{k,l}}{\partial b_0} &= \varepsilon \{ G^{b_0} \log(G) b (\alpha_b T(c)) S \hat{x}_{k-1} - G^{b_0} \log(G) (\alpha_b T(c)) \hat{x}_{k-1} \}, \\ \frac{\partial \tilde{x}_{k,l}}{\partial c} &= \varepsilon \{ b \left(\alpha_b G^{b_0} \frac{\partial T(c)}{\partial c} \right) S \hat{x}_{k-1,l} - \left(\alpha_b G^{b_0} \frac{\partial T(c)}{\partial c} \right) \hat{x}_{k-1,l} \} \end{aligned}$$

where we have the following expressions for

$$\frac{\partial \bar{e}(x)}{\partial x} = \frac{\partial}{\partial x} (\bar{e}(x))_{i,j} = \begin{cases} \frac{\partial e_1^{i,j}(x)}{\partial x}, & \text{if } i > j; \\ \frac{\partial e_2^{i,j}(x)}{\partial x}, & \text{if } i = j; \\ 0, & \text{if } i < j \end{cases}$$

where

$$\begin{aligned} \frac{\partial e_1^{i,j}(x)}{\partial x} &= - \left[1 - \frac{\max \left[\left(\frac{(L_i^3 - L_j^3)^{1/3}}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]}{\min \left[\left(\frac{(L_i^3 - L_j^3)^{1/3}}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]} \right]^2 \times \\ &\exp \left[-x \left[1 - \frac{\max \left[\left(\frac{(L_i^3 - L_j^3)^{1/3}}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]}{\min \left[\left(\frac{(L_i^3 - L_j^3)^{1/3}}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]} \right]^2 \right], \\ \frac{\partial e_2^{i,j}(x)}{\partial x} &= - \left[1 - \frac{\max \left[\left(\frac{(L_i^3 - (\frac{L_{i-1} + L_i}{2})^3)^{1/3}}{2} \right)^{Df}, \left(\frac{(\frac{L_{i-1} + L_i}{2})}{2} \right)^{Df} \right]}{\min \left[\left(\frac{(L_i^3 - (\frac{L_{i-1} + L_i}{2})^3)^{1/3}}{2} \right)^{Df}, \left(\frac{(\frac{L_{i-1} + L_i}{2})}{2} \right)^{Df} \right]} \right]^2 \times \\ &\exp \left[-x \left[1 - \frac{\max \left[\left(\frac{(L_i^3 - (\frac{L_{i-1} + L_i}{2})^3)^{1/3}}{2} \right)^{Df}, \left(\frac{(\frac{L_{i-1} + L_i}{2})}{2} \right)^{Df} \right]}{\min \left[\left(\frac{(L_i^3 - (\frac{L_{i-1} + L_i}{2})^3)^{1/3}}{2} \right)^{Df}, \left(\frac{(\frac{L_{i-1} + L_i}{2})}{2} \right)^{Df} \right]} \right]^2 \right] \end{aligned}$$

$$\frac{\partial e(x)}{\partial x} = \frac{(e(x))_{i,j}}{\partial x} = \frac{\partial e_3^{i,j}(x)}{\partial x}$$

where

$$\frac{\partial e_3^{i,j}(x)}{\partial x} = - \left[1 - \frac{\max \left[\left(\frac{L_i}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]}{\min \left[\left(\frac{L_i}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]} \right] \times \exp -x \left[1 - \frac{\max \left[\left(\frac{L_i}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]}{\min \left[\left(\frac{L_i}{2} \right)^{Df}, \left(\frac{L_j}{2} \right)^{Df} \right]} \right],$$

$$\frac{\partial \bar{m}(y)}{\partial y} = \frac{\partial (\bar{m}(y))_{i,j}}{\partial y} = \begin{cases} \frac{\partial m_1^{i,j}(y)}{\partial y}, & \text{if } i > j; \\ \frac{\partial m_2^{i,j}(y)}{\partial y}, & \text{if } i = j; \\ 0, & \text{if } i < j \end{cases}$$

where

$$\frac{\partial m_1^{i,j}(y)}{\partial y} = - \frac{G}{6} \frac{\left((L_i^3 - L_j^3)^{1/3} + L_j \right)^3}{\left[\left(\frac{(L_i^3 - L_j^3)^{1/3}}{2} \right)^{Df} \left(\frac{L_j}{2} \right)^{Df} \right]^y} \times \log \left[\left(\frac{(L_i^3 - L_j^3)^{1/3}}{2} \right)^{Df} \left(\frac{L_j}{2} \right)^{Df} \right]$$

and

$$\frac{\partial m_2^{i,j}(y)}{\partial y} = - \frac{G}{6} \frac{\left(\left(L_i^3 - \left(\frac{L_{i-1} + L_i}{2} \right)^3 \right)^{1/3} + \left(\frac{L_{i-1} + L_i}{2} \right) \right)^3}{\left[\left(\frac{(L_i^3 - (\frac{L_{i-1} + L_i}{2})^3)^{1/3}}{2} \right)^{Df} \left(\frac{(L_{i-1} + L_i)}{2} \right)^{Df} \right]^y} \times \log \left[\left(\frac{(L_i^3 - (\frac{L_{i-1} + L_i}{2})^3)^{1/3}}{2} \right)^{Df} \left(\frac{(L_{i-1} + L_i)}{2} \right)^{Df} \right]$$

$$\frac{\partial m(y)}{\partial y} = \frac{\partial (m(y))_{i,j}}{\partial y} = \frac{\partial m_3^{i,j}(y)}{\partial y}$$

where

$$\frac{\partial m_3^{i,j}(y)}{\partial y} = -\frac{G}{6} \frac{(L_j + L_j)^3}{\left[\left(\frac{L_i}{2}\right)^{Df} \left(\frac{L_j}{2}\right)^{Df}\right]^y} \log \left[\left(\frac{L_i}{2}\right)^{Df} \left(\frac{L_j}{2}\right)^{Df} \right],$$

and

$$\frac{\partial T(c)}{\partial c} = \frac{\partial (T(c))_{i,j}}{\partial} = \begin{cases} \frac{\partial T^i(c)}{\partial c}, & \text{if } i = j; \\ 0, & \text{if } i \neq j \end{cases}$$

where

$$\frac{\partial T^i(c)}{\partial c} = \left(\frac{L_i}{2}\right)^c \log \left(\frac{L_i}{2}\right).$$

4.4 Results

4.4.1 Simulation using the particular case (Silva 2010)

In [Sil+10], they studied the analytic solution of the Population Balance Equation given by Patil and Andrews (1998) [PA98] for a special case of aggregation and breakage kernels where the total number of particles is constant. For these case they obtained an analytical solution. We used this cases in order to study the behavior of the Extended Kalman Filter applied using as state equation the discretized model. We compared the number density function theoretically obtained with the predicted state recovered using the Extended Kalman Filter.

We use the theoretical solution of the PBE in order to simulate the measure and state processes. In this, we simulate samples from the theoretical number distribution adding to this distribution an additive Gaussian noise. From this noised state, we obtain an volume distribution function also adding to this process a Gaussian noise. This represents the simulated measure.

Definition 4.2. Discretization scheme for the Population Balance Equation. The solution of the PBE can be computed numerically by an iterative procedure like

$$\begin{aligned} n_{k+1} &= n_k + \varepsilon [B_k^a - D_k^a + B_k^b - D_k^b] \\ n_{k+1} &= n_k + \varepsilon \left[n_t \circ [(\bar{B} \circ \bar{S}) n_t] - n_t \circ (\beta S n_t) + (\Omega \circ \bar{S}) A n_t - A n_t \right] \end{aligned} \quad (4.10)$$

having $A, B, S, \bar{B}, \bar{S}, \bar{\bar{S}}$ and Ω as appropriate $N+1 \times N+1$ matrices and n_k is the number distribution vector at time instant k . The Measure discretized

equation is also defined as

$$v_k = \varphi' S n_k,$$

where v_k is the volume fraction and, the constant V_T included is the known Total Mass/Volume.

The definition of the number density distribution studied in [Sco67] and [Sil+10] is

$$n'(v; 0) = \mu_0(0) \left(\frac{\mu_0(0)}{\mu_1(0)} \right) e^{-\left(\frac{\mu_0(0)}{\mu_1(0)}\right)v} \quad (4.11)$$

or

$$n'(v; 0) = \mu_0(0) \left[2 \frac{\mu_0(0)}{\mu_1(0)} \right]^2 v e^{-2 \frac{\mu_0(0)}{\mu_1(0)} v} \quad (4.12)$$

where $\mu_0(0)$ and $\mu_1(0)$ are the initial zero and first-order moments. Due to mass conservation, μ_1 is constant for the considered problems. The aggregation and breakage kernels are:

- $\beta'(v, \epsilon) = C$, with C constant,
- $a'(v) = Sv$, with S constant,
- $b'(v | \epsilon) = \frac{2}{\epsilon}$.

For these choice of aggregation and breakage kernels, if the initial distribution is normalized ($\mu_0(0) = 1$), and $\mu_1 = 1$, and using $C = 1$ and $S = \frac{[\Phi(\infty)]^2}{2}$, and $\Phi(\infty) = 1$, the number of particles and the density distribution are constants. The analytic solution given by (Patil and Andrews, 1998)([PA98]), assuming no variation in the total number of particles (and using the constants defined before), is

$$n'_a(v; t) = \sum_{i=1}^2 \frac{K_1(t) + p_i(t) K_2(t)}{L_2(t) + 4p_i(t)} e^{p_i(t)v} \text{ for } \forall t > 0; \quad (4.13)$$

where

$$\begin{aligned} K_1(t) &= 7 + t + e^{-t} \\ K_2(t) &= 2 - 2e^{-t} \\ L_2(t) &= 9 + t - e^{-t} \\ p_{1,2}(t) &= \frac{1}{4} (e^{-t} - t - 9) \pm \sqrt{d(t)} \\ d(t) &= t^2 + (10 - 2e^{-t})t + 25 - 26e^{-t} + e^{-2t}. \end{aligned} \quad (4.14)$$

The simulation processes followed is:

1. We use the theoretical distribution in order to obtain the number density function at each time instant k ,
2. We generate a zero mean Gaussian process at constant variance matrix. This is a $N + 1$ vector.
3. We add to this Gaussian noise to the theoretical state. We don't let the simulated vector behave atypically. It means that we don't use values negatives or highly deviated, because the nature of the process to be simulated.
4. We use the Measure discretized model in order to obtain the volume fraction distribution.
5. We add a Gaussian noise to this distribution in order to obtain the simulated measure.
6. We use the Extended Kalman Filter in order to recover the original noiseless state from the simulated measure. We use a zero vector like initial condition in order to simulate the no-information case.

In the Figure (4.1) we can see the theoretical distribution simulated using the discretized model in definition (4.2) and the simulated state with additive Gaussian noise. Also, we can see the recovered state evolution in time from the no-information initial condition. We can see also the absolute error in all the time evolution and the comparison between the least theoretical state and the least recovered state.

4.4.2 Simulation from the case of (Vlieghe, 2016) theoretical study

We use the discretized model in order to simulate an state number distribution from the kernels and parameter vector identified in (Vlieghe, 2016)[Vli14]. Then we get from this simulated state the simulate volume distribution using also the Measure discretized model (definition 4.2). We observe the stability of the simulation and compare to the experimental data.

In Figure (4.2) we observe the real experimental data. Then we see the simulated state using the model, the kernels and parameter vector identified by (Vlieghe, 2016)[Vli14]. Then we observe the simulated volume fraction distribution obtained from the simulated number density.

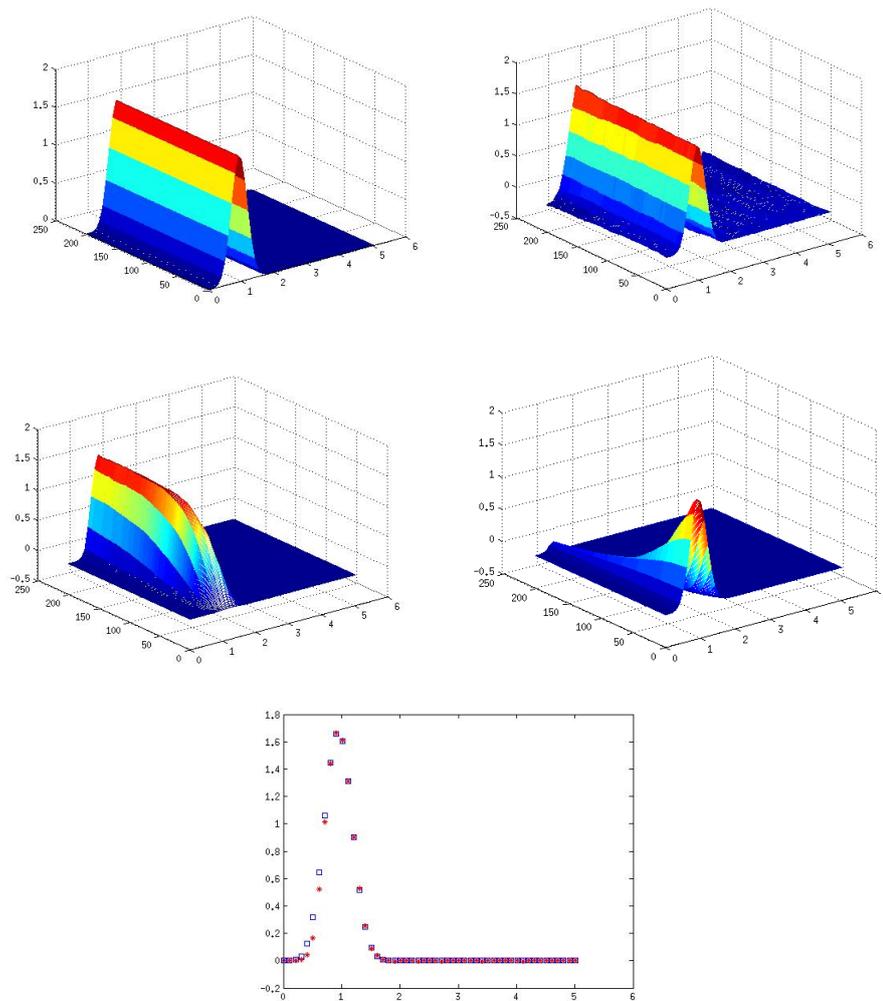


Figure 4.1: Behavior of the Extended Kalman Filter recovering the number density. a) Theoretical state. b) Simulated state with noise; c) State recovered by EKF. d) Absolute error; e) Error in the last state recovered

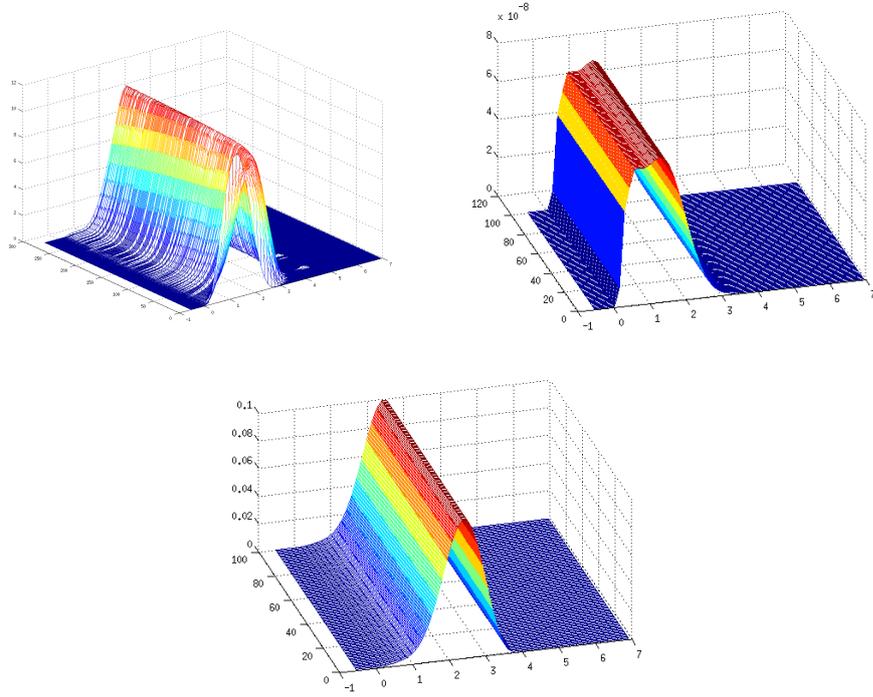


Figure 4.2: Behavior of the discretized model using the kernels and parameters identified in Vlieghe 2016. a) Experimental volume distribution function data. b) Simulated state using the kernels identified by Vlieghe 2016. c) Simulated measure from the state obtained using the discretized model.

In this case, we consider the following aggregation frequency and efficacy kernel and rupture frequency kernel, like

$$\beta(L, \lambda) = \alpha_m \exp \left[-x \left[1 - \left(\frac{\max \left[\left(\frac{L}{2} \right)^{Df}, \left(\frac{\lambda}{2} \right)^{Df} \right]}{\min \left[\left(\frac{L}{2} \right)^{Df}, \left(\frac{\lambda}{2} \right)^{Df} \right]} \right) \right] \right] \times$$

$$\frac{G}{6} \frac{(L + \lambda)^3}{\left[\left(\frac{L}{2} \right)^{Df} \left(\frac{\lambda}{2} \right)^{Df} \right]^y}$$

$$a(L) = \alpha_b G^{b_0} \left(\frac{L}{2} \right)^c,$$

this kernels are considered fixed and not time variant, and the vector of unknown parameters is $\theta = [\alpha_m, x, y, \alpha_b, b_0, c]^T \in \mathbb{R}^6$.

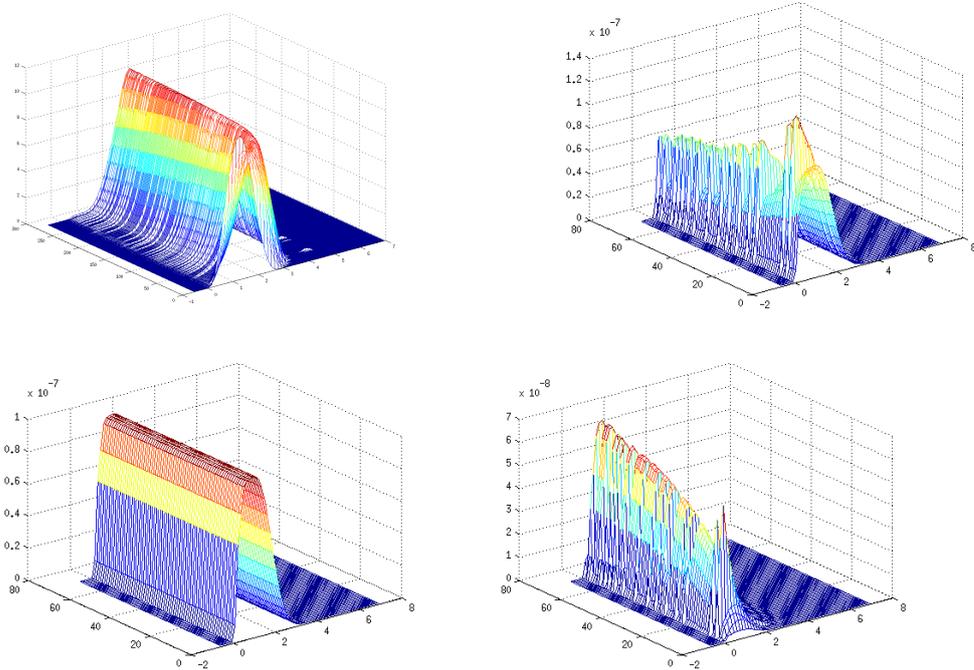


Figure 4.3: Behavior of the discretized model using the kernels and parameters identified in Vlieghe 2016. a) Experimental volume distribution function data. b) Empirical state identified by Vlieghe 2016. c) Recovered state using the kernels identified by Vlieghe 2016 and the EKF. d) Absolute error

4.4.3 Simulation from the case of Vlieghe 2016

We use the discretized model with the kernels and parameters set by Melody as state equation, in order to recover the unknown number density from the experimental measures obtained from Melody 2016.

The recovered state is then compared to the empirical number density obtained by Melody. We use the same procedure as in the section (4.4.2) and with the kernels shown in section (4.4.3) in order to compare the recovered number density with the empirical number density obtained from the real data.

In Figure (4.3) we can see the experimental volume fraction and the empirical number density obtained from them. Also, we can see the number density function recovered from the real volume fraction using the Extended Kalman Filter, and the absolute error between the empirical and the recovered number density.

4.4.4 Simulation using the EKF and LSE in a simulated case

We use the discretized model with the kernels used in Melody in order to test the behavior of the method in simulated data. We use the procedure described to obtain simulated measures using a vector of fixed parameters. Then we use different initialization vectors in order to test the sensibility of this choice.

4.4.5 Parameter estimation of parameters values of Guerin datasets

We use an initial vector of parameters, in order to apply the LSE to estimate the parameters values.

Chapter 5

Conclusion

Experimental empirical data was obtained in ([Vli14]). They analyzed the behavior of one micron (μm) bentonite primary particles in Jar test. They used several mixing speed in order to get different conditions for aggregation and breakage processes. The speed of mixing used were 30, 50, 70 and 90 revolutions per minute (rpm).

Bentonite is a kind of clay of type smectite, montmorillonite. The clay particles are the plates consisting of a stack of sheets separated by an inter-lacing space.

The experiments originating the empirical data implement the flocculation of Bentonite provided by CECA Chemicals. The mass concentration of the suspension is 30 mgL^{-1} . The suspension is diluted as much as allowing to take images.

The bentonite's mass needed for the experiments was putted in the water at least 24 hours before the experiments. This suspension had high intensity agitation by 45 *min*, before introduction into the reactor. Therefore, the suspension is constituted of primary agglomerates, compound of several elementary particles.

The initial population of this primary agglomerates shows that they have high variation about the Circle Equivalent Diameter (CED), with some values larger than $200 \mu m$. The most part of particles had size about some tens of micrometers and there are a small proportion of agglomerates with size larger than $200 \mu m$. The water used was demineralized water, the Ph is about 4.5 ± 0.1 and temperature between 20 and 25 °C. The experiments were performed in a Taylor-Couette Reactor.

In the experiments, they used a mechanism of aggregation by charges neutralization, using aluminum sulfate $Al_2(SO_4)_3$ as coagulant, which is used commonly in water treatment. It allows to produce flocs weak enough to interact with the hydrodynamic. The concentration of aluminum sulphate in

the reactor is $3.5 \times 10^{-5} \text{ molL}^{-1}$.

A continuation, the most important conclusions obtained from the descriptive data analysis are mentioned. The conclusions are divided in:

- Conclusions from the descriptive analysis of the data.
- Conclusions from the exploratory multivariate analysis of the data.

5.1 Conclusions of the descriptive analysis of the data.

The analysis was done using a natural classification for the variables in "size" variables and "shape" variables. Then, the description searched to interpret the data in that terms.

Initial population of Bentonite flocks.

Populations of Bentonite flocks under different hydrodynamics conditions.

Populations of Bentonite flocks under different hydrodynamics conditions thought time. We use the discretized model with the kernels and parameters seted by Melody as state equation, in order to recover the unknown number density from the experimental measures obtained from Melody 2016.

The recovered state is then compared to the empirical number density obtained by Melody. We use the same procedure as in the section (4.4.2) and with the kernels shown in section (4.4.3) in order to compare the recovered number density with the empirical number density obtained from the real data.

In Figure (4.3) we can see the experimental volume fraction and the empirical number density obtained from them. Also, we can see the number density function recovered from the real volume fraction using the Extended Kalman Filter, and the absolute error between the empirical and the recovered number density.

Bibliography

- [ADG14] Jean-Marc Azais, Yohann De Castro, and Fabrice Gamboa. “Spike detection from inaccurate samplings”. In: *Applied and Computational Harmonic Analysis* (2014).
- [Cab11] Hubert Cabana. *La coagulation, la floculation et l’agitation. Conception: usine de traitement des eaux potables*. 2011.
- [CDS98] S. S. Chen, D. L. Donoho, and M. A. Saunders. “Atomic decomposition by basis pursuit”. In: *SIAM J. Sci. Comput.* 20 (1998), pp. 33–61.
- [DG11] Yohann De Castro and Fabrice Gamboa. “Exact reconstruction using beurling minimal extrapolation”. In: *Journal of Mathematical Analysis and Applications* (2011).
- [DV10] Pierre Del Moral and Christelle Vergé. *Modèles et méthodes stochastiques. Une introduction avec applications*. Springer, 2010.
- [EC12] Carlos Fernandez-Granda Emmanuel Candès. “Towards a Mathematical Theory of Super-Resolution”. In: *Communications on Pure and Applied Mathematics* (2012).
- [FG17] Peter Filzmoser and Moritz Gschwandtner. *mvoutlier: Multivariate Outlier Detection Based on Robust Methods*. R package version 2.0.8. 2017. URL: <https://CRAN.R-project.org/package=mvoutlier>.
- [FMW08] Peter Filzmoser, Ricardo Maronna, and Mark Werner. “Outlier identification in high dimensions”. In: *Computational Statistics & Data Analysis* 52 (2008). Ed. by ELSEVIER, pp. 1694–1711. DOI: 10.1016/j.csda.2007.05.018.
- [FR02] Chris Fraley and Adrian E. Raftery. “Model-based Clustering, Discriminant Analysis and Density Estimation”. In: *Journal of the American Statistical Association* 97 (2002), pp. 611–631.

- [Fra+12] Chris Fraley et al. *mclust Version 4 for R: Normal Mixture Modeling for Model-Based Clustering, Classification, and Density Estimation*. 597. Technical Report. 2012.
- [Gor68] R.G. Gordon. “Error bounds in equilibrium statistical mechanics”. In: *Journal of Mathematical Physics* 9.5 (1968), pp. 655–663.
- [HS15] Wolfgang Karl Härdle and Léopold Simar. *Applied Multivariate Statistical Analysis*. Ed. by Springer. Fourth Edition. 2015. ISBN: 978-3-662-45171-7.
- [Lag07] Paulo L.C. Lage. “The quadrature method of moments for continuous thermodynamics”. In: *Computers and Chemical Engineering* 31 (2007), pp. 782–799.
- [Mar+03] Daniele Marchisio et al. “Quadrature Method of Moments for Population-Balance Equations”. In: *AIChE Journal* 49.5 (2003).
- [McG97] R. McGraw. “Description of aerosol dynamics by the quadrature method of moments”. In: *Aerosol Science and Technology* 27 (1997), pp. 255–265.
- [MF05] Daniele Marchisio and Rodney Fox. “Solution of population balance equations using the direct quadrature method of moments”. In: *Aerosol Science* 36 (2005), pp. 43–73.
- [Mor] *Morphologi G3 User Manual*. Malvern Instruments Ltd. 2008.
- [MVF03] Daniele Marchisio, Dennis Vigil, and Roney Fox. “Quadrature method of moments for aggregation–breakage processes”. In: *Journal of Colloid and Interface Science* 258 (2003), pp. 322–334.
- [Nak92] Shoichiro Nakamura. *Métodos Numéricos Aplicados con Software*. Spanish. Ed. by Pearson. Prentice-Hall Hispanoamericana, 1992.
- [PA98] D.P. Patil and J.R.G. Andrews. “An analytical solution to continuous population balance model describing floc coalescence and breakage. A special case.” In: *Chemical Engineering Science* 53.3 (1998), pp. 599–601.
- [RAG04] Branko Ristic, Sanjeev Arulampalam, and Neil Gordon. *Beyond the Kalman Filter. Particle Filters for tracking Applications*. DSTO, 2004.
- [Ram00] Doraiswami Ramkrishna. *Population Balances. Theory and Applications to Particulate Systems in Engineering*. San Diego, CA, USA: Academic Press, 2000. ISBN: 0-12-576970-9.

- [Ren02] Alvin Rencher. *Methods of Multivariate Analysis*. Ed. by Wiley Interscience. Second Edition. Wiley series in probability and statistics. John Wiley and Sons, 2002. ISBN: 0-471-41889-7.
- [RS85] Peter Ruymgaart and Tsu Soong. *Mathematics of Kalman-Bucy Filtering*. Springer-Verlag, 1985.
- [Saa] Dieter Saal. “Flocculation and Flocculation Filtration”. In: *Common fundamentals and unit operations in thermal desalination systems*. Vol. II.
- [Saf] Steven Safferman. *Fundamentals of Coagulation and Flocculation*. PennWell.
- [Sal07] Jean-Louis Salager. *Granulometría. Teoría*. Version 2. Escuela de ingeniería química. Facultad de ingeniería. Universidad de Los Andes. Mérida. Venezuela: Universidad de Los Andes, 2007.
- [Sco67] William Scott. “Analytic Studies of Cloud Droplet Coalescence”. In: *Journal of the Atmospheric Sciences* 25 (1967), pp. 54–65.
- [Sil+10] L.F.L.R. Silva et al. “Comparison of the accuracy and performance of quadrature-based methods for population balance problems with simultaneous breakage and aggregation”. In: *Computers and Chemical Engineering* 34 (2010), pp. 286–297.
- [Soo+07a] M. Soos et al. “Population Balance modeling of aggregation and breakage in turbulent Taylor-Couette flow”. In: *Journal of Colloid and Interface Science* (2007), pp. 433–446.
- [Soo+07b] M. Soos et al. “Population balance modeling of aggregation and breakage in turbulent Taylor–Couette flow”. In: *Journal of Colloid and Interface Science* 307 (2007), pp. 433–446.
- [SS66] Karlin Samuel and William J. Studden. *Tchebycheff Systems: with applications in analysis and statistics*. Interscience, 1966.
- [TD06] Tridib Tripathy and Bhudeb Ranjan De. “Flocculation : A New Way to Treat the Waste Water”. In: *Journal of Physical Sciences* 10 (2006), pp. 97–127.
- [Van00] Marco Vanni. “Approximate Population Balance Equations for Aggregation-Breakage Processes”. In: *Journal of Colloid and Interface Science* 221 (2000), pp. 143–160.
- [Vli+11] Melody Vlieghe et al. “Modélisation des processus d’agrégation et de rupture de floes par les bilans de population”. In: *Récents Progrès en Génie des Procédés* 101 (2011), Axx1–Axx6.

- [Vli14] Melody Vlieghe. “Agrégation et rupture de flocs sous contraintes turbulentes”. PhD thesis. Université de Toulouse ; INP,UPS ; LGC (Laboratoire de Génie Chimique), 2014.
- [WMpt] Xue Z. Wang and Cai Y. Ma. “Morphological Population Balance Model in Principal Component Space”. In: *AIChE Journal* 55.9 (September 2009).
- [WVF05] Liguang Wang, R. Dennis Vigil, and Rodney O. Fox. “CFD simulation of shear-induced aggregation and breakage in turbulent Taylor–Couette flow”. In: *Journal of Colloid and Interface Science* 285 (2005), 167–178.
- [R C17] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria, 2017. URL: <https://www.R-project.org/>.

Résumé

Nous concentrons notre intérêt sur l'Équation du Bilan de la Population (PBE). Cette équation décrit l'évolution, au fil du temps, des systèmes de particules en fonction de sa fonction de densité en nombre (NDF) où des processus d'agrégation et de rupture sont impliqués.

Dans la première partie, nous avons étudié la formation de groupes de particules et l'importance relative des variables dans la formation de ces groupes en utilisant les données dans (Vlieghe 2014) et des techniques exploratoires comme l'analyse en composantes principales, le partitionnement de données et l'analyse discriminante. Nous avons utilisé ce schéma d'analyse pour la population initiale de particules ainsi que pour les populations résultantes sous différentes conditions hydrodynamiques.

La deuxième partie nous avons étudié l'utilisation de la PBE en fonction des moments standard de la NDF, et les méthodes en quadrature des moments (QMOM) et l'Extrapolation Minimale Généralisée (GME), afin de récupérer l'évolution, d'un ensemble fini de moments standard de la NDF. La méthode QMOM utilise une application de l'algorithme Produit-Différence et GME récupère une mesure discrète non-négative, étant donnée un ensemble fini de ses moments standard.

Dans la troisième partie, nous avons proposé un schéma de discrétisation afin de trouver une approximation numérique de la solution de la PBE. Nous avons utilisé trois cas où la solution analytique est connue (Silva et al. 2011) afin de comparer la solution théorique à l'approximation trouvée avec le schéma de discrétisation.

La dernière partie concerne l'estimation des paramètres impliqués dans la modélisation des processus d'agrégation et de rupture impliqués dans la PBE. Nous avons proposé une méthode pour estimer ces paramètres en utilisant l'approximation numérique trouvée, ainsi que le Extended Kalman Filter. La méthode estime interactivement les paramètres à chaque instant du temps, en utilisant un estimateur de Moindres Carrés non-linéaire.

Mots-clés : Modélisation Stochastique, Méthode de la quadrature des moments, Extrapolation Minimale Généralisée, Filtre Étendu de Kalman, Estimateur de Moindres Carrés non-linéaire.

Abstract

We center our interest in the Population Balance Equation (PBE). This equation describes the time evolution of systems of colloidal particles in terms of its number density function (NDF) where processes of aggregation and breakage are involved.

In the first part, we investigated the formation of groups of particles using the available variables and the relative importance of these variables in the formation of the groups. We use data in (Vlieghe 2014) and exploratory techniques like principal component analysis, cluster analysis and discriminant analysis. We used this scheme of analysis for the initial population of particles as well as in the resulting populations under different hydrodynamics conditions.

In the second part we studied the use of the PBE in terms of the moments of the NDF, and the Quadrature Method of Moments (QMOM) and the Generalized Minimal Extrapolation (GME), in order to recover the time evolution of a finite set of standard moments of the NDF. The QMOM methods uses an application of the Product-Difference algorithm and GME recovers a discrete non-negative measure given a finite set of its standard moments.

In the third part, we proposed an discretization scheme in order to find a numerical approximation to the solution of the PBE. We used three cases where the analytical solution is known (Silva et al. 2011) in order to compare the theoretical solution to the approximation found with the discretization scheme.

In the last part, we proposed a method for estimate the parameters involved in the modelization of aggregation and breakage processes in PBE. The method uses the numerical approximation found, as well as the Extended Kalman Filter. The method estimates iteratively the parameters at each time, using an non-linear Least Square Estimator.

Key words : Stochastic Modelization, Quadrature Method of Moments, Generalized Minimal Extrapolation, Extended Kalman Filter, non-linear Least-Squares Estimator.