



SAPIENZA  
UNIVERSITÀ DI ROMA

# Exploring Authorial Identity through Data Mining Techniques

Analyzing Artistic Traits for Authorship Attribution

Faculty of Mathematics, Physics and Natural Sciences  
Master degree in Applied Mathematics

**Stefano Magrini Alunno**

ID number 1851728

Advisor

Prof. Gabriella Anna Pupo

Academic Year 2023/2024

Thesis defended on 00/00/00  
in front of a Board of Examiners composed by:

Prof. uno (chairman)

Prof. due

Prof. tre

Prof. ...

Prof. ...

Prof. ...

Prof. ...

---

**Exploring Authorial Identity through Data Mining Techniques**  
Master thesis. Sapienza University of Rome

© 2024 Stefano Magrini Alunno. All rights reserved

This thesis has been typeset by L<sup>A</sup>T<sub>E</sub>X and the Sapthesis class.

Version: November 21, 2024

Author's email: [stefanomagrini99@gmail.com](mailto:stefanomagrini99@gmail.com)

*Mater artium necessitas*

# Contents

<b>Acronyms</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Literature Review</b>	<b>2</b>
2.1 n-gram language model . . . . .	2
2.1.1 Insights . . . . .	3
2.1.2 Implications . . . . .	4
2.2 Fuzzy clustering . . . . .	5
2.2.1 Insights . . . . .	6
2.2.2 Implications . . . . .	10
2.3 Discrete Fourier Transform (DFT) . . . . .	12
2.3.1 Insights . . . . .	12
2.3.2 Implications . . . . .	14
2.4 Application . . . . .	16
2.4.1 Comparing works, the idea of clustering . . . . .	16
2.4.2 Fuzzy Clustering as a Noise Filtering Method . . . . .	18
2.4.3 Analysis of images with DFT . . . . .	19
<b>3 Methodology</b>	<b>22</b>
3.1 Data Set . . . . .	24
3.1.1 Description . . . . .	24
3.2 Pre processing . . . . .	26
3.2.1 Grayscale reduction . . . . .	27
3.2.2 Removing squares . . . . .	28
3.3 Synthesis . . . . .	31
3.4 Comparison . . . . .	32
<b>4 Results</b>	<b>39</b>
4.1 pre processing . . . . .	39
4.2 synthesis . . . . .	39
4.3 analysis . . . . .	39
<b>5 Discussion</b>	<b>41</b>
<b>6 Conclusion</b>	<b>42</b>
<b>A KMeans, GMM, FCM compare figures</b>	<b>43</b>

<b>B FCM implementation with CUDA</b>	<b>46</b>
<b>Glossary</b>	<b>49</b>
<b>Bibliography</b>	<b>50</b>

# List of Figures

2.2	example of data for clustering . . . . .	11
2.3	Dati in 2D sintetici . . . . .	18
3.1	BW transformation . . . . .	23
3.2	Pixel grid detail . . . . .	26
3.5	Synthetic FFT test . . . . .	30
3.6	Illustration of synthesis process . . . . .	32
3.8	We want to calculate the sum of the values in the first row. The idea is to divide the vector into 2 regions, sum the components, and repeat over the new vector with half the size of the previous vector. . . . .	38
4.1	Come si osserva i grigi si sono molto avvicinati tra loro dopo Fast Fourier Transform (FFT), ma è possibile ribilanciare i colori confrontandoli con l'immagine originale. . . . .	40
A.1	example of KMeans clustering . . . . .	43
A.2	example of Gaussian Mixture Models (GMM) clustering . . . . .	44
A.3	example of Fuzzy C-Means Clustering (FCM) clustering . . . . .	45

# List of Tables

2.1 Comparing distance results with/without clustering . . . . .	18
3.1 Summary of dataset . . . . .	25

# List of Theorems

2.1 Theorem (Update formula 'M')	9
2.2 Theorem (DFT)	14

# Special contents

To Review . . . . .	2
To Review . . . . .	22
To Do . . . . .	25
To Do . . . . .	25
To Review . . . . .	26
To Review . . . . .	27
To Review . . . . .	39

# Chapter 1

## Introduction

# Chapter 2

## Literature Review

This chapter reviews key literature that forms the theoretical foundation of this research. Specifically, the  $n$ -gram language model, fuzzy clustering, and Fourier transforms will be discussed. Additionally, specific details of these topics will be examined to demonstrate their practical applications within this project.

To Review

Relevant definitions and properties needed to understand these applications will be introduced, alongside examples and comparisons to provide a more comprehensive and nuanced understanding of the theory.

### 2.1 n-gram language model

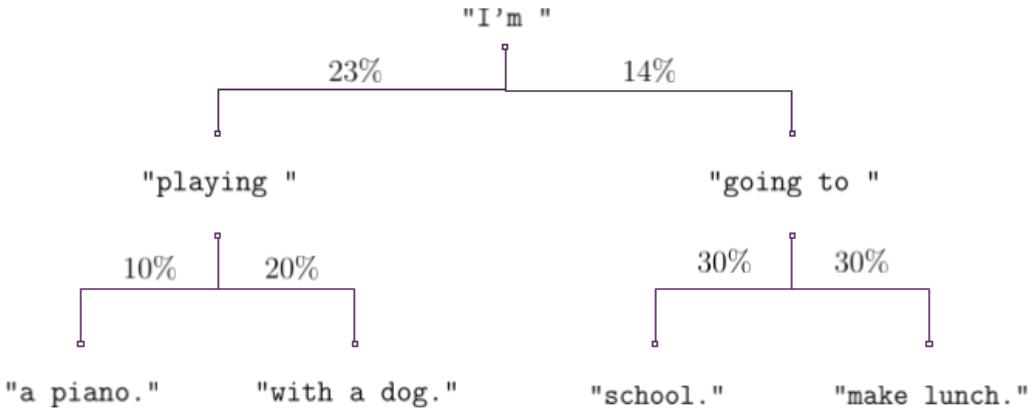
The  $n$ -gram language model analyzes a sequence of symbols by examining sets of frames, known as  $n$ -grams. This model is designed to construct a natural language model based on the assumption that each new symbol is statistically influenced by the preceding symbols. For example, the phrase 'I am playing' could be followed by 'a piano' or 'with a dog'. The  $n$ -gram model assumes that the continuation of this phrase depends solely on a finite window of preceding words or characters, and assigns 'a piano' a certain probability of being the most natural continuation. However, as language modeling techniques have advanced, the  $n$ -gram model has largely been replaced by more sophisticated models, such as transformers<sup>1</sup>.

The  $n$ -gram language model, first introduced by Shannon and described in Martin [5], was later applied in Basile et al. [1] for the purpose of author attribution. In this case, the goal was not to generate natural text but to identify the author of a given text. The research found that the best performance was achieved using 8-grams, and attribution was done by comparing each known author's work with the target text.

An additional development of this model extended its use to a completely different field: images. In the research thesis [4], this attribution method was employed to determine the authorship of a manuscript by treating the grams as square pixel tiles.

---

<sup>1</sup>For more information on transformers, see [https://en.wikipedia.org/wiki/Transformer\\_\(deep\\_learning\\_architecture\)](https://en.wikipedia.org/wiki/Transformer_(deep_learning_architecture))



### 2.1.1 Insights

As previously mentioned, the  $n$ -gram model predicts the next characters or words to be placed based on the preceding ones. For example, if the focus is on sequences of  $n$  printable characters, the model is trained on a dataset of text files, empirically deriving the probabilities  $\mathbb{P}[w_n|w_1, \dots, w_{n-1}]$ , where  $w_i$  represents the  $i$ -th character of the  $n$ -gram.

This probabilistic information can then be used to estimate the likelihood of subsequent characters appearing, with a certain margin of error. The model assumes that:

$$\mathbb{P}[w_{n+1}|w_1 \dots w_n] = \mathbb{P}[w_{n+1}|w_2 \dots w_n] \quad (2.1)$$

From this assumption we derive the following property:

**Proposition 2.1.1.** *Let  $w_1, \dots, w_{n-1}, w_n, \dots, w_{n+k}$  be a sequence of characters. We have that:*

$$\mathbb{P}[w_n, \dots, w_{n+k}|w_1, \dots, w_{n-1}] = \prod_{i=0}^k \mathbb{P}[w_{n+i}|w_{i+1} \dots w_{n+i-1}] \quad (2.2)$$

*Proof.* This property is proved by using Bayes' theorem and the assumption of the model eq. (2.1):

$$\begin{aligned} \mathbb{P}[w_n, \dots, w_{n+k}|w_1, \dots, w_{n-1}] &= \prod_{i=0}^k \mathbb{P}[w_{n+i}|w_1 \dots w_{n+i-1}] \\ &= \prod_{i=0}^k \mathbb{P}[w_{n+i}|w_{i+1} \dots w_{n+i-1}] \end{aligned}$$

□

The model is built on the assumption that considering  $n$ -grams larger than  $n$  is unnecessary. However, this limitation can introduce significant weaknesses in natural language processing, as it ignores information outside the  $n$ -character window. For instance, in a mystery novel, understanding the entire plot and its intricate details is crucial, something the model might overlook. Despite this, the  $n$ -gram model can

still be quite effective in simpler contexts like everyday conversation. For example, after the input 'Hi! How are ', the model might predict 'you?' as a natural continuation.

As mentioned earlier, this model was repurposed for a different goal by Basile et al. [1]. While the  $n$ -gram model might struggle to generate fully coherent natural language without encountering semantic inconsistencies or syntactic errors, it remains useful for capturing an author's writing style. This is achieved by training the model exclusively on texts by that author. The approach is to extract all  $n$ -grams from the target work and compare their distribution to that of known works by other authors. The proposed formula for comparing work  $A$  with work  $B$  is as follows:

$$d(A, B) = \frac{1}{|D_A| + |D_B|} \sum_x \left( \frac{f_A(x) - f_B(x)}{f_A(x) + f_B(x)} \right)^2 \quad (2.3)$$

where  $f_A(x)$  represents the frequency of the  $n$ -gram  $x$  in the work  $A$  and  $|D_A|$  is the variety of observed  $n$ -grams.

As highlighted in [4], two relevant properties of this method of comparing distributions can be mentioned:

- Each addendum of the summation takes on a value between 0 and 1.
- The external factor not only normalises the result so that it remains between 0 and 1, but also calls up the Jaccard index<sup>2</sup>, improving the arithmetic mean of the summation:

$$\frac{1}{|D_A| + |D_B|} = (1 + J_{D_A, D_B})^{-1} \frac{1}{|D_A \cup D_B|}$$

### 2.1.2 Implications

We conclude this section by discussing the implications of this theory. As seen in eq. (2.1), no distinction is made between one  $n$ -gram and another. This is because  $n$ -grams of length 3, or slightly longer, are typically used, which are insufficient to cover a full word. However, when using 7-grams, it becomes possible to consider the relationships between synonyms and distances between  $n$ -grams. For instance, we would expect some level of correlation between 'paper' and 'article' as they are synonyms.

In written texts, this rarely impacts the outcome, as there are relatively few synonyms for each term within the full vocabulary, and especially few long chains of synonyms. For example, it is unlikely that many synonyms for 'paper' begin with the string 'p' or 'pe'.

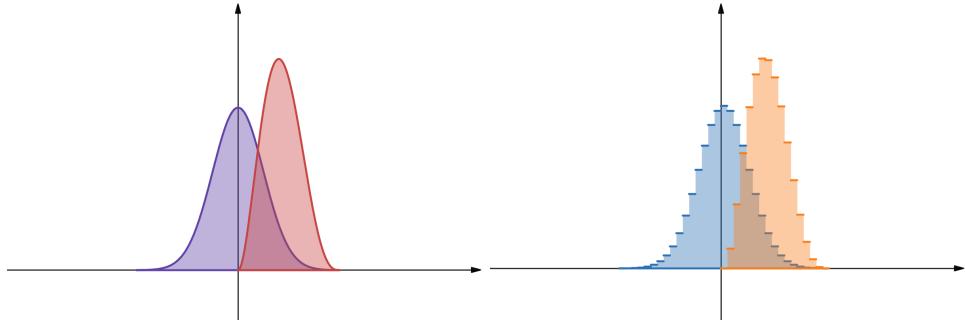
As noted in [4], the research deliberately set appropriate image *resolution* and *post-terization* to avoid discussions about the similarity between tiles. By controlling these parameters, the focus remained on the attribution process rather than getting lost in complex considerations of tile 'synonyms' and their potential concatenations.

We further examine the properties of the distance defined in eq. (2.3), particularly focusing on the problem of sparsity in  $n$ -grams.

---

<sup>2</sup>For more about Jaccard index, see [https://en.wikipedia.org/wiki/Jaccard\\_index](https://en.wikipedia.org/wiki/Jaccard_index)

Let  $N$  represent the samples drawn from two absolutely continuous distributions  $\mathcal{A}$  and  $\mathcal{B}$  defined on  $\mathbb{R}$ . We partition  $\mathbb{R}$  into intervals (boxes) of amplitude  $1 / C$ , thus transforming the original distributions into their discretized versions,  $\mathcal{A}(C)$  and  $\mathcal{B}(C)$ , along with their respective samples.



Now, let's consider eq. (2.3) and analyze the behavior of the estimated distance as the number of data points and the size of the boxes vary. By fixing the number of data points and decreasing the size of the boxes, the estimated distance gradually approaches 1. This phenomenon occurs because there exists a critical dimension for  $1 / C$ , such that with a sufficiently high probability, each sample becomes isolated within its own box. Consequently, this drives  $J_{D_A, D_B}$  toward 0, making each term in the summation equal to 1.

In [4], the continuous nature of the  $n$ -tiles was addressed by discretizing the continuous space in which they were defined, fixing a resolution and applying posterization. However, this approach results in a significant loss of information. Therefore, this paper will explore a method to achieve a cleaner and more dynamic discretization.

## 2.2 Fuzzy clustering

Fuzzy clustering, also known as soft k-means, is a data analysis technique that allows data points to be assigned to more than one cluster. Unlike traditional "hard" clustering, where each data point is assigned exclusively to a single cluster, fuzzy clustering employs fuzzy logic to determine the degree of membership for each data point within multiple clusters. Instead of a binary membership statement (1 or 0), fuzzy clustering provides a continuous measure of membership that ranges from 0 to 1. This is represented by a similarity function,  $\mu_x(C)$ , which quantifies the degree to which a data point  $x$  belongs to a cluster  $C$ .

Fuzzy clustering originated from the work of Dunn, who introduced this concept as an extension of the  $k$ -means clustering algorithm in 1973. Dunn emphasized the increased accuracy and robustness of fuzzy clustering compared to hard clustering, particularly in handling anomalous data (outliers) [2].

The algorithm proposed by [2] is known as the FCM and employs the Expectation-Maximisation (EM) technique to define the cluster membership function. This method is an iterative approach for estimating statistical parameters, such as cluster centroids. The EM process consists of two main steps:

1. Expectation (E): In this phase, a likelihood function is formulated. This function expresses the probability that a data point belongs to each cluster, based on its distance or similarity to the centroids.
2. Maximisation (M): In the maximization phase, an optimization procedure is used to identify the model parameters that maximize the likelihood, specifically the centroids of the clusters.

In summary, fuzzy clustering provides greater flexibility compared to hard clustering, allowing for a more accurate and detailed representation of the data, particularly in the presence of heterogeneous data or outliers.

### 2.2.1 Insights

To fully understand the FCM algorithm, we must first define the fundamental concepts on which it is based.

**Data points:** These are the observations or entities in our dataset, each characterized by a set of attributes or features. Data points can be viewed as points in a multidimensional space, where each dimension corresponds to a specific attribute.

**Clusters:** Once we have defined the data points, we can organize them into groups called clusters. A cluster is a collection of data points that share similar characteristics. The goal of the clustering algorithm is to group data points so that those within the same cluster are more similar to each other than to those in different clusters.

**Cluster Membership Function:** Each data point can be associated with one or more clusters through the cluster membership function, denoted as  $\mu_x$ , where  $x$  represents a specific data point and  $C$  represents a cluster. This function assigns each data point a degree of membership in each cluster, represented by a value between 0 and 1. For example, if  $\mu_x(C) = 1$ , it indicates that data point  $x$  belongs completely to cluster  $C$ . Conversely, if  $\mu_x(C) = 0.5$ , it suggests that  $x$  has some degree of uncertainty or ambiguity in its membership to cluster  $C$ .

**Notation.** Denote:

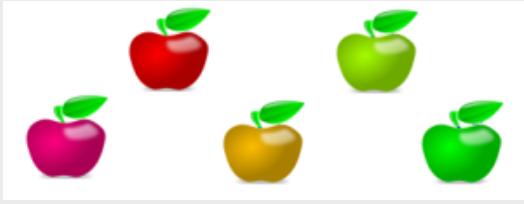
- the dataset with  
 $\mathcal{S} = (x_i)_{i=1:N}$  where  $x_i \in \mathbb{R}^K$
- the set of cluster centroids with  
 $\mathcal{C} = C_1, \dots, C_M$  where  $C_j \in \mathbb{R}^K$
- the weight matrix with  
 $U = (u_{ij})$  dove  $u_{ij} = \mu_{x_i}(C_j)$
- the quadratic Euclidean distance matrix with  
 $D^2 = (d_{ij}^2)$  where  $d_{ij} = \|x_i - C_j\|^2$

*Exempli Gratia.* To understand the meaning of data points and clusters, let us consider the following example:

Imagine we need to classify the apples transported by a lorry into two color categories: green and red. However, within this set, there might also be yellow apples. The  $k$ -means algorithm would treat a yellow apple as an anomaly compared to the red and green apples. In contrast, fuzzy clustering quantifies this anomaly by asserting that a yellow apple is a blend of red and green.

The colours of the apples can be represented in a one-dimensional space. In this space, we observe a distribution of data points showing green and red at the extremes, with yellow apples located in the center.

Clusters, in the context of fuzzy logic, are sets defined by centroids representing specific shades of colour, which are not necessarily present in the set of data points. For example, if there are red and green apples, there are actually two centroids representing specific shades of red and green, even though there are no apples with such colour shades.



*Remark.* The FCM algorithm recognizes that real-world data may contain errors or uncertainties. Instead of rigidly assigning each data point to a single cluster, it introduces a measure of uncertainty in the assignment. This approach results in assigning a fuzzy degree of membership to each point concerning each cluster, rather than relying on a binary classification.

**Definition 2.2.1** (membership measure). Given a point  $x$  and a cluster  $C$ , the membership measure  $\mu_x(C)$  is proportional to the inverse of the square of the Euclidean distance between the point  $x$  and the centroid of the cluster  $C$ , i.e.  $\frac{1}{\|x-C\|^2}$ . We assume that the sum of the membership measures of  $x$  to all clusters in the set  $\mathcal{C}$  is equal to 1:

$$\sum_{C \in \mathcal{C}} \mu_x(C) = 1 \quad \forall x$$

Consequently, the membership measure  $u_{ij}$  of point  $x_i$  to cluster  $C_j$  is computed as:

$$u_{ij} := \mu_{x_i}(C_j) = \frac{1}{\sum_k \frac{d_{ij}^2}{d_{ik}^2}}$$

*Remark.* The KMeans algorithm, which is based on Boolean logic, seeks to minimise the sum of the squared distances between each data point and the centroid of the cluster to which it is assigned. In other words, the goal of the KMeans algorithm is to minimise the sum of intra-cluster variances, where the variance of a cluster is defined as the sum of the squared distances between each point in the cluster and the cluster centroid.

Formally, the goal of the KMeans algorithm can be expressed as:

$$\operatorname{argmin}_C \left[ \sum_{C \in \mathcal{C}} \sum_{x \in C} \|x - C\|^2 \right]$$

In the case of fuzzy clustering, the membership of a point in a cluster is expressed by a fuzzy measure instead of a Boolean logic. The objective function of the KMeans method can be rewritten equivalently with membership measures:

$$\sum_{C \in \mathcal{C}} \sum_{x \in \mathcal{S}} \delta_x(C) \|x - C\|^2$$

where  $\delta_x(C)$  represents the membership measure of point  $x$  in cluster  $C$ . This objective function reflects the weighting of the squared distances by the membership measure of each point in the clusters.

**Definition 2.2.2** (Loss function and cluster weight). We define the fuzzy clustering loss function from the KMeans algorithm as follows:

$$L_{\mathcal{S}}(\mathcal{C}) := \sum_{C \in \mathcal{C}} p_{\mathcal{S}}(C) := \sum_{C \in \mathcal{C}} \sum_{x \in \mathcal{S}} \mu_x(C)^2 \|x - C\|^2 \quad (2.4)$$

where  $p_{\mathcal{S}}(C)$  is the sum of the squared distances between the points in the cluster  $C$  and its centroid. This loss function is coercive on the centroids, which means that it tends to infinity when the centroids recede to infinity, and has a finite lower bound when the centroids are bound to a compact set in the data space. For this reason we know that there is a global minimum although the uniqueness of local minima is not guaranteed.<sup>a</sup>

---

<sup>a</sup>for more about coercive functions, see  
<https://www1.mat.uniroma1.it/people/lanzara/OTT0708/Ottimizzazione0708.pdf>

*Remark.* The KMeans algorithm and fuzzy clustering via FCM share a similar approach in updating of cluster centroids. Both algorithms aim to find the optimal position of the centroids that minimizes the loss function based on the distances between the data points and the centroids themselves.

In the KMeans algorithm, the process of updating centroids occurs iteratively. Initially, centroids are randomly assigned to clusters. Data points are then assigned to the nearest clusters (based on Euclidean distance), after which the centroids are updated by averaging the points assigned to each cluster. This process is repeated until the centroids converge to stable positions.

Similarly, the FCM algorithm iteratively updates both the centroids of the clusters and the cluster membership measures for the points. In each iteration, the optimal centroids for a given configuration of membership measures are determined, and vice versa, until a stable configuration is reached that minimizes the loss function defined based on the membership measures and the distances between the data points and the centroids.

**Theorem 2.1 (Update formula 'M')**

Fixed the matrix  $U$ , the optimisation of eq. (2.4) finds its minimum in the centroids according to the following update formula:

$$C_j^{\text{new}} = \frac{\sum_{i=1}^N u_{ij}^2 x_i}{\sum_{i=1}^N u_{ij}^2} \quad \forall j \quad (2.5)$$

*Proof.* Fixed the matrix  $U$ , the loss function is differentiable, convex and coercive. Consequently, to find the global minimum of the function, it is sufficient to find the points where the gradient is null.

Calculating the partial derivatives of the gradient with respect to the centroids  $C_j$ , we obtain:

$$\frac{\partial}{\partial C_j} L_S(\mathcal{C}) = \frac{\partial}{\partial C_j} p_S(C_j) = \sum_{i=1}^N \frac{\partial}{\partial C_j} u_{ij}^2 \|x_i - C_j\|^2 = 2 \sum_{i=1}^N u_{ij}^2 (C_j - x_i)$$

The gradient is null when all partial derivatives are null, which implies:

$$\sum_{i=1}^N u_{ij}^2 C_j^{\text{new}} = \sum_{i=1}^N u_{ij}^2 x_i \quad \forall j$$

Since  $C_j^{\text{new}}$  does not depend on the index  $i$ , the thesis is proved.  $\square$

*Remark.* The calculation of the centroid update formula in fuzzy clustering is closely tied to the concept of parameter estimation through the maximum likelihood method. Specifically, it involves maximizing the joint density of the data and model parameters, conditional on the assignment of points to clusters. The density to be maximized is given by:

$$f_U(\mathcal{S}|\mathcal{C}) \propto \prod_{i,j} \exp \left[ -u_{ij}^2 \frac{\|x_i - C_j\|^2}{2} \right]$$

where  $f_U(\mathcal{S}|\mathcal{C})$  is the joint density of the points conditioned on the set of centroids where  $u_{ij}$  is the membership measure of point  $x_i$  in cluster  $C_j$ .

The density of the single sample  $x_i$  is:

$$f_U^i(x_i|\mathcal{C}) \propto \prod_j \exp \left[ -u_{ij}^2 \frac{\|x_i - C_j\|^2}{2} \right]$$

where  $f_U^i(x_i|\mathcal{C})$  is the density of the observed datum  $x_i$  conditioned on the centroids. The goal is to maximise the likelihood density in order to obtain accurate estimates of the model parameters, which in the context of fuzzy clustering are represented by the cluster centroids. For this reason, the algorithm is classified as EM because the phase E sets the densities for maximum likelihood when the matrix  $U$  is defined and the phase M solves the problem by finding the centroids that best explain the model.

The algorithm 1 computes the ideal centroids in fuzzy clustering and does not directly aim to minimise of eq. (2.4), but rather to iteratively updates the centroids and membership measures until a stable configuration is reached that best represents the input data.

---

**Algorithm 1** Fuzzy Clustering

---

INPUT

- $\mathcal{S}$ : set of data  $x_1, \dots, x_N$
  - $\mathcal{C}$ : centroids  $C_1, \dots, C_j$
- 

```

1: procedure FUZZYCLUSTERING( $\mathcal{S}, \mathcal{C}$ )
2:   while not reached the exit criterion do
3:      $D^2 \leftarrow (d_{ij}^2)$  with  $d_{ij} = \|x_i - C_j\|^2$ 
4:      $U^2 \leftarrow (u_{ij}^2)$  with  $u_{ij} = \left(\sum_k D_{ik}^2 / D_{ik}^2\right)^{-1}$             $\triangleright$  see theorem 2.1
5:     for  $j \leftarrow 1$  to  $M$  do
6:        $C_j^{\text{new}} \leftarrow \sum_{i=1}^N x_i U_{ij}^2 / \sum_{i=1}^N U_{ij}^2$ 
7:     end for
8:   end while
9: end procedure

```

---

## 2.2.2 Implications

This algorithm is particularly effective for clustering problems that involve noise. Two other well-known methods operate on different paradigms: KMeans and GMM.

**KMeans** is a fundamental clustering algorithm that partitions data points into disjoint groups, where the centroids represent the mean of each group. Its primary goal is to minimize the sum of variances within these groups, resulting in well-defined clusters.

**GMM** is a more sophisticated algorithm based on the assumption that the data points are generated by a specific statistical process:

- 1 Select with probability  $p_j$  the  $j$ -th cluster.
- 2 Generate a random data point from  $\mathcal{N}(C_j, \Sigma_j)$ .

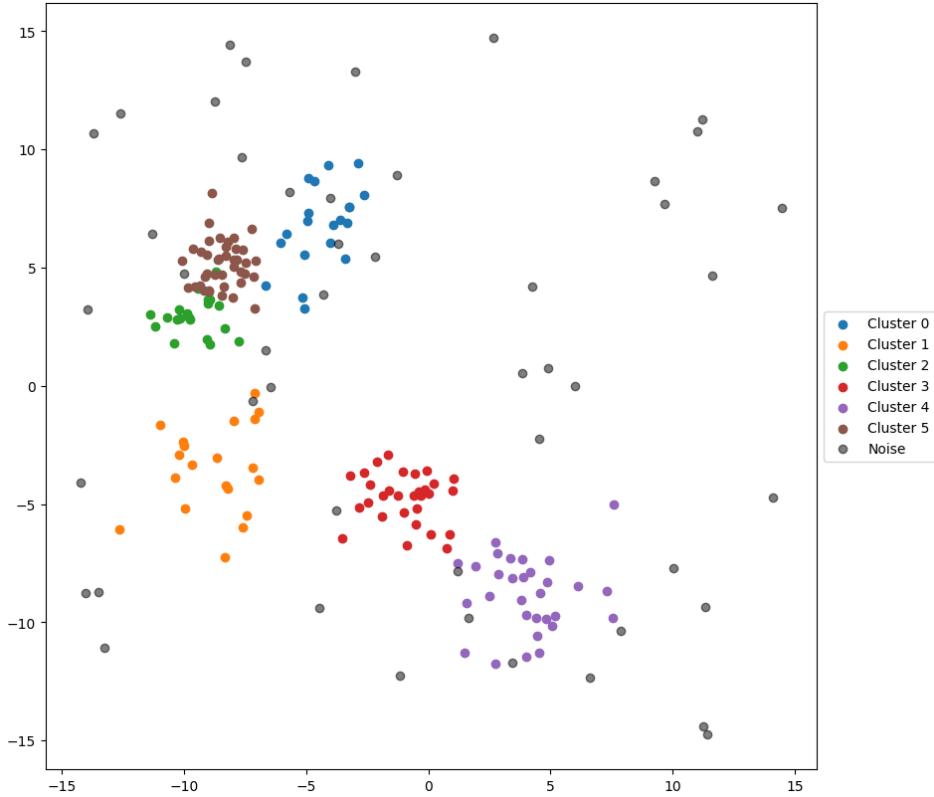
This method not only tries to find clusters and label data points, but also proposes normal distributions with averages  $\mathcal{C}$  and covariances  $\Sigma$ , assigning probabilities indicating the weight of each cluster. It is possible to label the data points using the Mahalanobis distance<sup>3</sup>.

Next, lets see how these methods behave in the presence of noisy data. Specifically, 6 random centroids were chosen in a two-dimensional space, from which 20, 20, 20, 30, 30 and 40 data points with different covariance matrices were generated,

---

<sup>3</sup>For more information on the Mahalanobis distance, see [https://en.wikipedia.org/wiki/Mahalanobis\\_distance](https://en.wikipedia.org/wiki/Mahalanobis_distance)

respectively. In addition, uniformly generated noisy data points were added, constituting 30% of the total (i.e. 48 noisy data points). A representation of this data can be seen in fig. 2.2. As can be seen, noise can create *outliers*, i.e. distant points that would be better ignored to achieve good clustering.



**Figure 2.2.** Data points are coloured according to their label, while noise is shown in grey.

Now let will try clustering with KMeans using the right number of centroids, i.e. 6. In fig. A.1 we can see that the noise has created an additional cluster and two of the original clusters have merged into one. Although the result is quite satisfactory, it is not entirely clean.

Next, let us apply clustering with GMM again using 6 centroids. In fig. A.2 it can be seen that the weight of the cluster generated by the noise is very small. Furthermore, the covariance matrices allow the position of the centroids to be calculated more precisely, providing information on the nature of the clusters and creating a hierarchical structure. This is because, in reality, these are not real clusters, but normal distributions.

Finally, let clustering be applied with FCM using 6 centroids. In fig. A.3 we observe that the centroids are less affected by noise, making it possible to identify which data are potentially noisy or shared between several clusters.

We have observed that in the presence of noise, the algorithm FCM can be very useful. However, the algorithm GMM, although computationally onerous, also provides very valuable information about the nature of the data. In this thesis, we

will use the FCM algorithm to reduce the amount of information to be processed and the KMeans algorithm for efficient initialisation of centroids.

## 2.3 DFT

The DFT is a mathematical transformation used to analyse the internal periodicities of a time sequence of data. Specifically, it decomposes a periodic time series as a sum of sine waves at different frequencies, giving information on the amplitude and phase of each frequency present in the signal. The more data available, the more detailed DFT will be.

This procedure, known as spectral analysis, makes it possible to study signals, waves, vibrations, sounds and images. In addition to these areas, DFT also has major scientific applications, such as the precise estimation of sunspot cycles, helping to predict and control phenomena such as geomagnetic storms<sup>4</sup>.

### 2.3.1 Insights

The DFT comes about as a discrete version of the Continue Fourier Transform (CFT), which represents a continuous function as a potentially infinite sum of sine waves. More rigorously, given a function  $f \in L^1(\mathbb{R})$  (i.e., a function integrable over the whole  $\mathbb{R}$ ), its Fourier transform is a functional operator, denoted by  $\mathcal{F}$ , which associates the function  $f$  with its frequency representation:

$$\mathcal{F}(f) : \omega \mapsto \int_{\mathbb{R}} e^{-i2\pi\omega x} f(x) dx$$

The **Fourier inversion theorem** states that, if both  $f$  and  $\mathcal{F}(f)$  belong to  $L^1(\mathbb{R})$  (i.e. they are both integrable), then for almost any  $x \in \mathbb{R}$  it is possible to recover  $f(x)$  via the following inverse relation:

$$f(x) = \int_{\mathbb{R}} \mathcal{F}(f)(\omega) e^{i2\pi x \omega} d\omega$$

In other words, the Fourier transform and its inverse allow switching back and forth between the time (or space) and frequency domains.

The discrete version of this transform, namely DFT, is used to analyse signals sampled at regular intervals. Thus, while CFT works on continuous signals, DFT applies to finite and discrete signals, making it suitable for digital signal processing. To better understand how DFT gives information about the amplitudes and phases of the different frequencies that make up a sequence, it is useful to introduce the Fourier coefficients. These coefficients make it possible to decompose a data sequence into sinusoids associated with different frequencies.

---

<sup>4</sup>For more about geomagnetic storms and them correlation with sunspot, see [https://en.wikipedia.org/wiki/Geomagnetic\\_storm](https://en.wikipedia.org/wiki/Geomagnetic_storm)

**Definition 2.3.1.** We define Fourier coefficients of the sequence  $(x_n)_{n=0}^{N-1}$  as those complex numbers  $(X_k)_{k=0}^{N-1}$  such that:

$$x_n = \frac{1}{N} \sum_{k=0}^{N-1} X_k e^{i2\pi \frac{k}{N} n} \quad (2.6)$$

The  $k$ -th coefficient gives information about the amplitude and phase of the frequency  $\frac{2\pi i}{N} k$ . We present some examples to better understand the meaning of the Fourier coefficients.

*Exempli Gratia.* Take as our first example a time series consisting of  $N$  complex numbers, described by the following coefficients:

$$X_0 = 0, X_1 = A, X_2 = 0 \dots X_{N-2} = 0, X_{N-1} = A$$

From these coefficients we obtain the time series:

$$x_n = \frac{1}{N} A \left( e^{i2\pi \frac{n}{N}} + e^{i2\pi \frac{n}{N}(N-1)} \right) = \frac{2A}{N} \cos \left( 2\pi \frac{n}{N} \right)$$

We note that the series is purely periodic, with amplitude  $\frac{2A}{N}$  and frequency  $\frac{2\pi}{N}$ .

Let us now consider a series of different coefficients, where all values are zero except:

$$X_p = A, X_{N-p} = A$$

The generated time series will be:

$$x_n = \frac{1}{N} A \left( e^{i2\pi \frac{n}{N}p} + e^{i2\pi \frac{n}{N}(N-p)} \right) = \frac{2A}{N} \cos \left( 2\pi \frac{n}{N}p \right)$$

As can be seen, the coefficients refer to the amplitude of a certain frequency. In addition, to obtain a real time series, symmetry on the coefficients is required.

At last, let us see how to derive the phase from a time series. Let us consider the following coefficients:

$$X_0 = 0, X_1 = A e^{i\theta}, X_2 = 0, \dots, X_{N-2} = 0, X_{N-1} = A e^{-i\theta}$$

The resulting time series will be:

$$x_n = \frac{2A}{N} \cos \left( \frac{2\pi}{N} n + \theta \right)$$

Here, the modulus of the coefficient describes the amplitude of the frequency, while its argument  $\theta$  describes its phase. Again, to ensure that the series is real, the opposite coefficient must be the conjunct of  $X_1$ .

We now ask ourselves whether it is possible to estimate the Fourier coefficients from a time series.

**Proposition 2.3.1.** *The matrix  $\frac{1}{\sqrt{N}} (e^{i2\pi \frac{k}{N} n})_{n,k=0}^{N-1}$  is unitary.*

*Proof.* We prove the thesis by studying the scalar product between two rows of the matrix, one of which is hermitian

$$\frac{1}{N} \sum_{k=0}^{N-1} e^{i2\pi \frac{k}{N} n} \cdot e^{-i2\pi \frac{k}{N} m} = \frac{1}{N} \sum_{k=0}^{N-1} e^{i2\pi \frac{k}{N} (n-m)}$$

If  $n = m$ , the result of the sum is 1. Else if  $n \neq m$ , we obtain:

$$\frac{1}{N} \sum_{k=0}^{N-1} e^{i2\pi \frac{k}{N} (n-m)} = \frac{1}{N} \frac{e^{i2\pi \frac{N}{N} (n-m)} - 1}{e^{i2\pi \frac{1}{N} (n-m)} - 1} = 0$$

This prove that the matrix is unitary.  $\square$

Proving that the matrix  $\frac{1}{\sqrt{N}} (e^{i2\pi \frac{k}{N} n})_{n,k=0}^{N-1}$  is unitary is extremely relevant to the computation of Fourier coefficients.

### Theorem 2.2 (DFT)

The Fourier coefficients of the time series  $(x_n)_{n=0}^{N-1}$  are given by the formula:

$$X_k = \sum_{n=0}^{N-1} x_n e^{-i2\pi \frac{n}{N} k} \quad (2.7)$$

*Proof.* Define  $U$  as matrix  $\frac{1}{\sqrt{N}} (e^{i2\pi \frac{k}{N} n})_{n,k=0}^{N-1}$

By definition of the discrete Fourier transform, we can write  $\vec{X} = \frac{1}{\sqrt{N}} U \vec{x}$ , where  $\vec{X}$  is the vector of Fourier coefficients and  $\vec{x}$  is the vector of the original time series.

As shown in proposition 2.3.1, the matrix  $U$  is unitary, so it admits an inverse, which is its transposed conjugate  $U^H$ .

Therefore, we can invert the transformation and obtain:

$$\vec{x} = \sqrt{N} U^H \vec{X}$$

This relation allows us to get the Fourier coefficients from the time series.  $\square$

### 2.3.2 Implications

The possibility of calculating Fourier coefficients directly by means of a simple matrix-vector product not only makes the analysis very informative, but also easily achievable.

---

**Algorithm 2** FFT Cooley-Tukey

---

```

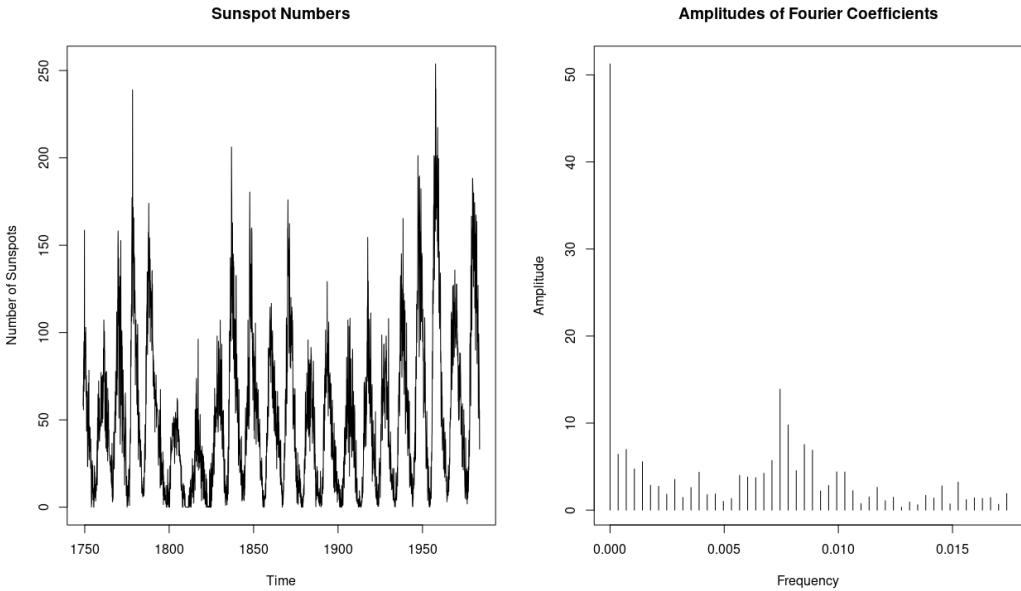
1: procedure FFT( $x$ )
2:   if  $N$  is 1 then return  $x$ 
3:   end if
4:   even  $\leftarrow$  FFT( $[x_0, x_2, \dots, x_{N-2}]$ )
5:   odd  $\leftarrow$  FFT( $[x_1, x_3, \dots, x_{N-1}]$ )
6:    $X$  is an array of size  $N$ 
7:   for  $k \in 0, \dots, \frac{N}{2} - 1$  do
8:      $t \leftarrow \exp(-2\pi ik/N) \cdot \text{odd}[k]$ 
9:      $X[k] \leftarrow \text{even}[k] + t$ 
10:     $X[k + N/2] \leftarrow \text{even}[k] - t$ 
11:   end for
12: end procedure

```

---

**FFT** The computational cost of the matrix-vector product for computing Fourier coefficients is  $O(N^2)$ , which is not excessive in itself. However, there is a specific algorithm for Fourier series that reduces this cost to  $O(N \log N)$ , comparable to a sorting algorithm. This algorithm, named FFT, was designed by Cooley and Tukey. It is possible to further optimise this algorithm using Graphics Processing Unit (GPU), which exploit the high parallelism of processors to further reduce the computational cost. Although algorithm 2 seems to be able to achieve a cost of  $O(\log N)$ , there are physical limitations of the GPU and in data structures that make its execution very fast, but asymptotically it remains of cost  $O(N \log N)$ .

**Periodogram** A periodogram graphically represents the amplitude of the different frequencies that make up a time series, highlighting the dominant components of the frequency spectrum. This tool is particularly useful for distinguishing between noise and the significant frequencies of a signal. For example, in the fields of climatology and astrophysics, it is often used to identify periodic cycles in data. We see an application of DFT on the `sunspots` dataset from the R package, which contains the number of sunspots observed each month from 1700 to 1988. Sunspots are correlated with the Sun's magnetic activity and one of the goals of the analysis is to identify any periodic cycles in the data.



Looking at the periodogram obtained with FFT, we observe that some frequencies are much more pronounced than others. In particular, the frequency  $\approx 0.0075$  Hz (corresponding to about a cycle of 135 months) indicates a cycle of about 11 years, which confirms what is already known about the cyclic nature of solar activity. This example demonstrates how the DFT is a powerful tool for detecting and analysing periodic patterns in observational data. The use of the periodogram not only makes it possible to identify the main frequencies, but also to evaluate their relative importance respect to the background noise.

## 2.4 Application

In this section, we will develop specific applications of the topics discussed in the chapter to see how they can be used in the context of the thesis. The discussion will focus on the use of fuzzy clustering to address the problem of continuity in colour space, the practical implementation of the FCM algorithm, and the use of FFT for image analysis.

### 2.4.1 Comparing works, the idea of clustering

As we have seen, a problem of eq. (2.3) becomes apparent when the continuous space in which the  $n$ -tiles are defined turns out to be too fit. The proposed idea was to fix the fewest possible boxes and work with the idea that 2 any tiles are as distinct as any other pair. However, this is not a credible assumption due to the loss of information, for this reason in this thesis we would like to approach the problem more dynamically by increasing the variety of tiles without running into sparsity problems.

The idea is to use the clustering of the merged data between the two samplings as a basis for fitting the continuous space. The expectation is that the result will

certainly be more accurate, however, it is necessary to consider the problem of the shape of the clusters (which will no longer be boxes).

*Exempli Gratia* (Distance between distributions reformulated). In this paragraph we will adopt the theoretically easy clustering KMeans. Let us take 10 000 samples from  $\mathcal{N}(-1, 0.25)$  and 40 000 from  $\mathcal{N}(+1, 1)$ , the two distributions will be denoted by  $\mathcal{A}$  and  $\mathcal{B}$  respectively.

The two samplings will be merged and clustered with KMeans using 32 centroids, resulting in a predictor  $\mathcal{P}$ . Each cluster will be viewed as a region with a certain measure that will be the mean square of the distances of each datum in the cluster from its centroid. Given a cluster of centroid  $c$ , we want to estimate:

$$\mu(c) = \sqrt{\mathbb{E}_{\mathcal{L}} [\|x - c\|^2 | \mathcal{P}(x) = c]}$$

where  $\mathcal{L}$  is the law of the two merged samples.

We can see in the first image that regions have been captured and that automatically the regions are denser where the frequency is higher. Furthermore, the image shows the weights of each cluster  $c$  defined as  $\mathbb{P}_{\mathcal{L}} [\mathcal{P}(x) = c]$

We now study the two samples separately over this clustering. The densities will be weighted on the new measure  $\mu$ , so the density on the centroid  $c$  of measure  $\mu(c)$  respectively for the distribution  $\mathcal{A}$  will be:

$$d_{\mathcal{A}}(c) = \frac{\mathbb{P}_{\mathcal{A}} [\mathcal{P}(x) = c]}{\mu(c)}$$

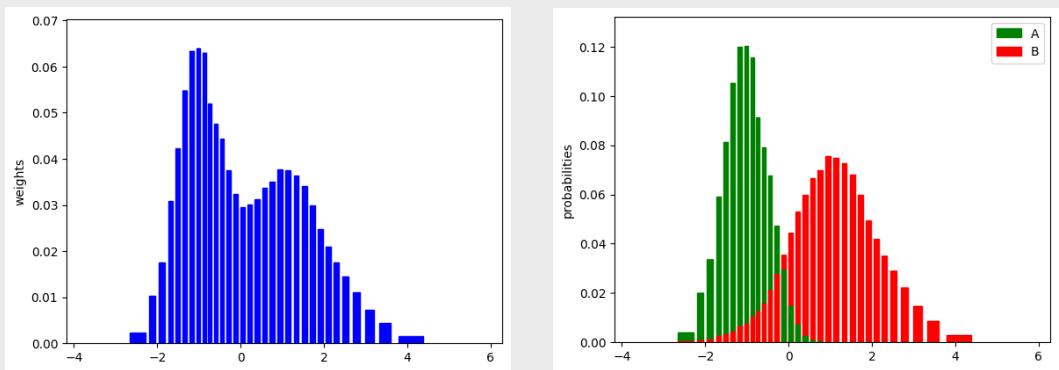
similarly for  $\mathcal{B}$ .

In the second figure we show the cluster membership probabilities of the two densities  $\mathcal{A}$  and  $\mathcal{B}$ .

Let us try to calculate the distance formulated eq. (2.3) considering the measure:

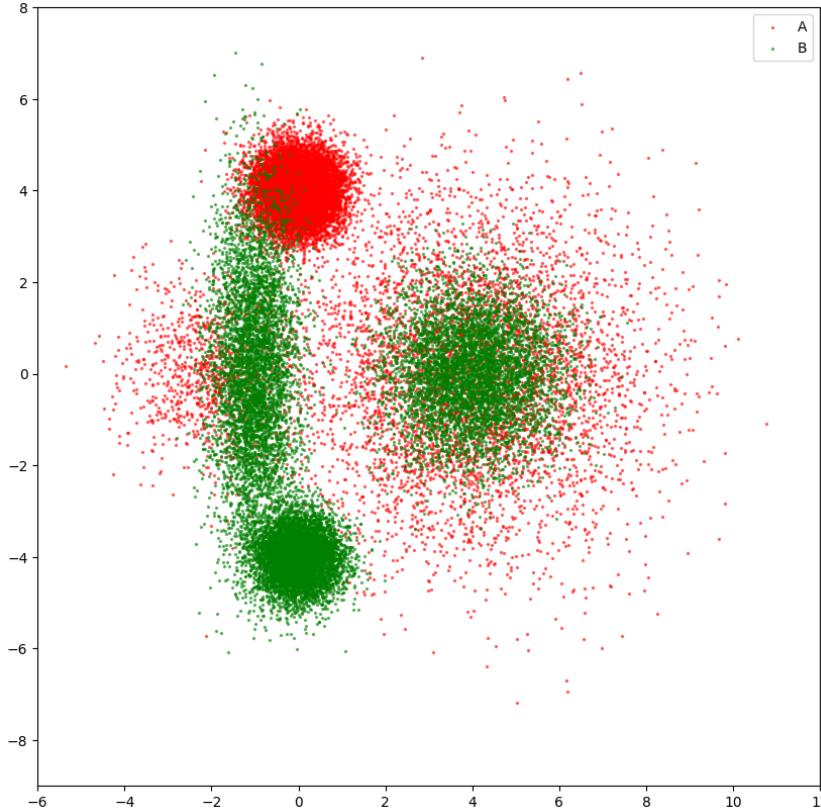
$$d(\mathcal{A}, \mathcal{B}) = (1 + J_{D_{\mathcal{A}}, D_{\mathcal{B}}})^{-1} \frac{1}{\mu(D_{\mathcal{A}} \cup D_{\mathcal{B}})} \int d\mu(x) \left( \frac{r(x) - 1}{r(x) + 1} \right)^2$$

where  $D_{\mathcal{A}}$  is the support for the discretisation of  $\mathcal{A}$  and similarly for  $D_{\mathcal{B}}$ . The result for 32 centroids is 0.54.



In the example, it is basically shown that the number of nodes can be reduced and can be adapted to the distributions of the tiles, and can also be very effective

at high dimensions by counteracting the curse of dimensionality. In fig. 2.3 we take an example very similar to the previous one but in 2 dimensions and with less data. Let us try clustering using more and more nodes, repeating the calculation carried



**Figure 2.3.** In the figure, the data from the set  $\mathcal{A}$  are shown in red and the data from the set  $\mathcal{B}$  in green. They were obtained by merging normal Gaussians with different representability.

out in the example, studying how the distance changes by increasing the number of nodes. In table 2.1 we can see how the clustering shows more stable results.

Number of nodes	1	4	16	64	256	1024
Distance with clustering	0.00	0.13	0.21	0.30	0.34	0.37
Distance without clustering	0.00	0.40	0.49	0.45	0.49	0.48

**Table 2.1.** Comparison between KMeans clustering and box tessellation with same number of nodes.

#### 2.4.2 Fuzzy Clustering as a Noise Filtering Method

In this thesis, a variant of the algorithm FCM will be used to compare two datasets. The variant will be used both to consider data of different weights and to consider machine error in the computation of centroids.

**data's weight** We introduce a vector  $w$  indicating the weight of the data as a positive real value. As stated in theorem 2.1 the update follows the following law:

$$C_j^{\text{new}} = \frac{\sum_{i=1}^N u_{ij}^2 x_i}{\sum_{i=1}^N u_{ij}^2} \quad \forall j$$

Since it is a weighted average over  $u_{ij}^2$ , if a datum has a higher weight then it should increase its influence, thus having the following result:

$$C_j^{\text{new}} = \frac{\sum_{i=1}^N u_{ij}^2 w_i x_i}{\sum_{i=1}^N u_{ij}^2 w_i} \quad \forall j$$

**machine error** The use of FCM in this thesis involves millions of pieces of data, and it is possible that the classical algorithm will find itself making serious machine errors that must be kept under control. An example that might exist is a value  $D_{ik}^2 \approx 0$  that may be null or so small that when  $D_{ij}^2/D_{ik}^2$  is computed, it is infinity or a number so large that when added to other numbers it is infinity.

For this reason algorithm 3 proposes a more robust approach. We also remark that the computational cost in a sequential algorithm will be  $O(N C k)$  however scalability allows us to reduce this cost to  $O(\log(kC))$  by exploiting the independence of each cycle and scalable binary operations (details in chapter 3).

$N$  is the number of data and  $C$  is the number of centroids,  $k$  is the size of the tiles

---

### Algorithm 3 Membership update stable computation.

INPUT

$\mathcal{S}$ : set of data  $x_1, \dots, x_N$   
 $\mathcal{C}$ : centroids  $c_1, \dots, c_C$

---

```

1: procedure MEMBERSHIPUPDATESTABLE( $\mathcal{S}, \mathcal{C}$ )
2:    $D^2 \leftarrow (d_{ij}^2)$  with  $d_{ij} = \|x_i - C_j\|^2$ 
3:   for  $i \leftarrow 0$  to  $N$  do
4:      $l \leftarrow \min_k \{D_{ik}^2\}$ 
5:     if  $l = 0$  then
6:       where  $D_{ij}^2 = 0$  do  $u_{ij} \leftarrow 1$ 
7:       where  $D_{ij}^2 \neq 0$  do  $u_{ij} \leftarrow 0$ 
8:     else
9:        $u_{ij} \leftarrow \frac{l}{D_{ij}^2} \quad \forall j$ 
10:    end if
11:     $S_i \leftarrow \sum_j u_{ij}$ 
12:     $u_{ij} \leftarrow u_{ij}/S_i \quad \forall j$ 
13:  end for
14: end procedure

```

---

### 2.4.3 Analysis of images with DFT

It was seen in the introductory section of DFT that the algorithm FFT can only be applied to time series or, more generally, to a sequence of complex numbers.

However, it is also possible to extend this concept to the analysis of the spectrum of a matrix with periodic behaviour.

Consider a matrix  $x \in \mathbb{C}^{N \times M}$  defined as follows:

$$x = (\cos(\omega_r n + \omega_c m))_{n,m}$$

It can be shown that its DFT results in a matrix of the same shape as  $x$ , where the only nonzero component is in the row corresponding to  $\omega_r$  and in the column corresponding to  $\omega_c$ .

This is analogous to the application of CFT on  $\mathbb{R}^2$ . In particular, we would like to obtain the following inverse relation to reconstruct  $x$  from its Fourier coefficients:

$$x_{n,m} = \frac{1}{NM} \sum_{r,c} X_{r,c} e^{i2\pi(\frac{rn}{N} + \frac{cm}{M})} \quad (2.8)$$

Where the Fourier coefficients  $X_{r,c}$  are given by:

$$X_{r,c} = \sum_{n,m} x_{n,m} e^{-i2\pi(\frac{nr}{N} + \frac{mc}{M})} \quad (2.9)$$

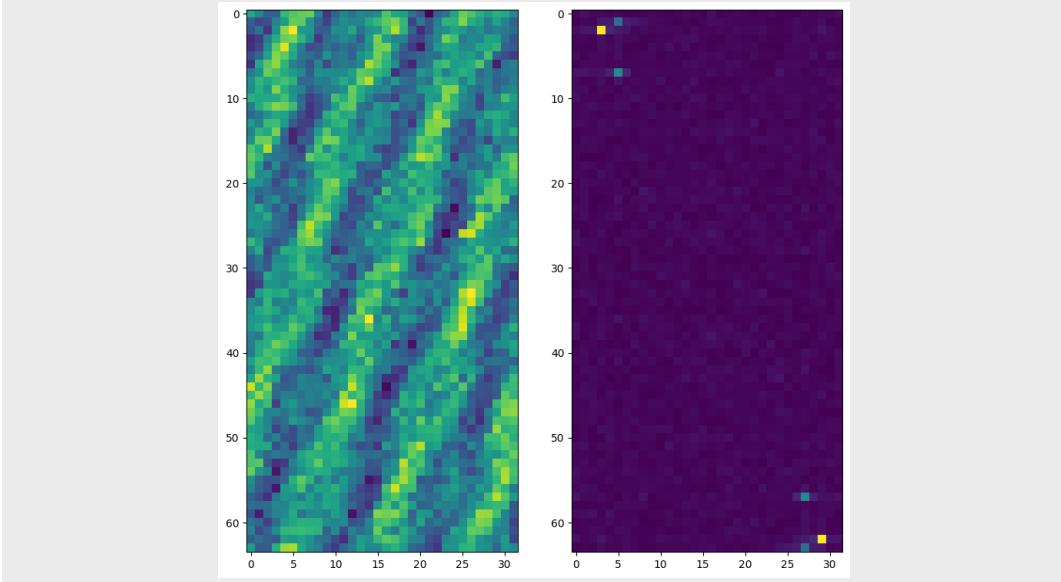
At the algorithmic level, the two-dimensional FFT is obtained by applying the algorithm to the columns first and to the rows of the original matrix  $x$ . In particular:

$$\begin{aligned} y_{r,m} &= \sum_n x_{n,m} e^{-i2\pi \frac{nr}{N}} && \text{over each column apply FFT} \\ X_{r,c} &= \sum_m y_{r,m} e^{-i2\pi \frac{mc}{M}} && \text{over each row apply FFT} \end{aligned}$$

*Exempli Gratia.* We analyse a matrix composed of several overlapping frequencies and Gaussian noise with variance 1.

$$\begin{aligned}x_{n,m} = & 2 \cos(2\pi(2n + 3m) + 3) \\& + 0.8 \cos(2\pi(n + 5m) + 2) \\& + \cos(2\pi(7n + 5m)) + \mathcal{N}_{n,m}\end{aligned}$$

In the figure, the Fourier coefficients are shown. They are computed from the matrix  $x$  using a `viridis` colour scale.



The two-dimensional FFT can be used to analyse periodic patterns in matrices that represent images or signals. Typical applications include filtering, image compression and pattern recognition, where the spectral decomposition allows significant components to be distinguished from the noise.

## Chapter 3

# Methodology

This research aims to analyse the attribution of graphic works, taking its inspiration from the methodologies used for the attribution of literary works [see 4]. In particular, it examines the application of  $N$ -grams, i.e. sequences of  $N$  adjacent symbols, as a method of representing works.

To Review

To fully understand the representation of graphic works, it is appropriate to start from the method used for literary works. For example, in the expression "Hello world!", a 3-gram can be represented by "llo", which corresponds to a sequence of 3 consecutive characters. Importantly, spaces and punctuation are also considered characters, so "o w" and "ld!" are also valid 3-grams.

Research in this field is extensive and has led to various applications. Representations based on  $N$ -grams are used not only in literary attribution, but also in natural language models in which these analyses are integrated with the use of recurrent neural networks. However, little is known so far about the application of this technique to graphic works, and this constitutes the main focus of this research.

**Discrete Automatic Drawings' Analysis (DADA)** In the thesis research about the application of  $N$ -gram analysis on calligraphic works, [see 4], graphics are considered as matrices of colours and an  $N$ -gram was defined as a square submatrix of size  $N$ , also called 'tile'.

The research focused on the analysis of images that have been reduced to matrices, in which the only colours present are black and white (see fig. 3.1). This process, known as 'posterisation', aims to reduce the amount of colours present in the work (known as 'depth'). This methodological choice was motivated by the need to manage the size of the alphabet; indeed, while in a literary work there may exist up to  $32^N$  distinct  $N$ -grams (considering the 26 letters of the alphabet plus 6 punctuation symbols), in a pictorial work there may be  $256^{3N^2}$  distinct tiles with side  $N$ . This considerable difference caused significant challenges in image analysis, rendering the application of tools designed for literary works on graphic works ineffective.

The research concluded with promising results, showing that the automatic analysis of handwriting samples allows very precise author attributions. However, it turned out that applying the same analysis to pictorial works such as panels and drawings does not produce satisfactory results. It was conjectured that this is due to the intrinsic nature of pictorial works, where the presence of colour plays an important



**Figure 3.1.** Example of a conversion from a coloured image to an image with only black and white (bw) pixels. The image is transformed to greyscale by reducing the saturation, finally half of the lightest pixels are rendered white and the remaining are rendered black.

role, and posterisation, instead of helping attribution, generates noise and new information that make the representation of data unreliable.

**Continuous Automatic Drawings' Analysis (CADA)** Following the same line of research as the thesis of Magrini Alunno [4], we have proposed a re-adaptation of the methodology in order to investigate what might be the most effective strategy for the automatic analysis of images characterised by a large alphabet of colours. The aim is to construct a CADA that can analyse colours in their natural space, i.e. in a continuous space. This will allow us to overcome the discrete constraint of DADA and obtain more accurate and detailed results.

Image analysis in DADA is a process in five stages: acquisition, pre-processing, synthesis, comparison and attribution. In this paper we are going to focus on the first four phases, detailing the process up to image comparison, while attribution will be briefly discussed as the concluding phase.

The acquisition phase digitises two-dimensional works of art by scanning or photography. Scanning, in particular, allows a high resolution and accurate representation of the original document, while photography allows greater flexibility when the size of the work or physical conditions make scanning complex. The result of this phase is a digital image that accurately represents the work to be analysed.

The pre-processing phase prepares the image for analysis by removing irrelevant elements through pixel-wise and work-size transformations. Pixel-wise techniques, such as greyscale conversion, act on each individual pixel. Work-size transformations, on the other hand, act on the image as a whole, e.g. with compression techniques or normalisation of colour scales.

The synthesis phase deconstructs the image into a list of tiles of size  $N \times N$ , generating a sequence of discrete elements representing portions of the image. This process, aligned with the idea of n-gram language models, allows stylistic features to be analysed in a  $\mathbb{R}^{N^2}$  space, preparing the data format for later comparison.

In the comparison phase, we are going to address the real computational and theoretical challenge. Clustering techniques will be used to dynamically discretise the space of tiles and organise them according to similarity criteria. This allows us both to efficiently process the large number of tiles in a very high dimension and also to compare the target work with known works in the database.

Finally, in the attribution phase, the similarities between the target work and known works are examined. The aim is to identify the author of the target work by studying the closest works in terms of stylistic and compositional characteristics. This phase focuses on identifying the works that show the greatest stylistic affinity with the target work, suggesting possible attributions.

This research not only aims to tackle technical and computational challenges, but also to explore new perspectives in the analysis of artworks, thus opening up new horizons in the field of art criticism and art attribution.

## 3.1 Data Set

### 3.1.1 Description

In this study, the data set consists of a set of graphic works for which attribution is established and unambiguous and for which a uniform resolution expressed in dots per inch (DPI) is known. It was decided to use calligraphic works, as evidenced in [4]. The use of calligraphic works represents a promising first way for artistic attribution.

The dataset used in the reference article consisted of authentic and unaltered handwriting samples. In this study, the dataset consists of university notes, which do not have a sufficiently high standard and therefore present some significant challenges:

- **The squares:** Note sheets often contain small squares that could be mistaken for the author's handwriting. This visual interference can complicate the attribution process.
- **Use a variety of writing tools:** University notes can be written using a variety of writing tools, such as highlighters, markers, pencils or white-out. These different writing tools produce different lines and colours, which can affect the attribution of authorship.
- **Different graphical contexts:** Notes can contain a variety of elements such as mathematical formulae, graphs and erasures. These features represent heterogeneous graphical contexts that add complexity to the analysis.
- **Scanning contamination:** The quality of the page scan can introduce blemishes into the record. These blemishes, such as stains or blurring, can affect the clarity and legibility of the works.

In addition, the use of a pen with a thickness of 0.4 mm on the sheet was considered a reasonable first choice since in [4] the stroke had the same thickness. This particular aspect is also important in ensuring a certain uniformity and consistency in the visual attributes of the works under examination, thus facilitating a more accurate attribution of authorship.

The dataset is organised into (indicare numero di autori) 5 authors, and the included images are as follows:

- The images are saved in the png format.
- Each image is composed of three colour channels Red Green Blue (RGB), each with a depth of 8 bits.
- The original sheets were scanned at a resolution of 400 DPI.
- The thickness of the pen used to write on the original sheets is 0.4 mm.

The images were then cut to remove large impurities or large white spaces. Finally leaving a total of: (scrivere quante immagini sono state raccolte e il peso complessivo del dataset)

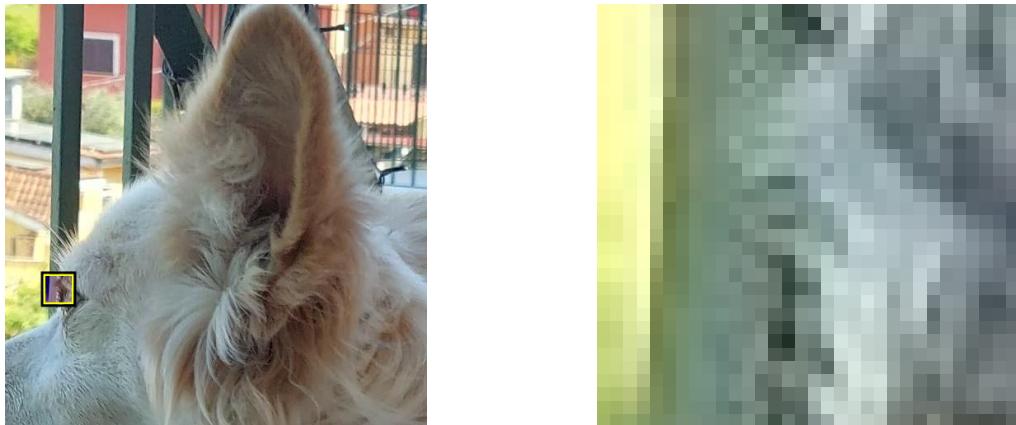
AuthorID	Num of items	Size	Mean side length
1	270	554 MB	1.43 Kpx
2	43	133 MB	1.8 Kpx
3	56	227 MB	2.0 Kpx
4	51	247 MB	2.2 Kpx
Total	420	1.2 GB	1.7 Kpx

**Table 3.1.** Each item is a cleaned area of the original images. The size represent the total number of pixels. The last column is the geometric mean of the side length, it's compute as follow:  $\sqrt{\frac{\text{Num of pixels}}{\text{Num of items}}}$

In the dataset creation, full clipboard images were collected and then cut to exclude non-pen strokes, highlights, and erasures. In addition, some cuts are particularly small and in a few cases are slightly overlapping.

$$\begin{aligned}
 m^*(B) \leq |B| \text{ per definizione di inf. dato che } B \subseteq \bigcup_{j=1}^{+\infty} B_j. \\
 \text{Ora devo dim. } |B| \leq m^*(B), \varepsilon > 0 \exists \{A_i\}_{i \in \mathbb{N}} \text{ box disgiunti tc } B \subseteq \bigcup_{i \in \mathbb{N}} A_i, \sum_{j=1}^{+\infty} |A_j| \leq m^*(B) + \varepsilon \\
 \forall j \in \mathbb{N}, \text{ sia } D_j \text{ un box aperto tc } |D_j| \leq |A_j| + \frac{\varepsilon}{2^j} \Rightarrow B \subseteq \bigcup_{j=1}^{+\infty} D_j \text{ aperti} \Rightarrow B \subseteq \bigcup_{j=1}^{+\infty} D_j \Rightarrow \\
 \Rightarrow |B| \leq \sum_{j=1}^{+\infty} |D_j| \leq \sum_{j=1}^{+\infty} |D_j| \leq \sum_{j=1}^{+\infty} |A_j| + \varepsilon \leq m^*(B) + 2\varepsilon
 \end{aligned}$$

## 3.2 Pre processing



**Figure 3.2.** Image detail.

Digitising images is a fundamental step in preparing the data set for analysis. To fully understand how these transformations work, it is important to understand the process by which a work of art is captured by a camera and then digitised in an electronic device. To Review

**Definition of image: from the real world to the virtual world** Artists of antiquity used various techniques, such as the camera obscura and the principles of projective geometry, to represent landscapes realistically. During the Renaissance, they further refined these methods with tools such as the camera lucida, showing advanced mathematical understanding in works of art such as Raphael's School of Athens.<sup>1</sup>

The invention of photography in 1826, in conjunction with the art movement of Realism, marked a fundamental moment in the history of visual art. Photography, with its ability to represent reality objectively, became a tool increasingly used by artists. The advancement of knowledge in optics and chemistry led to the creation of analogue photography.<sup>2</sup>

During the 1980s, with the advent of computer technology, digital photography became a reality. The basic operation remained similar to previous techniques: instead of being printed on a photosensitive surface, the image was captured by a grid of optical sensors and digitised as a matrix of colours, with each component called a pixel.

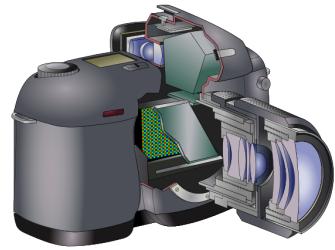
The concept of colour reproduction is not new and as early as 1855 James Clerk Maxwell worked on this subject, introducing the colour model RGB, which is still widely used today (see [3]). This model encodes each colour with three real numbers in an interval between 0 and 1. However, for computing purposes, it was agreed to represent the 3 values of a colour with integers between 0 and 255. In this way,

<sup>1</sup>for more about camera obscura, see  
[https://en.wikipedia.org/wiki/Camera\\_obscura](https://en.wikipedia.org/wiki/Camera_obscura)

<sup>2</sup>for more about history of photography, see  
[https://en.wikipedia.org/wiki/History\\_of\\_photography](https://en.wikipedia.org/wiki/History_of_photography)



(a) 'View from the Window at Le Gras', Joseph Nicéphore Niépce, c.1826, Heliography, Harry Ransom Center, University of Texas at Austin, USA[6].



(b) Illustration of digital camera[7].

a digital image becomes a matrix of triples, representing the pixel values the three colour channels RGB.

### 3.2.1 Grayscale reduction

The quantisation of colours is expressed in 3 channels RGB, each with an integer value between 0 and 255. The choice of these 3 colours is not random or arbitrary, but derives from a medical peculiarity of the human eye. The human eye senses colours through two types of photoreceptor cells: cones and rods. The cones are responsible for colour vision, while the rods are activated during night vision, which is monochromatic. In our case, we are interested in cones, which are only activated above a certain brightness threshold and are divided into 3 types: S-cones, M-cones and L-cones. These cones perceive frequencies close to the colours blue, green and red respectively.

When the colour yellow arrives, the cones responsible for green and red are activated. The same happens when there are two lights, one green and one red, very close together: our brain interprets them as a single yellow light. This mixture of lights deceives the human eye by mistaking them for pure colours. A remarkable example is magenta: this colour should not exist in its pure form. However, combinations of red and blue are still interpreted as pure colours, even though they are actually colours that do not exist in the physical world.

The human brain's interpretation of cone stimulus makes the shape of the colour space more a neurological than a physical question. The difference is such that while colour space is topologically a solid (often represented by cylinders, spheres or cones), physically existing colour space is topologically a polygon. Furthermore, the interpretation of colours is not the same for every human. For example, a colour-blind person may have difficulties due to cone quantity problems or neurological problems. However, in the absence of colour-blindness or differences in the number of cones, there are populations in the world that can distinguish shades of certain hues very well. For example, the Himba of Namibia easily distinguish different

To Review

shades of green, but have difficulty distinguishing green from blue<sup>3</sup>.

The need to classify colours must inevitably pass through subjective considerations. For this reason, in addition to the colour space RGB, which is related to medicine and physics, there are also colour spaces such as Hue Saturation Lightness (HSL), which consider hue, saturation and brightness. In an artistic context or where the choice of colour depends on the subjectivity of the author, it is preferable to speak in terms of HSL rather than RGB.

Greyscale reduction consists of nullifying the saturation of a colour, rendering the hue useless and leaving only the brightness. This process can be done by various techniques, as there is no absolute definition of ‘saturation’. For example, as a human being ages, he tends to perceive colours as less saturated due to the crystalline lens becoming duller, thicker or more yellowed. To reduce greyscale, a convex combination of the three channels RGB is used. Other techniques include averaging between the most intense and the least intense channel.



### 3.2.2 Removing squares

Una delle principali sfide nell’elaborazione e raccolta del dataset riguarda la pulizia delle immagini, che possono presentare impurità significative, come la presenza di quadretti sui fogli. È essenziale rimuovere o attenuare questi elementi affinché si confondano con lo sfondo e non interferiscano con l’analisi. Inizialmente, si è considerato l’uso di tecniche statiche, come il thresholding o l’applicazione di kernel specifici per il riconoscimento dei quadretti. Tuttavia, tali tecniche rischiavano di compromettere la scrittura dell’autore cancellando dettagli. Per questo motivo, si è scelto di ricorrere a tecniche di compressione per individuare i pattern, poiché i quadretti costituiscono un pattern ben visibile e più facilmente distinguibile dalla macchina rispetto alla scrittura umana.

**Singular Value Decomposition (SVD)** Un primo test ha coinvolto l’esame dei pattern predominanti tramite una decomposizione SVD. L’immagine può essere

---

<sup>3</sup>for more about Himba people, see [https://en.wikipedia.org/wiki/Himba\\_people](https://en.wikipedia.org/wiki/Himba_people)

vista come una matrice  $H \times W$  di valori reali e può essere decomposta nel prodotto:

$$U\Sigma V^H = M$$

dove  $M$  rappresenta l'immagine,  $U$  e  $V$  sono matrici unitarie ( $UU^H = I$ , stesso per  $V$ ),  $V^H$  è l'hermitiana di  $V$ , e  $\Sigma$  è una matrice  $H \times W$  con solo valori reali lungo la diagonale, chiamati valori singolari.

È possibile ottenere una versione meno dettagliata di  $M$  annullando alcuni valori singolari di  $\Sigma$ , ignorando così il contributo di determinati pattern. Si è osservato che il pattern predominante (associato al valore singolare maggiore) riguarda i quadretti, ma la sua rimozione non risolve completamente il problema, richiedendo l'eliminazione di ulteriori pattern. Il risultato ottenuto è incoraggiante, ma comporta una perdita di informazioni dalla scrittura umana, poiché anch'essa è presente tra i principali pattern identificati.

**FFT** Un'idea alternativa è effettuare un'analisi armonica delle immagini utilizzando la FFT. I quadretti si ripetono con una frequenza specifica lungo le due dimensioni del foglio e quindi si può sfruttare la loro natura regolare rispetto alla scrittura umana, che è meno ripetitiva. Così si sfrutta sia la differenza tra i pattern regolari dei quadretti e i tratti irregolari della scrittura, sia una maggiore capacità di compressione rispetto a una semplice analisi pixel per pixel.

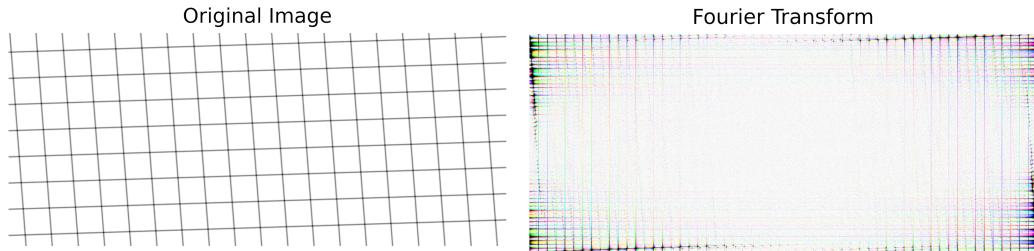
Per verificare questa idea, si esegue una FFT su immagini di prova: un foglio virtuale con solo quadretti e lo stesso foglio con dei disegni sopra. Nel caso del foglio senza segni, le frequenze con ampiezza maggiore si trovano principalmente allineate lungo gli assi  $x$  e  $y$  (fig. 3.5a). Aggiungendo un disegno, le frequenze principali restano pressoché le stesse (fig. 3.5b). Per rimuovere i quadretti, si eliminano queste frequenze significative e si ricostruisce l'immagine (fig. 3.5c).

Empiricamente, analizzando 10 immagini reali selezionate casualmente, si è osservato che i quadretti sono stati rimossi con successo senza danneggiare le scritte, annullando le ampiezze corrispondenti al primo 5% percentile delle frequenze più significative.

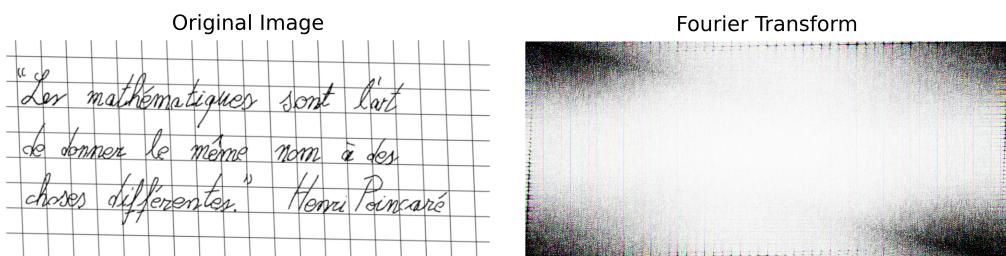
Tuttavia, questa operazione distorce parzialmente i colori. Infatti, in un'immagine originale, la distribuzione dei livelli di grigio è piuttosto intuitiva: la maggior parte dei pixel è bianca, mentre solo una piccola percentuale è composta da tonalità di grigio più scure. Dopo la ricostruzione, si nota la presenza di molte tonalità di grigio intermedie, che interessano sia lo sfondo che le scritte.

Si è osservato che, nelle immagini originali, i pixel con un livello di grigio inferiore a 0.2 appartengono certamente alle scritte, mentre quelli con un livello di grigio superiore a 0.8 corrispondono allo sfondo bianco della pagina. Di conseguenza, si contano i pixel delle scritte e dello sfondo bianco nell'immagine originale. Nell'immagine ripulita, i pixel più scuri (con un valore di grigio inferiore a 0.2 nell'immagine originale) vengono assegnati al nero, mentre i pixel più chiari (con un valore di grigio superiore a 0.8 nell'immagine originale) vengono assegnati al bianco.

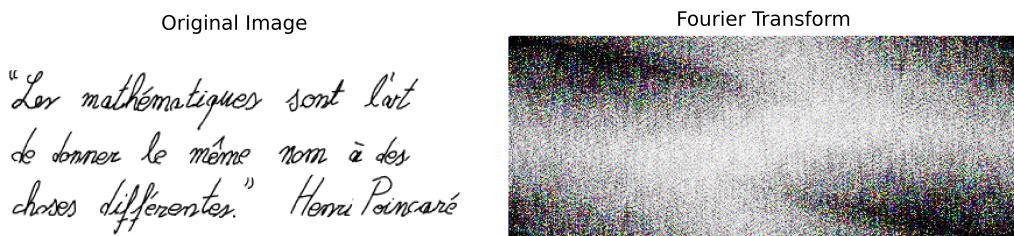
Questo processo viene realizzato tramite una "normalizzazione con morsetto": i colori vengono espansi in modo tale che i pixel che devono essere neri assumano un valore di grigio inferiore a 0, mentre i pixel bianchi superano il valore di 1. Successivamente, si imposta a 0 qualsiasi valore negativo e a 1 qualsiasi valore maggiore di 1. Un risultato di questa procedura è mostrato in fig. 4.1



- (a) Un reticolato lievemente trasformato, causato da una scansione imperfetta. Le frequenze più significative sono evidenziate in nero.



- (b) Un reticolato deformato con una scrittura sovrapposta. Le frequenze del reticolato restano visibili.



- (c) La scrittura senza reticolo non mostra coefficienti di Fourier significativi, poiché mancano pattern ricorrenti.

**Figure 3.5.** Coefficiente di Fourier dell'immagine bidimensionale. I quadretti mostrano coefficienti di Fourier simili a quelli delle onde quadre:  $\text{sqwave}(t) = \frac{4}{\pi} \sum_{k=1}^{\infty} \frac{\sin(2\pi(2k-1)t)}{2k-1}$ . Le intensità più significative sono a intervalli regolari con intensità decrescenti.

---

**Algorithm 4** Cleaning procedure using FFT.

INPUT

 $\mathcal{I}$ : Image as matrix  $W \times H$  in gray scale $p$ : percentile of significant magnitudes $a$ : thresholds (min, max)OUTPUT  $\mathcal{I}_{\text{new}}$ : Image as matrix  $W \times H$  in gray scale

---

```

1: procedure CLEANFFT( $\mathcal{I}, p, a$ )
2:    $\hat{\mathcal{I}} \leftarrow \text{normalize}(\mathcal{I})$ 
3:    $F_{\hat{\mathcal{I}}} \leftarrow \text{fft2d}(\hat{\mathcal{I}})$ 
4:    $f = (f_1, \dots, f_{W \times H}) \leftarrow \text{argsort}(\text{flatten}(F_{\hat{\mathcal{I}}}), \text{rule}=\text{abs})$ 
5:    $f_{\text{best}} \leftarrow f[\text{last } \lceil p \cdot \text{length} \rceil \text{ values}]$             $\triangleright$  components with high magnitude
6:   for  $k \in f_{\text{best}}$  do
7:      $F_{\hat{\mathcal{I}}}[k] = 0.0$ 
8:   end for
9:    $\hat{\mathcal{I}} \leftarrow \text{ifft2d}(F_{\hat{\mathcal{I}}})$                                  $\triangleright$  rebuild normalized image
10:   $w_{\text{cnt}} \leftarrow \text{count } \mathcal{I} \geq a(\max)$            $\triangleright$  normalization using dilation and clamp
11:   $b_{\text{cnt}} \leftarrow \text{count } \mathcal{I} \leq a(\max)$ 
12:   $h_w \leftarrow \max_{w_{\text{cnt}}} \hat{\mathcal{I}}$ 
13:   $h_b \leftarrow \min_{b_{\text{cnt}}} \hat{\mathcal{I}}$ 
14:   $\mathcal{I}_{\text{new}} \leftarrow \frac{\hat{\mathcal{I}} - h_b}{h_w - h_b}$ 
15:  where  $\mathcal{I}_{\text{new}} > 1$ ,  $\mathcal{I}_{\text{new}} = 1$  do
16:  where  $\mathcal{I}_{\text{new}} < 0$ ,  $\mathcal{I}_{\text{new}} = 0$  do
17:  return  $\mathcal{I}_{\text{new}}$ 
18: end procedure

```

---

### 3.3 Synthesis

The tessellation phase extracts tiles  $N \times N$  from each image. Starting from a matrix of grey pixels, represented as a matrix in  $[0, 1]^{h \times w}$ , where  $h$  is the height and  $w$  the width of the image.

Al termine di questa procedura, ogni immagine viene rappresentata come un elenco di tessere, motivo per cui questa fase è chiamata "sintesi": infatti, risulta difficile ricostruire l'immagine originale a partire dalle tessere. Questo processo complica notevolmente qualsiasi tentativo di falsificazione dell'opera, poiché ricostruire manualmente un'immagine con una distribuzione di tessere simile a quella originale è estremamente complicato. Inoltre, la dimensione delle tessere nella rappresentazione è ridotta, ma sufficientemente grande da catturare piccoli automatismi della mano dell'autore, dettagli che sono difficilmente riproducibili in modo volontario.

**Algorithm 5** Algorithm for tile extraction

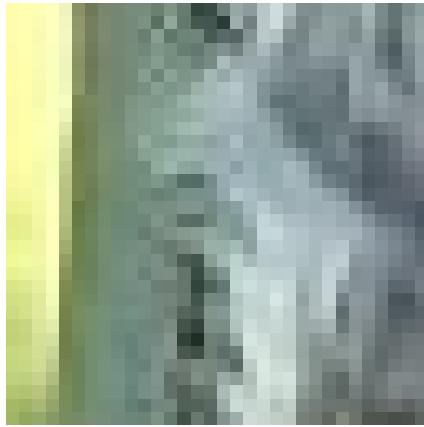
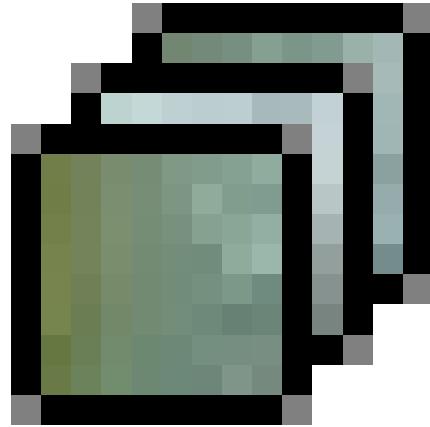
---

```

1: function TILESEXTRACTION( $M, h, w$ )
2:   declare  $v$  as matrix  $N \times N$                                  $\triangleright N$  is size of tiles
3:    $L \leftarrow$  empty list            $\triangleright$  will have  $(h - N + 1) \times (w - N + 1)$  elements
4:   for  $row \leftarrow 0$  to  $h - N$ ,  $col \leftarrow 0$  to  $w - N$  do
5:     for  $i \leftarrow 1$  to  $N$ ,  $j \leftarrow 1$  to  $N$  do
6:        $v[i][j] \leftarrow M[row + i][col + j]$ 
7:     end for
8:     call Append( $L, v$ )
9:   end for
10:  end function

```

---

(a) Image detail  $32 \times 32$ .(b) List of tiles  $8 \times 8$ .**Figure 3.6.** The synthesis process convert an image in a list of its tiles.

## 3.4 Comparison

Per comparare due opere si è deciso di usare la clusterizzazione FCM in grado di dipendere meno da eventuale rumore, soprattutto perché la comparazione tra due opere farà leva sulle densità più ridotte e quindi il rumore potrebbe inquinare i risultati.

In questa sottosezione vedremo in dettaglio l'algoritmo per la comparazione tra opere e si vedrà un esempio di applicazione per indagare al meglio le sue caratteristiche.

**Definizioni** Come si è già accennato in chapter 2, servirà fornire una definizione di misura del cluster e di probabilità di appartenenza ad un cluster da parte di una distribuzione. Ad esempio per KMeans la misura di un cluster è stata definita come:

$$\mu(c) = \sqrt{\mathbb{E}_{\mathcal{L}} [\|x - c\|^2 | \mathcal{P}(x) = c]}^k \approx \sqrt{\frac{1}{|\{x : \mathcal{P}(x) = c\}|} \sum_{x:\mathcal{P}(x)=c} \|x - c\|^2}^k$$

dove  $k$  è la dimensione dello spazio dei dati,  $\mathcal{P}$  è il predittore di KMeans,  $\mathcal{L}$  è la distribuzione fusa delle due distribuzioni,  $c$  è il centroide.

Nel nostro caso non esiste un predittore poiché è stata definita una misura di appartenenza di ogni dato ad ogni cluster. Tuttavia si ricorda che FCM nasce da una generalizzazione di KMeans e che dunque si può riprendere la formula minimizzata dall'algoritmo EM (eq. (2.4))

**Definition 3.4.1** (misura di un cluster). Applicando FCM al data set  $\mathcal{S}$  con pesi  $w$ , otteniamo la misura di appartenenza  $\mu$  e i centroidi  $\mathcal{C}$ .  $\mu_x(c)$  sarà così l'appartenenza del dato  $x \in \mathcal{S}$  al cluster  $c \in \mathcal{C}$ . Chiamato  $k$  la dimensione dello spazio dove sono istanziati i dati di  $\mathcal{S}$ , si denota con  $\mu(c)$  il peso del cluster  $c$ . Il peso sarà il volume della sfera cui raggio è la radice della media quadratica delle distanze del cluster:

$$u_{ij}^2 := w_i \mu_{x_i}(C_j)^2 \quad \text{or also} \quad u_{xc}^2 := w_x \mu_x(c)^2 \quad (3.1)$$

$$\mu(c) := \sqrt{\frac{\sum_{x \in \mathcal{S}} u_{xc}^2 \|x - c\|^2}{\sum_{x \in \mathcal{S}} u_{xc}^2}}^k \quad (3.2)$$

Definita la misura dei cluster sarà necessario poi parlare del peso di appartenenza della distribuzione  $\mathcal{A}$  (o analogamente  $\mathcal{B}$ ) ad un cluster  $c$ . In KMeans questa appartenenza è stata facilmente attribuita alla probabilità:  $\mathbb{P}_{\mathcal{A}}[\mathcal{P}(x) = c]$ . Come già detto in precedenza, FCM non propone un predittore, quindi bisognerà fornire una definizione più generale. In KMeans si può riformulare la probabilità di appartenenza come segue:

$$\mathbb{P}_{\mathcal{A}}[\mathcal{P}(x) = c] = \int d\mu_{\mathcal{A}}(x) \delta_x(c) = \int d\mu_{\mathcal{A}}(x) \mu_x(c)^2 \approx \frac{1}{|A|} \sum_{x \in A} \mu_x(c)^2$$

Ne consegue la definizione per FCM:

**Definition 3.4.2** (peso di un cluster). Si denota con  $\omega_A(c)$  il peso del set  $A$  sul cluster  $c$  ottenuto con FCM.

$$\omega_A(c) := \frac{\sum_{x \in A} u_{xc}^2}{\sum_{x \in A} w_x} \quad (3.3)$$

Ora bisognerà definire l'indice di Jaccard, obiettivo non semplice poiché nella logica fuzzy la cardinalità di un insieme non è definita come nella logica booleana, poiché l'appartenenza stessa non ha un valore di verità binario. Inoltre si fa notare che è possibile fornire una definizione dettagliata dell'indice di Jaccard che non dipende da quanto detto fino ad ora, vedendo i due valori moltiplicati di natura differente. Ad esempio, possiamo essere portati a pensare che

$$|D_A \cup D_B| = \sum_c \mu(c)$$

a giudicare da quanto detto in definition 3.4.1, tuttavia questa definizione è molto scomoda dal momento che non ci suggerisce niente in merito a  $|D_A \cap D_B|$ . Trattiamo quindi questa uguaglianza come una coincidenza nel caso di KMeans che mette d'accordo un'eventuale definizione dell'indice di Jaccard con la definizione dell'integrale approssimato con Monte Carlo method.

Quindi pensiamo a  $D_A$  e  $D_B$  come dizionari, ossia l'elenco dei centroidi che riguardano  $A$  e  $B$ .

**Definition 3.4.3.** (Indice di Jaccard) Si definiscono le operazioni di unione e intersezione in logica fuzzy come segue:

$$\begin{aligned}\mu_c(D_A \cup D_B) &= \max\{\mu_c(D_A), \mu_c(D_B)\} \\ \mu_c(D_A \cap D_B) &= \min\{\mu_c(D_A), \mu_c(D_B)\}\end{aligned}$$

e si definisce la cardinalità di un insieme  $S$  come:  $|S| := \sum_c \mu_c(S)$   
L'appartenenza di un cluster  $c$  al dizionario  $D_A$  sarà così definita:

$$\mu_c(D_A) = \max_{x \in A} (u_{xc}^2) \quad (3.4)$$

analogamente per  $\mu_c(D_B)$ .

Per questa ragione la definizione di indice di Jaccard che ne segue sarà:

$$J_{D_A, D_B} := \frac{|D_A \cap D_B|}{|D_A \cup D_B|} = \frac{\sum_c \min\{\max_{x \in A} (u_{xc}^2), \max_{x \in B} (u_{xc}^2)\}}{\sum_c \max\{\max_{x \in A} (u_{xc}^2), \max_{x \in B} (u_{xc}^2)\}} \quad (3.5)$$

definitions 3.4.1 to 3.4.3 sono sufficienti per formulare la distanza tra due campionamenti proposta a seguito di FCM.

**L'algoritmo** L'algoritmo prevede inizialmente di applicare FCM al set fuso di dati delle due opere, al fine definire la misura e la discretizzazione dello spazio delle tessere. Non sarà necessario aver memoria dei centroidi, sarà altresì utile poter calcolare:

- $\mu(c) = \sqrt{\frac{\sum_{x \in S} u_{xc}^2 \|x - c\|^2}{\sum_{x \in S} u_{xc}^2}}^k$
- $\omega_A(c) = \frac{\sum_{x \in A} u_{xc}^2}{\sum_{x \in A} w_x}$  e analogamente  $\omega_B(c) = \frac{\sum_{x \in B} u_{xc}^2}{\sum_{x \in B} w_x}$
- $J_{D_A, D_B}$  definito in definition 3.4.3

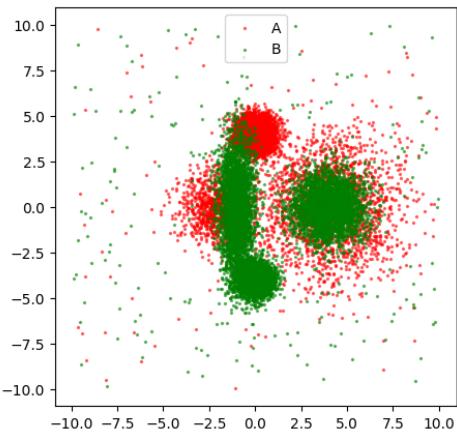
La distanza tra opere sarà così definita come:

$$d_{FCM}(A, B) = (1 + J_{D_A, D_B})^{-1} \frac{\sum_c \left( \frac{r(c)-1}{r(c)+1} \right)^2 \mu(c)}{\sum_c \mu(c)} \quad (3.6)$$

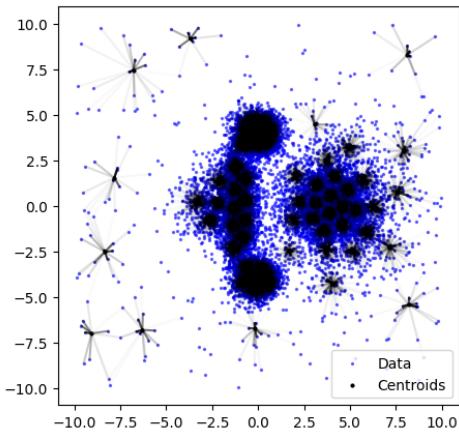
dove  $r(c)$  è il rapporto tra le densità ossia tra  $\omega_A(c)$  e  $\omega_B(c)$ .

Qui di seguito è mostrata una semplice applicazione per un set sintetico di dati. Come fatto in chapter 2 si generano dei dati in uno spazio bidimensionale e si clusterizzano con FCM in 64 cluster.

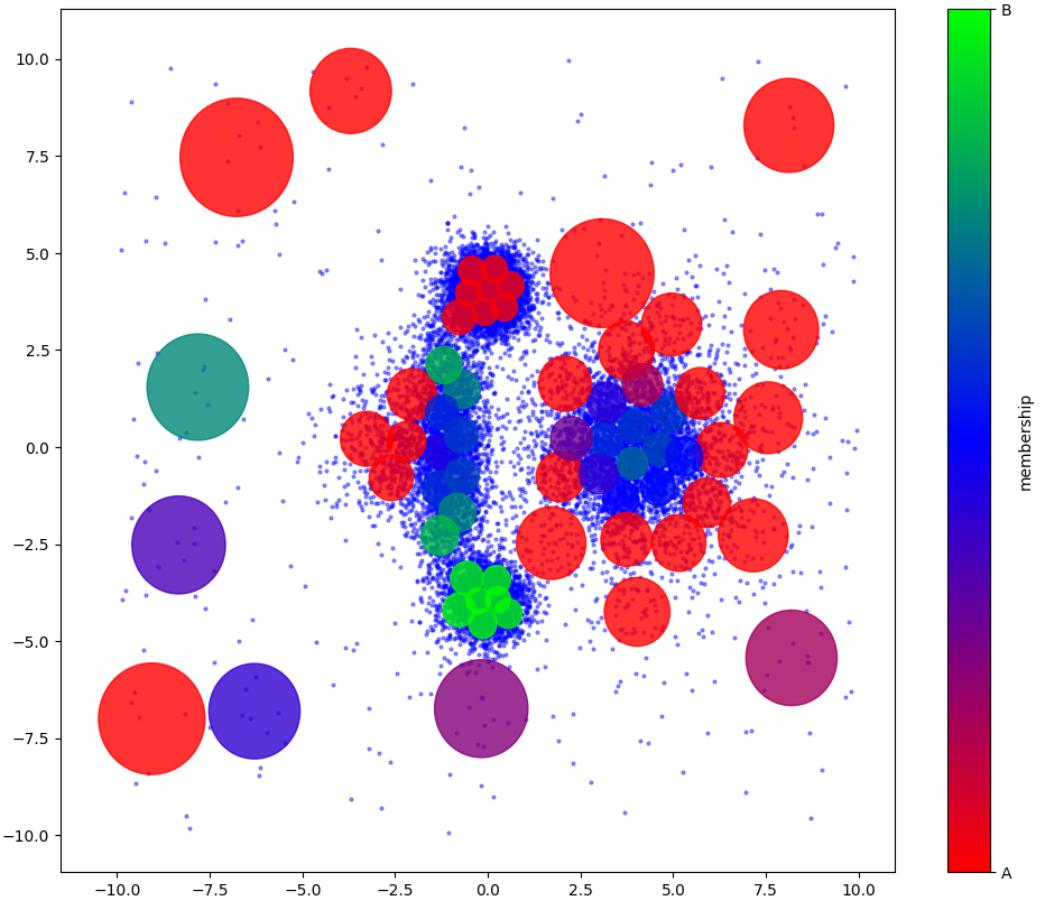
In fig. 3.7a sono mostrati i dati usati, si tratta di normali gaussiane con una certa rappresentanza e del rumore uniforme. I due set, chiamati  $A$  e  $B$ , sono stati poi fusi e clusterizzati con FCM (in fig. 3.7b una rappresentazione). In fig. 3.7c sono raffigurati i cluster con le rispettive rappresentabilità tra il set  $A$  e  $B$  attraverso una mappa di colori.



(a) In figura sono mostrati i due set di immagini.



(b) In figura è mostrata una clusterizzazione FCM con 64 centroidi, attraverso dei segmenti sono rappresentate le relazioni significative tra dati e cluster.



(c) In figura sono mostrati i cluster aventi un'area che copre la misura del cluster e un colore che ne indica la caratterizzazione, in particolare la scala di colori è mappata su  $\frac{\omega_A - \omega_B}{\omega_A + \omega_B}$ .

Calcoliamo l'indice di Jaccard:

- $|D_A| = 53.17$  e  $|D_B| = 52.26$
- $|D_A \cap D_B| = 43.91$  e  $|D_A \cup D_B| = 61.53$
- $J_{D_A, D_B} = 0.71$

Stimo l'integrale e di conseguenza la distanza tra le distribuzioni:

- $\frac{\sum_c \left( \frac{r(c)-1}{r(c)+1} \right)^2 \mu(c)}{\sum_c \mu(c)} = 0.31$
- $(1 + J_{D_A, D_B})^{-1} = 0.58$
- $d(A, B) = 0.18$

Osserviamo che il risultato pur essendo in presenza di rumore non è molto dissimile dal risultato ottenuto con KMeans in assenza di rumore (vedere table 2.1). Un altro esempio è con 256 nodi dove la distanza sarà 0.38 quindi ancora una volta non ci sono problemi legati al rumore.

Verifichiamo la potenzialità di FCM esaminando le distanze computate da KMeans in presenza di rumore:

64: la distanza stimata è 0.12 contro la precedente 0.30.

256: la distanza stimata è 0.23 contro la precedente 0.34

La riduzione della distanza è un fenomeno atteso, poiché il rumore è identico per entrambi i campionamenti rendendoli quindi più simili. FCM può lievemente pulire i set dal rumore e quindi avere una stima più precisa.

**GPU** The algorithm 3 shown is a sequential solution for FCM. It works by iterating through the data and centroids to compute the membership matrix. However, this approach has a computational cost of  $O(NCk)$ .

To improve computational efficiency, the use of GPU boosting can be employed. This kind of operation is known as General-Purpose computing on Graphics Processing Units (GPGPU). Exploiting the parallel computing power offered by a GPU can greatly accelerate the process of comparison.

The GPU is an electronic component present in every computer, able to perform a large number of operations in parallel. Originally designed to handle the graphical interface in video games, the GPU is capable of handling billions of pixels on any computer screen at speeds that the Central Processing Unit (CPU) cannot achieve. This processor is made up of thousands of threads, organised hierarchically at the hardware level to maximise performance:

- i. Stream Multi-Processing (SM): Runs a kernel and consists of numerous warps;
- ii. warp: Runs the kernel of its stream and has shared memory between its threads, usually 32;

- iii. thread: Executes the kernel of its warp by synchronising with the other threads of the same warp and has its own reserved memory in its registers.

The memory that a GPU can access is divided into different categories, namely *global*, *shared*, *cache* and *register* memory. Access to these memories by the processor's threads depends on the hierarchy of the threads themselves. For example, the *global* memory is accessible by every thread, while the *shared* memory is only accessible by threads in the same warp.

In a high-end computer, it is common to find a GPU equipped with 14 SM, each of which contains 1024 threads divided into 32 warps. This total of 14336 threads can execute the exact same code in parallel, permitting extremely fast processing of images and other operations requiring a high degree of parallelism.

Since the 2000s, the use of GPU has extended to the field of scientific computing, introducing important concepts such as scalability and High Performance Computing (HPC). Since 2020, dedicated GPU are available on the market for artificial intelligence operations.

In Python, there are useful frameworks for the utilisation of GPU, such as *torch* and *TensorFlow*, which are widely employed in the field of computer vision. However, also languages such as C++ offer dialects that allow these powerful computing units to be exploited. In this paper, the Compute Unified Device Architecture (CUDA) dialect will be used.<sup>4</sup><sup>5</sup>

The integration of GPGPU techniques would allow the workload to be distributed over several cores of the GPU, thus reducing the time needed for clustering operations. This method is particularly advantageous when handling large amounts of data, as the GPU can perform many operations in parallel, speed up computation to be 2000 times faster than the CPU could have done.

The sum of  $N$  numbers can be performed in with computational cost  $O(\log(N))$ . This is because in parallel the GPU threads sum one half of the vector over another at the same instant and then repeat until they get a single component that will have only one number. This operation is called *reduction* and we can see it in the algorithm 6. This is just one detail of how the GPU can reduce the asymptotic computational cost of an algorithm. Suffice it to say that thanks to reductions and strong parallelism, it is possible to multiply two  $N \times K$  and  $K \times M$  matrices with cost  $O(\log(K))$  instead of  $O(NMK)$ . In clustering many operations can be parallelised and FCM in particular requires many sums and linear operations.

The limitations of GPU are not only related to the execution of the same operations on all threads, but also to the nature of these operations. Normally, an instruction takes much longer to be executed by a GPU than by a CPU. Arithmetic instructions are the most efficient, while the use of conditions tends to be avoided.

---

<sup>4</sup>for more details about GPU architecture, see  
<https://researchcomputing.princeton.edu/support/knowledge-base/gpu-computing>

<sup>5</sup>for more details about CUDA language, see  
<https://docs.nvidia.com/cuda/>

---

**Algorithm 6** Parallel algorithm sum reduction

---

```

1: procedure KERNEL i, SUMREDUCTION(v,N)
2:   Let  $S$  a shared vector with  $2^k \geq N$  component
3:    $S[i] \leftarrow v[i]$  if  $i < N$  else 0
4:   for  $L \leftarrow 2^k/2, 2^k/4, \dots, 1$  do
5:     if  $i < L$  then
6:        $S[i] \leftarrow S[i] + S[i + L]$ 
7:     end if
8:     require synchronisation between threads
9:   end for
10: end procedure

```

---

1	2	3	4	5	6	7
1+5=6	2+6=8	3+7=10	4	5	6	7
6+10=16	8+4=12	10	4	5	6	7
16+12=28	12	10	4	5	6	7

**Figure 3.8.** We want to calculate the sum of the values in the first row. The idea is to divide the vector into 2 regions, sum the components, and repeat over the new vector with half the size of the previous vector.

# Chapter 4

## Results

Mostrare i risultati delle varie fasi e infine le analisi utili per l'attribuzione

To Review

### 4.1 pre processing

Mostrare i risultati del pre processing.

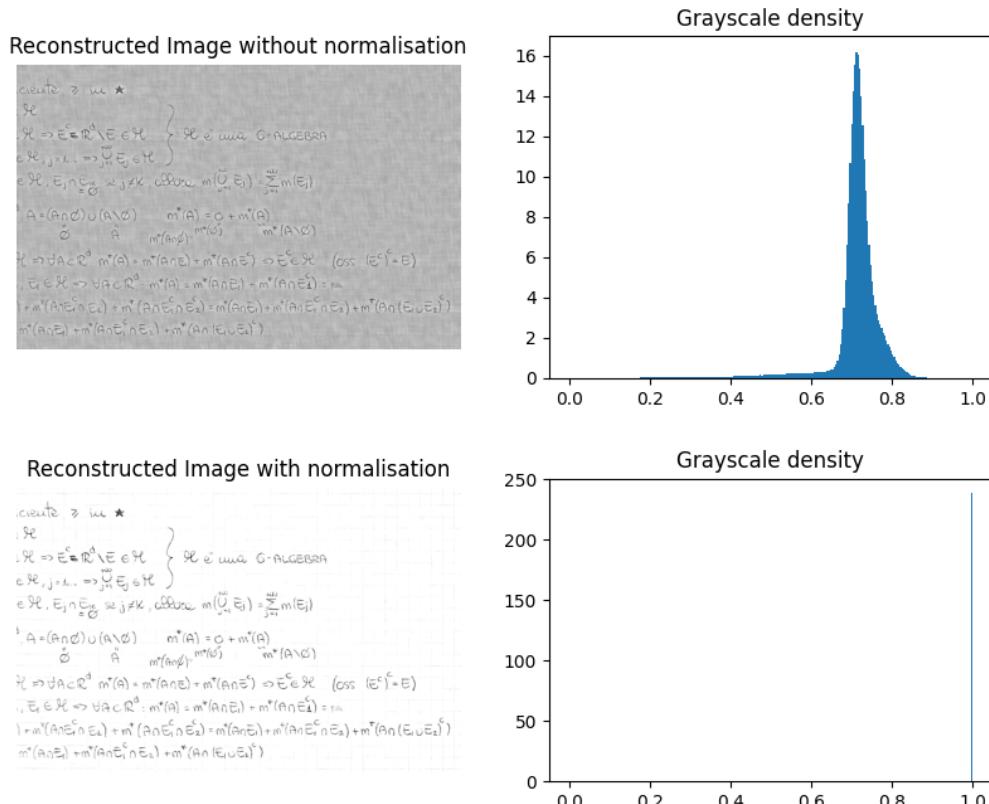
- serie di fourier applicate a griglie sintetiche per dedurre quali sono le frequenze di nostro interesse: griglie dritte e ruotate, griglie chiare e scure, griglie con elementi inquinanti, confronto tra teoria e aspettativa
- serie di fourier applicate a immagini vere, confronto con le immagini sintetiche
- risultati della pulizia dei quadretti

### 4.2 synthesis

Mostrare una piccola analisi delle sintesi ottenute e della velocità del programma

### 4.3 analysis

Mostrare i risultati delle analisi con i vari clustering



**Figure 4.1.** Come si osserva i grigi si sono molto avvicinati tra loro dopo FFT, ma è possibile ribilanciare i colori confrontandoli con l'immagine originale.

# Chapter 5

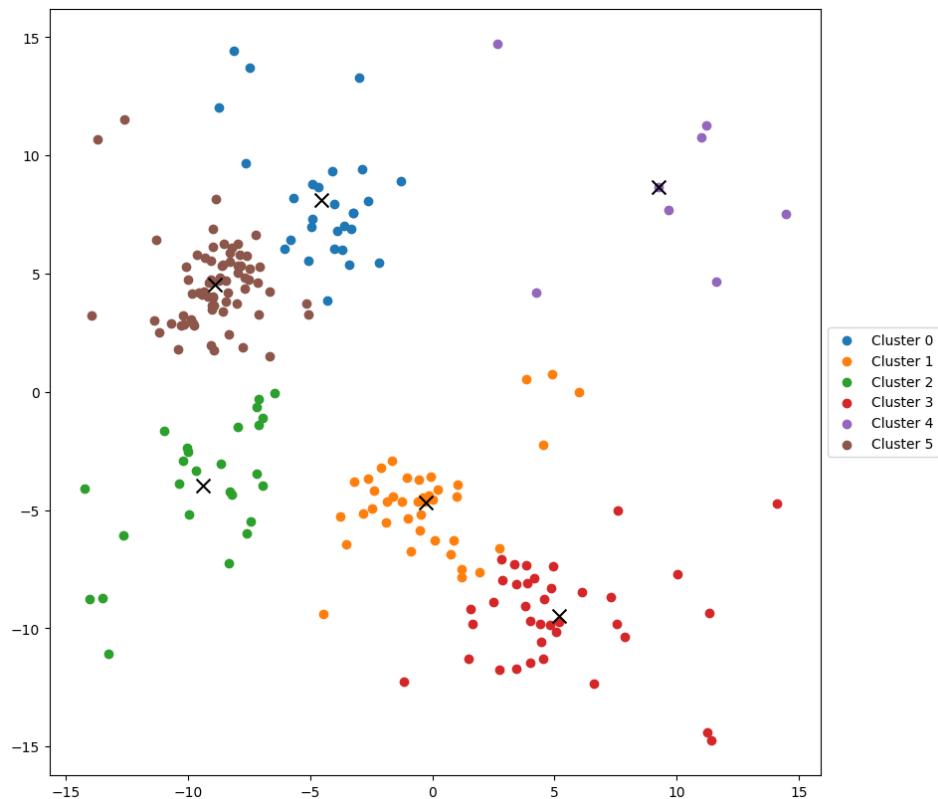
## Discussion

# Chapter 6

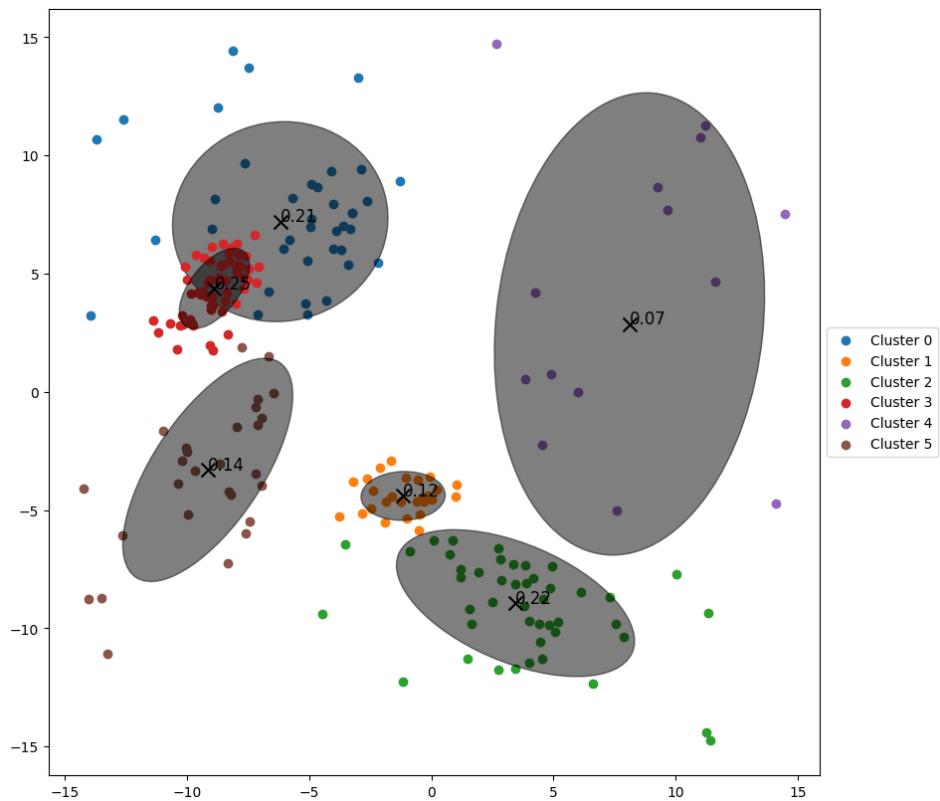
# Conclusion

## Appendix A

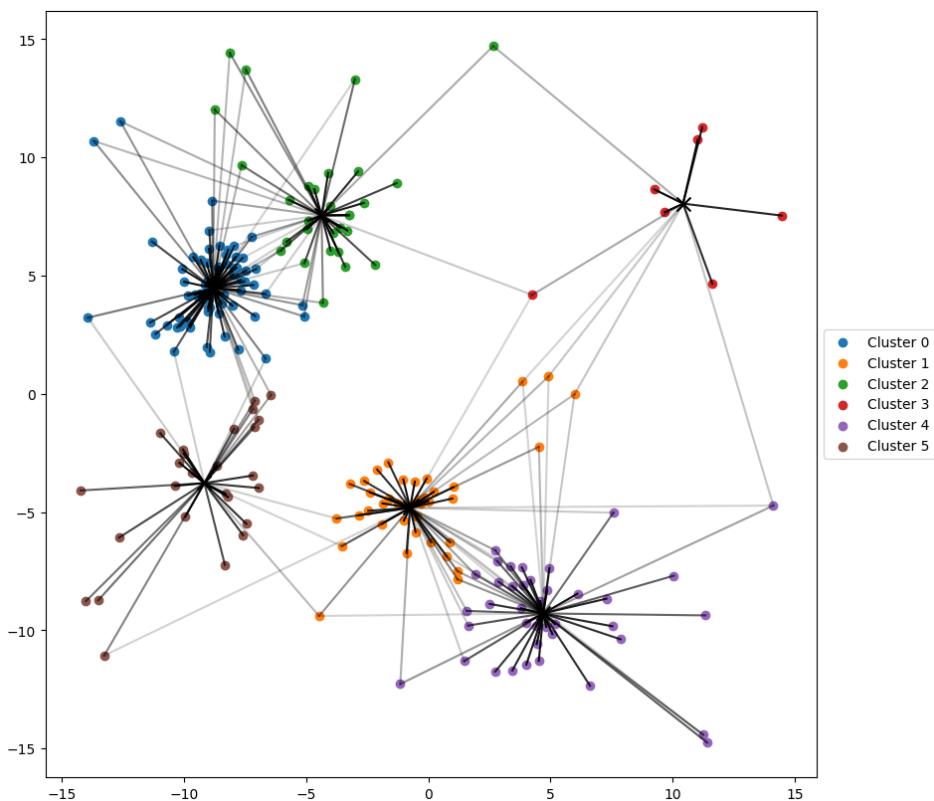
### KMeans, GMM, FCM compare figures



**Figure A.1.** The data points are coloured according to the calculated label and the estimated centroid is indicated with a X.



**Figure A.2.** The data points are coloured according to the Mahalanobis distance and the estimated centroid is indicated with a X. The ellipse of the normal distribution represents the covariance and is a confidence region of 95%.



**Figure A.3.** The data points are coloured according to the most probable label and the estimated centroid is indicated with a  $\times$ . The lines indicate the assignments with the greatest degree of affiliation of each point to the clusters; the darker the lines, the stronger the assignment.

## Appendix B

# FCM implementation with CUDA

CUDA è un'architettura hardware supportata dai processori grafici di *NVIDIA*, un'azienda statunitense. E' possibile realizzare il codice che definisce la comunicazione tra CPU e GPU in un dialetto del C++. Il codice che effettivamente esegue il processore grafico può essere già implementato con funzioni speciali delle librerie thrust e cuBLAS o realizzato personalmente attraverso dei kernel. In questa appendice si mostra il kernel usato per realizzare il calcolo della matrice  $U^2$  in FCM. In particolare l'algoritmo è ottimizzato per compiere il calcolo su al più MAX\_THREAD\_PER\_BLOCK centroidi, poiché nel medesimo

```

1  /**
2   * @brief This kernel computes the matrix U2 of membership between
3   * data points and centroids
4   *
5   * @param[in] d_data : the i-th is d_data[i * n_dimensions + k]
6   * for k = 0, ..., n_dimensions - 1
7   * @param[in] d_weights : the weight of the i-th data point is
8   * d_weights[i]
9   * @param[in] d_centroids : the j-th is
10  * d_centroids[j * n_dimensions + k] for k = 0, ..., n_dimensions - 1
11  * @param[out] d_matrix : the membership between the i-th data point
12  * and the j-th centroid is stored in d_matrix[i * n_centroids + j]
13  * @param n_data : number of data points
14  * @param n_dimensions : dimensions of data points
15  * @param n_centroids : number of centroids
16  *
17  * @details This kernel requires a grid of blocks with n_data blocks
18  * and MAX_THREADS_PER_BLOCK threads for each block.
19  *
20  * @note This kernel synchronize threads at the end of the computation
21  */
22 __global__ void
23 kernel_compute_U2 (const float *const d_data, const float *const d_weights,
24                     const float *const d_centroids, float *const d_matrix,
25                     size_t n_data, size_t n_dimensions, size_t n_centroids)
26 {
27     __shared__ float sdata[MAX_THREADS_PER_BLOCK];
28     size_t i = blockIdx.x; // i-th data
29     size_t j = threadIdx.x; // j-th centroid
30     float value = 0;
31     float reducted = 0;
32
33     // compute the distance between the i-th data point and the j-th

```

```

34 // centroid
35 if (i < n_data && j < n_centroids)
36 {
37     for (size_t k = 0; k < n_dimensions; k++)
38     {
39         float diff = d_data[i * n_dimensions + k]
40             - d_centroids[j * n_dimensions + k];
41         value += diff * diff;
42     }
43 }
44 // synchronize threads of this block
45 __syncthreads ();
46
47 // compute the min value of the block
48 if (j < n_centroids)
49     sdata[j] = value;
50 else
51     sdata[j] = FLT_MAX;
52 __syncthreads ();
53 for (size_t s = MAX_THREADS_PER_BLOCK / 2; s > 0; s >= 1)
54 {
55     if (j < s && sdata[j] > sdata[j + s])
56         sdata[j] = sdata[j + s];
57     __syncthreads ();
58 }
59 reducted = sdata[0];
60 // synchronize threads of this block
61 __syncthreads ();
62
63 // prepare the row to a stable normalization
64 if (reducted == 0.0)
65 {
66     // let to 1 the components that are 0 and to 0 the others
67     if (i < n_data && j < n_centroids)
68         value = value == 0.0 ? 1.0 : 0.0;
69 }
70 else
71 {
72     // for each component of the row, assign min/value
73     if (i < n_data && j < n_centroids)
74         value = reducted / value;
75 }
76 // synchronize threads of this block
77 __syncthreads ();
78
79 // compute the sum of the row
80 if (j < n_centroids)
81     sdata[j] = value;
82 else
83     sdata[j] = 0.0;
84 __syncthreads ();
85 for (size_t s = MAX_THREADS_PER_BLOCK / 2; s > 0; s >= 1)
86 {
87     if (j < s)
88         sdata[j] += sdata[j + s];
89     __syncthreads ();
90 }
91 reducted = sdata[0];
92 // synchronize threads of this block
93 __syncthreads ();
94
95 // assign the value to the matrix
96 if (i < n_data && j < n_centroids)
97     value /= reducted;
98     d_matrix[i * n_centroids + j] = value * value * d_weights[i];
99 // synchronize threads of this block

```

```
100     __syncthreads();
101 }
```

---



# Bibliography

- [1] Chiara Basile, Dario Benedetto, Emanuele Caglioti, and Mirko Degli Esposti. An example of mathematical authorship attribution. *Journal of Mathematical Physics*, 49, 12 2008. doi: 10.1063/1.2996507.
- [2] J. C. Dunn. A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters. *Journal of Cybernetics*, 3(3):32–57, 1973. doi: 10.1080/01969727308546046. URL <https://doi.org/10.1080/01969727308546046>.
- [3] I. B. HOPELY. Clerk maxwell’s experiments on colour vision. *Science Progress (1933- )*, 48(189):46–66, 1960. ISSN 00368504, 20477163. URL <http://www.jstor.org/stable/43424831>.
- [4] Stefano Magrini Alunno. Analisi automatica di disegni per attribuzione d’autore. Bachelor’s thesis, La Sapienza, March 2021. Available at <https://www.linkedin.com/feed/update/urn:li:activity:7182499324092694529/>.
- [5] Daniel Jurafsky & James H. Martin. An introduction to natural language processing, computational linguistics, and speech recognition. 2024. URL [https://web.stanford.edu/~jurafsky/slp3/ed3bookfeb3\\_2024.pdf](https://web.stanford.edu/~jurafsky/slp3/ed3bookfeb3_2024.pdf).
- [6] Joseph Nicéphore Niépce. View from the window at le gras. [https://upload.wikimedia.org/wikipedia/commons/5/5c/View\\_from\\_the\\_Window\\_at\\_Le\\_Gras%2C\\_Joseph\\_Nic%C3%A9phore\\_Ni%C3%A9pce.jpg](https://upload.wikimedia.org/wikipedia/commons/5/5c/View_from_the_Window_at_Le_Gras%2C_Joseph_Nic%C3%A9phore_Ni%C3%A9pce.jpg), c.1826. public domain.
- [7] Jean François WITZ. Illustrazione di una fotocamera digitale di tipo reflex. [https://upload.wikimedia.org/wikipedia/commons/3/31/Reflex\\_camera\\_numeric.svg](https://upload.wikimedia.org/wikipedia/commons/3/31/Reflex_camera_numeric.svg), 2008. CC BY-SA 3.0 Deed.