



universität  
wien

# MASTERARBEIT / MASTER'S THESIS

Titel der Masterarbeit / Title of the Master's Thesis

„Framework for Visual Cluster Analysis and  
Consensus Clustering“

verfasst von / submitted by  
Christian Permann, BSc.

angestrebter akademischer Grad / in partial fulfilment of the  
requirements for the degree  
MSc.

Wien, 2020 / Vienna, 2020  
Studienkennzahl lt. Studienblatt /  
degree programme code as it appears on  
the student record sheet

A 066 921

Studienrichtung lt. Studienblatt: /  
degree programme as it appears on  
the student record sheet

Masterstudium Scientific Computing

Betreut von / Supervisor:

Univ.-Prof. Dipl.-Inform.Univ. Dr. Claudia Plant



# Danksagungen

Ich möchte meiner Betreuerin Univ.-Prof. Dipl.-Inform.Univ. Dr.Claudia Plant für Ihre Unterstützung bei dieser Arbeit danken. Unsere offenen und konstruktiven Gespräche haben mir sehr geholfen diese Arbeit zielstrebig und motiviert zu vollenden.



# Contents

## I Introduction and Preliminaries

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Related Work</b>	<b>3</b>
2.1	Clustering Algorithms . . . . .	3
2.1.1	New Clustering Algorithms . . . . .	5
2.2	Visual Clustering Frameworks . . . . .	5
2.3	Consensus Clustering . . . . .	6
2.3.1	Object Co-occurrence . . . . .	6
2.3.2	Median Partitioning . . . . .	6
2.3.3	Consensus Clustering Frameworks . . . . .	6

## II Theory 7

<b>3</b>	<b>Methods</b>	<b>9</b>
3.1	Clustering . . . . .	9
3.1.1	OPTICS . . . . .	9
3.2	Consensus Clustering . . . . .	9
3.2.1	DICLENS . . . . .	9
3.3	Other Methods . . . . .	9
3.3.1	Hungarians Method . . . . .	9

## III Implemented Tool 11

<b>4</b>	<b>Using the Tool</b>	<b>13</b>
4.1	Data Visualization and Manipulation . . . . .	13
4.1.1	Plots and I/O . . . . .	13
4.1.2	Modifying the data . . . . .	16
4.2	Selection of Algorithms and Parameters . . . . .	17
4.3	Meta-View and Consensus Clustering . . . . .	17

<b>5</b>	<b>Implementation</b>	<b>19</b>
5.1	Programming Language and Tools . . . . .	19
5.2	File Input and Output . . . . .	19
5.2.1	Clustering and Point Data . . . . .	19
5.2.2	Clustering Workflows . . . . .	20
5.2.3	Pre-Clustered Data . . . . .	20
5.3	Methods and Interfaces . . . . .	20
5.3.1	Data Manipulation . . . . .	21
	Generators . . . . .	21
	Reducers . . . . .	21
	Normalizers . . . . .	22
5.3.2	Clustering Algorithms . . . . .	22
	Clustering Algorithms . . . . .	22
	ELKI Clustering Algorithms . . . . .	23
	Custom Clustering Algorithms . . . . .	23
5.3.3	Meta-Clustering . . . . .	23
	Clustering Distance Measures . . . . .	24
5.3.4	Consensus Clustering . . . . .	24
	Consensus Functions . . . . .	24
5.4	Visualization . . . . .	25
5.4.1	Scatter Plot . . . . .	25
5.4.2	Switching between Clusterings . . . . .	26
5.4.3	Scatter Plot Matrix . . . . .	26
5.4.4	OPTICS Plot . . . . .	26
5.4.5	Heat Map . . . . .	26
5.4.6	Filter Panel . . . . .	26
5.4.7	Cluster comparison . . . . .	26
<b>IV</b>	<b>Testing</b>	<b>27</b>
<b>6</b>	<b>Experiments</b>	<b>29</b>
6.1	Solutions better than Base-Clusterings . . . . .	29
6.2	Multiple Solutions . . . . .	29
<b>V</b>	<b>Concluding Thoughts</b>	<b>31</b>
<b>7</b>	<b>Future Work</b>	<b>33</b>
7.1	Improvements for Tool . . . . .	33
7.2	Research Consensus Clustering . . . . .	33
7.3	Visual Frameworks . . . . .	33

<b>8 Conclusion</b>	<b>35</b>
8.1 Lessons learned . . . . .	35
8.2 Reflection of Work . . . . .	35
<b>Appendices</b>	<b>37</b>
<b>A Abstract</b>	<b>39</b>
A.1 English abstract . . . . .	39
A.2 Deutsche Zusammenfassung . . . . .	40





# List of Figures

4.1	Data-View, the starting window of the application . . . . .	14
4.2	Import Dialog . . . . .	15
4.3	Data-View with data-set (including cluster labels) . . . . .	16



## List of Tables



## Part I

# Introduction and Preliminaries



# Chapter 1

## Introduction

The generation of clusterings is quite a difficult task when little knowledge on the data is available. There exist a large number of clustering algorithms that all aim to find a fitting partitioning of the data, but they all come with different assumptions on the data. Taking k-Means as an example, the data is assumed to be Gaussian distributed which might not actually be the case for many data-sets, making it unsuitable for clustering such data.

As a solution for this, visual frameworks have been proposed that aim to help finding the right clustering by evaluating the results of different algorithms and parameter settings via quality metrics and ranking the clustering solutions. One of these frameworks is Clustervision [20] whereby different algorithms are run, their solutions are visualized and ranked according to the average of five different quality metrics. The user is then able to look at the candidates through different visualizations and choose a solution as final result. A problem with this approach is though, that each metric is biased towards a specific resulting cluster structure, possibly ranking it higher than others even if it does not fit the data properly. Even the authors of Clustervision themselves mention that “the effectiveness of these metrics in gauging the quality of the clustering is also difficult to determine due to the lack of ground truth”. Also, as many solutions were calculated and only one of them is chosen in the final decision, a lot of calculated knowledge is ignored.

To overcome this problem, I propose to not use quality metrics for the choice of the best clustering, but to rely on robustness indicators for this selection. As robustness is difficult to evaluate by itself many different algorithms with different parameters should be used, analyzing where multiple methods or parameter sets find a similar result. Clusterings that are similar to many other results can be seen as robust and finding a consensus of those can result in a better overall clustering. To find such similar clusterings or groups of similar clusterings meta-clustering can be used. To then compute final candidates, each group of similar results can be merged to a consensus

result, capturing the essence of the group.

Based on this idea I created a tool which will be further described in Chapter 4. Additionally, the tool aims to include the expert user into the process of finding the result. Most research on consensus clustering, as also seen in the survey of Vega-Pons and Ruiz-Shulcloper [27], only briefly mention how generating or pre-selecting base clusterings can impact the result of consensus clustering, even though this impacts the quality of the final solution. For this reason, the tool allows the user to visually explore the clusterings created by the simple clustering algorithms, facilitating the choice of which ones to merge for computing a final result.

The rest of this thesis is organized as follows: [TEXT]



## Chapter 2

# Related Work

### 2.1 Clustering Algorithms

Clustering is the task of finding groups of points in data such that points within a group have some kind of relationship, while points in different groups do not relate. Defining how such a relationship can be expressed and how these groupings can be found are non-trivial tasks, for which reason a lot of research has been put towards this topic. Different ways of thinking about this problem, the expected results and the performance aspects of the algorithms resulted in a large number of different methods. To summarize which methods are available many surveys were written over time [5, 24, 29, 30]. Another aspect of clustering is which kinds of data can be clustered and even for more specific tasks like text clustering many methods are available [13], though these will not be mentioned further on as this paper only works with numeric data. Also, work has been put into optimizing clustering algorithms for different types of architectures and machines, e.g. for distributed systems [15] or optimized for energy consumption and load balancing [19]. Such aspects will also not be further considered as this paper and the created tool do not aim to beat existing methods in run-time or energy efficiency.

To give an overview over existing clustering algorithms Xu and Tian [30] generalize the idea of algorithms to nine categories. These categories are defined by which ideas the algorithms are based on, which are the following:

- Partition
- Hierarchy
- Fuzzy Theory
- Distribution
- Density

- Graph Theory
- Grid
- Fractal Theory
- Model

Partition based methods aim to find central points in the data to represent clusters. Data points are then assigned to the central point closest to them and each central point forms a cluster. The most famous method in this group is K-Means which has been implemented in different ways, an example for a newer way for this was proposed by Kanungo et al. [18].

Hierarchy based methods assume at first that each point represents a cluster and merges points to new clusters according to some logic. The opposite can also be done, first assuming all points belong to one cluster and then splitting them into multiple clusters. When starting with each point in one cluster typical linking strategies are single link, joining clusters according to the distance between the closest points, average link, joining according to the average distance of all points from one cluster to the other, or maximum link, whereby the maximum distance between points of the clusters is considered.

Fuzzy Theory based methods do unlike other methods not produce a discrete value of label ownership but a value in the interval  $[0, 1]$  representing how strongly a point belongs to a cluster. The usual method to compute such an ownership value is to optimize toward an objective function and seeing the values as likelihood of belonging to a given cluster.

Distribution based methods aim to find clusters with the idea that any cluster should have its own mathematical distribution. In other words, each distribution found in the data should have a cluster with its generated points assigned to it. In most cases Gaussian distributions are assumed, reducing the complexity of the search.

Density based methods assume that points with a dense neighborhood should belong to a common cluster. Typical parameters for such clustering algorithms include the size of the neighborhood to consider and the number of points that should be found in their, such that the points are considered densely connected. Example algorithms for this are DBSCAN [12] and OPTICS [6], which was introduced as an improvement of DBSCAN.

Graph Theory based methods work by regarding each point in the data as a vertex of a graph with their distances representing the edges. The algorithms then aim to partition the graph into  $k$  sub-graphs each representing a cluster, whereby the edges connecting sub-graphs should have a large sum (of the distances) or small sum (of similarities).

Grid based methods look at the data space as a multidimensional grid. The resulting rectangular spaces are then used in combination with hierarchical and density based methods to both improve the result and the computational performance of the algorithms.

Fractal Theory based methods aim to optimize the intrinsic quality of the fractal dimension [30]. In this context fractal represents the geometry divisible into multiple parts with common characters with the whole [30].

Model based methods work by first assuming a model for the data and then optimizing the clustering to reflect the choice of model. Many of those are based on statistical learning methods whereby a classification tree is built hierarchically under the assumption that the distribution for each attribute is independent [30].

### 2.1.1 New Clustering Algorithms

Newer methods for clustering have also been proposed, either with completely novel ideas or by combining ideas from the previous section. To briefly summarize [30], those algorithms are for example based on:

- Kernel
- Ensemble (more on this in Section 2.3)
- Swarm Intelligence
- Quantum Theory
- Spectral Graph Theory
- Affinity Propagation
- Density and Distance

Additionally algorithms were developed for special data structures [30]. These include:

- Spatial Data
- Data Streams
- Large-Scale Data

For a more in depth overview Xu and Tian [30] published a survey paper, on which the categories in this paper were based on.

## 2.2 Visual Clustering Frameworks

Visual clustering frameworks and tools aim to help users to perform clustering in a simpler way and to visualize results, such that they can make more educated decisions on which algorithms to use and to analyze the quality of their result.

Clustervision [20]

VISTA [9] ClusterMap and user result evaluation

iVIBRATE [10] extension VISTA?

simple visualizations ELKI implementation [25]

Frameworks with specific data-sets in mind:

Propose: “An Evolutionary and Visual Framework for Clustering of DNA Microarray Data” [8], tool for clustering DNA data with visualization

Even Memes [11]

In addition to visual frameworks for clustering, visual parameter space analysis is also relevant, as finding the best parameters for a clustering algorithm may be difficult. For such use-cases Sedlmair et al. [26] propose a framework which can help in finding parameters.

## 2.3 Consensus Clustering

Consensus clustering aims to combine multiple clustering solutions for the same data-set into one result. This result should be representative of the combined knowledge acquired by all base clusterings and therefor be more robust and accurate than the average result at least. As the consensus result should lie in the center of the solution space spanned by the base clusterings, if there are sufficiently many good clusterings it should be able to overcome problems arising from outliers. In literature the problem of consensus clustering is often found under different names, those being cluster ensemble [7], clustering aggregation [14] and ensemble of partitions [22].

Consensus clustering includes two parts. Firstly base clusterings must be generated, which can be done with different strategies [7]. One may choose to perform multiple runs of the same algorithm with the same parameters but different random seeds, the same algorithm with different parameters or even multiple different algorithms. All of these possibilities can be while always clustering the full data, a different sample of the data points (Bootstrapping) [23] each run (e.g. 80%) or even different subspaces of the data. After the generation of base clusterings, the next step is to combine those with a consensus function. In general, consensus functions can be separated into two groups, Object Co-occurrence (O-Co) methods and Median Partitioning (MP) methods [27]. There exist a large number of methods in both of those groups, which can be categorized according to what ideas they are based on.

### 2.3.1 Object Co-occurrence

### 2.3.2 Median Partitioning

### 2.3.3 Consensus Clustering Frameworks

consensus cluster plus [28]

# **Part II**

# **Theory**



## Chapter 3

# Methods

The newly created tool uses a large amount of existing methods which are described in this chapter. This description aims to give an overview of what exact methods were used, the theory behind them and in which scenarios they can be applied.

### 3.1 Clustering

The clustering algorithms implemented were used to create base clusterings for meta-clustering. The algorithms described in this section are the ones used for this step, except for OPTICS which was used for meta-clustering. For additional information on why OPTICS was used on the meta level, see Section [REFERENCE].

#### 3.1.1 OPTICS

### 3.2 Consensus Clustering

#### 3.2.1 DICLENS

### 3.3 Other Methods

#### 3.3.1 Hungarians Method





# **Part III**

## **Implemented Tool**



## Chapter 4

# Using the Tool

The implemented tool aims to facilitate the computation and selection of good clustering solutions for a given data-set. To do this both powerful visualizations and algorithms are needed, which are combined here. When working with the created application the user can first import and edit point data, then select clustering algorithms and parameters for meta-clustering and finally use the meta-view to analyze the results or create a consensus results. The final goal after all of this is either having found a clustering solution that is satisfactory or created one using consensus clustering in the final step. This simple workflow can be described with three steps, each with its own view and functionality. In the following sections the views and possibilities within each view are described. For more details on implementation details see Chapter 5

### 4.1 Data Visualization and Manipulation

To start with the tool, the user needs to import or create a data-set on which the clustering should be performed. As the data may have too many dimensions or include unnecessary columns, it needs to be sanitized as well. For this reason, the first window shown implements functionality for both of this steps. When the application is started a scatter plot without data is shown, as well as a side menu with different buttons. This starting window will further on also be called data-view as it visualizes point data, it can be seen in Figure 4.1.

#### 4.1.1 Plots and I/O

The scatter plot in the center left shows data that was imported or created and allows for changing the shown dimensions by clicking on the center of the corresponding axis. The range of values may also be adjusted this way, clicking on the upper or lower bound of the range on any axis opens a small

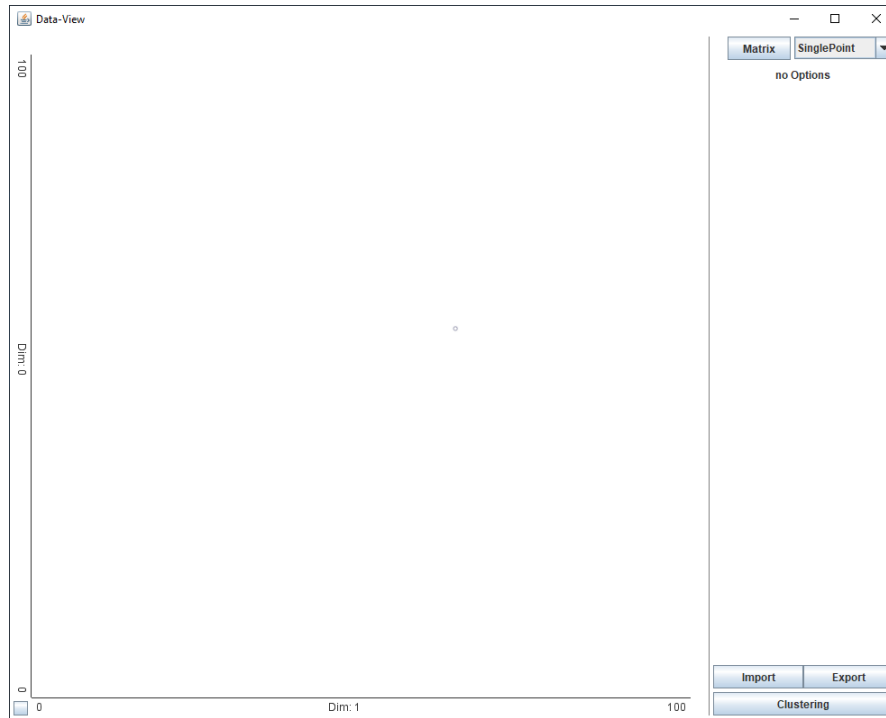


Figure 4.1: Data-View, the starting window of the application

input box. To automatically set this range, such that the minimum and maximum values of the shown data are just at the border, the small square button in the bottom left can be clicked. This is also a useful button to return to a full view of the data if the range was made smaller to have a zoomed view of the data. The right part of the tool contains the functionality for importing/exporting data, manipulating the data, showing it in a scatter plot matrix and opening the workflow-view which will be explained in the next step.

Firstly, to import a data-set for visualization and later use the import button in the bottom right can be selected, opening a file dialog which can load either CSV or ARFF files. These formats are supported as they are most commonly used to store such data. This file dialog can be seen in Figure 4.2.

After selecting a file the index of cluster labels can be defined, if labels should be loaded from the file. Otherwise the value can be left at  $-1$  indicating that no columns should be used for label information. “ $-2$ ” may also be specified, it is used as a shortcut for selecting the last column as cluster labels, removing the need to count columns if it is the last one. The chosen column is then not inserted as a dimension in the data space but used as the cluster labels, which the color of the loaded points will reflect in the scatter plot of the data-view.

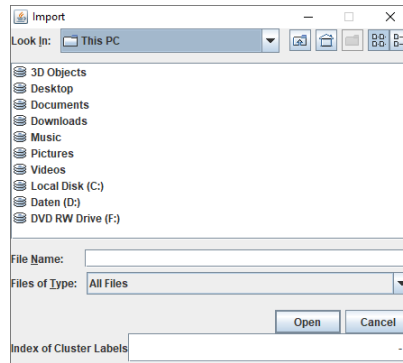


Figure 4.2: Import Dialog

To export data the export button can be used. It will open a similar dialog as for importing, though the index of the labels does not need to be defined, as data with cluster labels will be stored such that the labels for all points will be included in the first column. If the data does not contain label information, this first column is completely omitted. The supported format for exporting data is CSV.

Alternatively, the user can also create a data-set himself from within the tool. It is possible to draw points onto the canvas when the “SinglePoint” option is selected in the top right, as in Figure 4.1. For this, the user only needs to click wherever he wants points placed. It is important to note that if there are more than two dimensions all values for the not shown dimensions will not be set, which may cause a problem for the clustering algorithms later on. Another possibility is the use of the “ELKIGenerator”, which accepts a XML-style input for defining how to generate data. The definition for this format can be found in corresponding definitions file on the ELKI Github repository [1] or in the tools lib folder. This generator allows defining distributions on dimensions, the number of points for each generated distribution and many additional things like plane rotations. This can be quite useful when wanting to work with a quick sample data-set. In Figure 4.3 one can see how the data-view looks when data with cluster labels is shown.

In the top right the “Matrix” button can be used to display all dimensions at once, as it will open a new window showing the current data in a scatter plot matrix. Each scatter plot visualizes a different combination of dimensions and since the plots on the diagonal of this view are not useful, they were replaced with kernel density estimation plots. These kernel density estimation plots show how the values of the points in the corresponding dimension are distributed, overlapping the total distribution with the distribution of points with different labels. This may also allow seeing how clusters differ in each dimension.

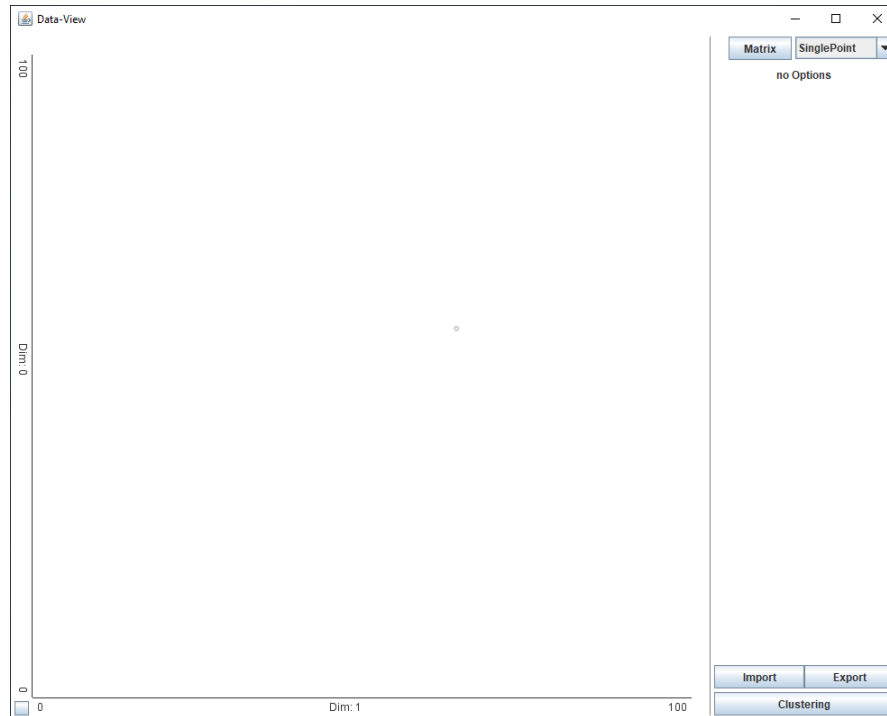


Figure 4.3: Data-View with data-set (including cluster labels)

### 4.1.2 Modifying the data

As the data may still contain unwanted dimensions after loading cleanup may need to be done. To reduce the dimensionality the simplest method is to select the “Dim. Remover” tool in the top right and define the index of the dimension to be removed. As more sophisticated methods PCA and t-SNE [CITATION] are also available. PCA calculates the principal components of the data and projects it onto the  $k$  strongest ones leaving us with  $k$  dimensions.  $k$  can easily be defined in the settings of PCA by filling in the dimensions field, as seen in Figure 4.3. t-SNE works similarly to PCA, has more options though and even a parallel implementation. For more information on those methods see Section [REFERENCE].

Other than reducing the dimensionality of the data the application allows to normalize the data, such that it can be more easily used for clustering. If there are different dimensions with strongly different meaning, unit and range this normalization step can benefit the accuracy of clustering greatly [CITATION?]. The available tools are “Normalize” and “Standardize”. Normalize changes the value range for each dimension to the interval  $[0, 1]$  while preserving relative distances, the Standardize tool adapts these ranges such that the average is 0 and the standard deviation is 1.

With these two pre-processing steps done, the actual clustering of the data

can be thought of next. To work on this the “Clustering” button opens the next major view of the tool.

## **4.2 Selection of Algorithms and Parameters**

By clicking the “Clustering” button in the data-view the window managing clustering workflows, further on called workflow-view, opens. Here clustering algorithms, their parameters and the parameters for meta-clustering can be selected. The workflow-view can be seen in Figure [REFERENCE]

## **4.3 Meta-View and Consensus Clustering**

bla





## Chapter 5

# Implementation

### 5.1 Programming Language and Tools

For the implementation of my tool my language of choice was Java. With Java there are lots of available implementations for different clustering methods available, including the Frameworks ELKI [4] and WEKA [16], as well as the machine learning library Smile [3]. To improve the performance in some parts of the tool Java 1.8 was chosen, making it the lowest required Java version to run it. This was done, as it enables the use of parallel streams which make it simple to parallelize functions. Additionally, it is easy to create custom visualizations using the low level `draw()` method of `JComponents` in Swing [2]. All graphs and plots were created using a custom `draw()` method as this leads to much better performance than using small components for the graphs.

### 5.2 File Input and Output

Loading and saving data is needed for the tool to make it more usable. For this reason it is possible to not only import data for clustering but also to save workflows and whole groups of clustering results for later reuse.

#### 5.2.1 Clustering and Point Data

Besides being able to create new data sets from within the tool, importing data from CSV or ARFF files is possible. Those file types are supported as they are very common, especially CSV is usually used to store such information. ARFF is commonly used, as it is also the preferred data type of WEKA [16]. For this reason the ARFF importer of WEKA was adapted to load in these files such that they can be used with this tool.

When importing data, it is also possible to define an index of where the cluster labels are located. When this is done the tool will not load the

specified column as data but use it to define clusters for the ground truth, which is then also handled differently by the different views of the tool.

The starting window or data-view, in which clusterings can also be loaded back into from the meta-view, also allows for exporting data to CSV. When doing this, the first column in the resulting file will contain the cluster identifiers starting with one. This way computed results can be exported and used in other tools or saved for later.

### 5.2.2 Clustering Workflows

As many runs of clustering algorithms with different parameter ranges can be performed and repeatedly defining this workflow can be tedious, the tool allows exporting and importing these workflows from the cluster workflow window. For this the underlying clustering objects are serialized and saved into .CWF (Cluster Workflow Files). It is important to note, that these files do not include the parameter settings for random seed or the meta-clustering. This is intentional as I do not consider them to be part of the workflow and they should be set by the user for every run, if the default values are not satisfactory.

### 5.2.3 Pre-Clustered Data

Due to the clustering algorithms possibly being time consuming, either because of their complexity or just sheer number of runs, the results of these algorithms can be saved from the meta-view. To do this, the clustering results are serialized and saved as .CRF (Clustering Result Files). To load those files one can define the different parameters of the meta-clustering and then import the .CRF. This way, only the meta-clustering is computed again. The meta-clustering is not saved in this file such that the user can more easily try different settings for it, by just saving the result and running the meta-clustering again. As this is usually much faster than the clustering algorithms that generate the base clusterings, using this may save a considerable amount of time.

## 5.3 Methods and Interfaces

To ensure that new or missing methods can easily be added to the tool, all methods where a selection can be made from within the tool can be added to by implementing onto the corresponding interface. To add methods one needs to first implement it and if needed it's OptionsPanel and then add it to the corresponding static method found in the view using it.

The following sections describe which interfaces and standard methods are implemented, as well as which functions and classes must be edited to add a new custom methods.

### 5.3.1 Data Manipulation

Data manipulation can be performed as a first step in the data-view to create data-sets or modify them such that they can be used for clustering or visualization. It is also of relevance when loading clustering results back into the data-view to analyze results using dimensionality reduction techniques. For this step three interfaces are of interest:

- Generators: IGenerator
- Reducers: IReducer
- Normalizers: INormalizer

#### Generators

Generators allow the user to add or replace data points. Two initial implementations for generators are available, firstly the SinglePointGenerator which allows adding points to the data by clicking the scatter plot. When a point is added this way its coordinates are set to where the user clicked on the screen translated into the coordinates of the data space, such that the point appears below the cursor. Secondly, the ELKIGenerator is available which uses an XML-style input, defining how the clusters should be generated. This class uses the GeneratorXMLDatabaseConnection of ELKI [4] and is able to create data points according to different random distributions. The definitions file for this input can be found in ELKI's JAR file, their GIT repository [1] or in the /lib folder of the implementation.

#### Reducers

Reducers allow the user to manipulate the data such that the dimensionality or potentially (in the future) the number of points can be reduced. Per default three dimensionality reduction techniques are available. The simplest one, DimensionRemover, enables the user to just define a dimension by its index and delete it. This can be useful, if the data is not well cleaned beforehand and irrelevant or possibly disrupting columns are present in the data. For more advanced dimensionality reduction, the PCAReducer and tSNEReducer are available. The PCAReducer, which is adapted from the implementation from Smiles [3] uses Principal Component Analysis [CITATION] to reduce the dimensionality of the data to what the user defined. In contrast the tSNEReducer adapted from Jonsson [17] uses t-Distributed Stochastic Neighbor Embedding [CITATION] and also requires a perplexity parameter.

## Normalizers

Normalizers provide functionality for normalizing the data points. This may be useful whenever different dimensions contain values in strongly differing ranges, while their concrete values do not have any direct relation to the other dimensions. Here, two default implementations are available, `Normalize` and `Standardize`. The `Normalize` class adapts the values such that in each dimension the range is set to the interval  $[0, 1]$  while preserving relative distances. The `Standardize` class adapts these ranges such that the average is 0 and the standard deviation is 1.

## Adding Data Manipulation Methods

To additional Methods they must be implemented using the corresponding interface and added to the class `DataView`. To add them to this class the static methods `initGenerators()`, `initReducers()` and `initNormalizers()` must add the new method to the list of methods. As an example adding a new generator (`newGeneratorX`) would require adding “`generators.add(new newGeneratorX());`” to the function `initGenerators()`.

### 5.3.2 Clustering Algorithms

Clustering algorithms can be selected and tuned in the clustering workflow window such that the algorithm is used to create base clustering results for the meta-view. For additional information on the algorithms themselves, see Section [REFERENCE] The following interfaces are relevant for clustering:

- Clustering Algorithms: `IClusterer`
- ELKI Clustering Algorithms: `IELKIClusterer`
- Custom Clustering Algorithms: `ICustomClusterer`

## Clustering Algorithms

Clustering algorithms can be used in the workflow-view to generate base clusterings for meta-clustering and visualization in the meta-view. The interface `IClusterer` is used as helper interface such that all clustering algorithms can use the same high level calls for compatibility. The actual implementations implement either `IELKIClusterer` if they use the ELKI database objects in their logic or `ICustomClusterer` if they use a data matrix. To further help not needing to re-implement generally needed methods the class `AbstractClustering` should also be extended by the actual implementation. With these class and interfaces the actual implementations can focus more on just the clustering logic. Additionally options panels should be created for each new algorithm allowing the user to set parameters or parameter ranges from the tool.

## ELKI Clustering Algorithms

IELKIClusterer is used for clustering algorithms that use ELKIs database object to store the data. To run algorithms from ELKI this interface must be implemented and the database object can be used with ELKIs algorithms. The pre-implemented methods include three variants of KMeans, DBScan and EM clustering. All of these methods use ELKIs high level initialization with `“ClassGenericsUtil.parameterizeOrAbort(AlgorithmClass.class, params);”` and other than this only require the logic for managing parameters, random sampling and collection of the result.

## Custom Clustering Algorithms

To enable other clustering algorithms than the ones implemented in ELKI the interface ICustomClusterer can be used. It receives the two-dimensional data array and the headers as input and allows for the use of any custom clustering logic. One implementation is already available here, namely Spectral Clustering which is implemented with the Smile library [3].

## Adding Clustering Algorithms

To add new clustering algorithms for creating base clusterings, first one must choose IELKIClusterer if the method uses ELKIs database object or ICustomClusterer otherwise. This interface then needs to be extended, the AbstractClustering class can help here with utility functions. When extending the corresponding interface an option panel must be implemented and the parameters must be managed within the cluster function. These interfaces also provide functionality for seeded randomness and progress indicators. To properly calculate progress, getCount() must return the proper number of needed executions and cluster() must call addProgress() whenever progress is made. Finally, the newly implemented class must be added to the initClusterers() function in ClusterWorkflow.

### 5.3.3 Meta-Clustering

The algorithm used for meta-clustering is OPTICS [CITATION]. One of the requirements for this method is that the used distance measure must be metric, otherwise the algorithm can only produce an approximate clustering at best. The distance measure used by the algorithm can be changed and additional distance measures can also be added. For this, the following interface is of interest:

- Clustering Distance Measure: IMetaDistanceMeasure

## Clustering Distance Measures

The interface `IMetaDistanceMeasure` is used to define distance measures that can be used for OPTICS for meta-clustering. These distance measures should represent how different two clustering solutions are from each other while needing to fulfill the requirements of a metric. For a distance measure, here indicated with  $d(x, y)$  where  $x$  and  $y$  are objects, to be considered metric the following conditions must hold:

- $d(x, y) \geq 0$ : distances must be non-negative
- $d(x, y) = 0$  iff  $x = y$ : the distance between two objects is zero if and only if both objects are equal
- $d(x, y) = d(y, x)$ : the distance must fulfill the symmetry condition
- $d(x, z) \leq d(x, y) + d(y, z)$ : the triangle inequality must hold

The default implemented distance measures are Variation of Information and Clustering Error.

## Adding Distance Measures

To add distance measures for OPTICS the interface `IMetaDistanceMeasure` must be extended and the above mentioned conditions on the distance measure must be fulfilled for the algorithm to produce a satisfying result. After implementing the distance it must be added to `initDistances()` within `ClusterWorkflow`, for it to be usable from within the tool.

### 5.3.4 Consensus Clustering

Consensus Clustering can combine the results from different base clusterings to a common clustering. This can be done in the meta-view and different algorithms can be used for this step. When run, the algorithm computes a consensus result for all base clusterings that are in a common meta-cluster in the OPTICS reachability plot and are not filtered out. This means that if the used defined three meta-clusters with some cut-off value in the reachability plot, there will be three resulting consensus clusterings, one created from each group. The interface of interest for consensus functions is:

- Consensus Functions: `IConsensusFunction`

## Consensus Functions

The consensus functions implement the logic for combining the base clusterings. For this the interface `IConsensusFunction` is of interest and can be extended to add new functions. Per default two main implementations are

available. Both implementations allow for the user to either define a target value for the number of resulting clusters  $k$  within each result, or for the algorithm to guess this parameter itself. Firstly, `CoAssociationMatrixAverageLink` can be used, which computes a Co-Association (CA) Matrix and combines objects via the average link strategy. When  $k$  is defined the algorithm joins together sets of points until only  $k$  sets remain, which is used as final answer. If no  $k$  is defined the algorithm performs two runs, calculating the lifetime [31] of the resulting tree and choosing the level to cut as the one resulting with maximum lifetime. As second method `DICLENS` is available, which was generously provided by Mimaroglu and Aksehirli [21].

### Adding Consensus Functions

To add consensus functions the interface `IConsensusFunction` must be extended. It is also important to note that weights may be supported by the function if implemented, although at this point there is no way for the user to set weights. As for the two default implementations the function may either allow defining the number of resulting clusters per consensus result or not. To add the newly created function to the tool, within `MetaViewer` the function `initConsensusFunctions()` must include the new method.

## 5.4 Visualization

To visualize the data in different ways and provide a functional Graphical User Interface (GUI) the Java Swing [2] Toolkit was used. It allows for both high level calls creating objects within windows, e.g. panels, buttons or combo-boxes, as well as low level calls by directly changing the behavior of the `draw()` function. For the basic GUI the panels and buttons of the tool were created with high level calls, while the different graphs and plots were implemented by overriding the `draw()` function. Overriding the `draw()` function immensely improved performance as the already available high level objects are computationally costly by comparison, when many of them are needed. This section only includes additional information on the implementation and functionality, information on the usage of the visual elements in the tool can be found in Section [REFERENCE].

### 5.4.1 Scatter Plot

The scatter plot is the simplest visualization of the data within the tool. It is implemented through two axis and a canvas object. The axis include logic for the draw the value range and selected dimension, including the mouse listeners for modifying it, as well as the translation from model coordinates to viewer coordinates. The canvas implements the drawing of data points, including different shapes and colors depending on their cluster identifier,

whether or not they are filtered out or if they represent special objects like the ground truth in the meta-view. Additionally the canvas can support mouse listeners enabling interactivity with the points.

#### **5.4.2 Switching between Clusterings**

bla bla bla

#### **5.4.3 Scatter Plot Matrix**

bla bla bla

#### **5.4.4 OPTICS Plot**

bla bla bla

#### **5.4.5 Heat Map**

bla bla bla

#### **5.4.6 Filter Panel**

bla bla bla

#### **5.4.7 Cluster comparison**

bla bla bla



# Part IV

# Testing



## Chapter 6

# Experiments

To show the value of the new tool different scenarios created. The goals of the tool are to improve on the ease of finding complex cluster structures and improve on the result in comparison to only using simple clustering methods. Another aim is to visualize the data and clustering results such that an expert user can make a well informed decision on which solutions best fit the data.

### 6.1 Solutions better than Base-Clusterings

bla

### 6.2 Multiple Solutions

bla



**Part V**

**Concluding Thoughts**



## Chapter 7

# Future Work

### 7.1 Improvements for Tool

bla

### 7.2 Research Consensus Clustering

bla

### 7.3 Visual Frameworks

bla





## Chapter 8

# Conclusion

### 8.1 Lessons learned

bla

### 8.2 Reflection of Work

In this article I showed



# Appendices



# Appendix A

## Abstract

### A.1 English abstract

Finding a good clustering solution for an unexplored data-set is a non-trivial task. Due to the large number of clustering algorithms which usually have lots of parameters, clustering results may differ strongly from each other and the underlying ground truth. With only little knowledge on the data the evaluation of which result best represents the underlying cluster structure is difficult. To find a fitting selection for the result, there exist visual frameworks that aim to simplify this choice by ranking the results according to quality measures. As those measures also have the downside of being biased towards specific structures (whether or not they fit the data) they are problematic for selecting a final result. For this reason, I propose to purely use indicators of robustness for the creation of a clustering result. This is done by meta-clustering results from different clustering algorithms and results and calculating consensus clusterings from each group of similar results. Additionally this process is supported through visualizations, giving the expert user the possibility to use his knowledge to further improve on the final result.

## A.2 Deutsche Zusammenfassung

Eine gute Clustering Lösung für wenig erforschte Daten zu finden ist eine komplexe Aufgabe. Wegen der großen Anzahl an Clustering Algorithmen, welche meist auch viele verschiedene Parameter benötigen, können sich die Ergebnisse stark untereinander, aber auch von dem richtigen Ergebnis, unterscheiden. Mit nur wenig Wissen über die Daten ist auch die Evaluierung welches Ergebnis am nächsten zu der der unterliegenden Wahrheit, beziehungsweise am besten der Struktur der Daten entspricht eine schwere Aufgabe. Um eine solche Auswahl besser treffen zu können wurden visuelle Frameworks erschaffen, die mittels Qualitäts-Metriken die verschiedenen Ergebnisse bewerten und gereiht anzeigen. Da diese Metriken aber auch das Problem haben gewisse Strukturen in Ergebnissen zu bevorzugen zeigen sie sich wiederum bei der Entscheidung über das endgültige Ergebnis als problematisch. Aus diesem Grund schlage ich vor die Eigenschaft wie robust ein Ergebnis ist für die finale Entscheidung heranzuziehen. Um dies zu tun werden die Clusterings auf Meta-Ebene nochmals geclustert, wobei ähnliche Ergebnisse in einer Gruppe mittels Consensus Clustering zu einer Lösung zusammengeführt werden. Dieser Prozess wird weiters durch Visualisierungen unterstützt, so dass ein Experte mit Hilfe seines Wissens die Lösung möglicher Weise noch weiter verbessern kann.

# Bibliography

- [1] ELKI xml generator definitions file. <https://github.com/gaoch023/ELKI/blob/master/src/main/java/de/lmu/ifi/dbs/elki/application/GeneratorByModel.xsd>. Accessed: 2020-01-30.
- [2] Java Swing package definition. <https://docs.oracle.com/javase/8/docs/api/index.html?javax/swing/package-summary.html>. Accessed: 2020-02-01.
- [3] Smile - statistical machine intelligence and learning engine. <http://haifengl.github.io/>. Accessed: 2020-01-30.
- [4] ACHTERT, E., KRIEGEL, H.-P., AND ZIMEK, A. Elki: A software system for evaluation of subspace clustering algorithms. In *Proceedings of the 20th International Conference on Scientific and Statistical Database Management* (Berlin, Heidelberg, 2008), SSDBM '08, Springer-Verlag, p. 580–585.
- [5] AHALYA, G., AND PANDEY, H. M. Data clustering approaches survey and analysis. In *2015 International Conference on Futuristic Trends on Computational Analysis and Knowledge Management (ABLAZE)* (Feb 2015), pp. 532–537.
- [6] ANKERST, M., BREUNIG, M. M., KRIEGEL, H.-P., AND SANDER, J. Optics: Ordering points to identify the clustering structure. *SIGMOD Rec.* 28, 2 (June 1999), 49–60.
- [7] BOONGOEN, T., AND IAM-ON, N. Cluster ensembles: A survey of approaches with recent extensions and applications. *Computer Science Review* 28 (2018), 1 – 25.
- [8] CASTELLANOS-GARZÓN, J., AND DÍAZ, F. An evolutionary and visual framework for clustering of dna microarray data. *Journal of Integrative Bioinformatics* 10 (12 2013).
- [9] CHEN, K., AND LIU, L. A visual framework invites human into clustering process. pp. 97 – 106.

- [10] CHEN, K., AND LIU, L. Ivibrate: Interactive visualization-based framework for clustering large datasets. *ACM Trans. Inf. Syst.* 24, 2 (Apr. 2006), 245–294.
- [11] DANG, A., MOH'D, A., GRUZD, A., MILIOS, E., AND MINGHIM, R. *An Offline–Online Visual Framework for Clustering Memes in Social Media*. Springer International Publishing, Cham, 2017, pp. 1–29.
- [12] ESTER, M., KRIEGEL, H.-P., SANDER, J., AND XU, X. A density-based algorithm for discovering clusters a density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining* (1996), KDD'96, AAAI Press, p. 226–231.
- [13] FASHENG LIU, AND LU XIONG. Survey on text clustering algorithm-research present situation of text clustering algorithm. In *2011 IEEE 2nd International Conference on Software Engineering and Service Science* (July 2011), pp. 196–199.
- [14] GIONIS, A., MANNILA, H., AND TSAPARAS, P. Clustering aggregation. *21st International Conference on Data Engineering (ICDE'05)* (2005), 341–352.
- [15] HAI, M., ZHANG, S., ZHU, L., AND WANG, Y. A survey of distributed clustering algorithms. In *2012 International Conference on Industrial Control and Electronics Engineering* (Aug 2012), pp. 1142–1145.
- [16] HALL, M., FRANK, E., HOLMES, G., PFAHRINGER, B., REUTEMANN, P., AND WITTEN, I. H. The weka data mining software: An update. *SIGKDD Explor. Newsl.* 11, 1 (Nov. 2009), 10–18.
- [17] JONSSON, L. t-SNE Java implementation. <https://github.com/lejon/T-SNE-Java>. Accessed: 2020-02-01.
- [18] KANUNGO, T., MOUNT, D. M., NETANYAHU, N. S., PIATKO, C. D., SILVERMAN, R., AND WU, A. Y. An efficient k-means clustering algorithm: analysis and implementation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 24, 7 (July 2002), 881–892.
- [19] KRISHNAKUMAR, A., AND ANURATHA, V. Survey on energy efficient load-balanced clustering algorithm based on variable convergence time for wireless sensor networks. In *2016 3rd International Conference on Advanced Computing and Communication Systems (ICACCS)* (Jan 2016), vol. 01, pp. 1–5.
- [20] KWON, B. C., EYSENBACH, B., VERMA, J., NG, K., DEFILIPPI, C., STEWART, W. F., AND PERER, A. Clustervision: Visual supervision



- of unsupervised clustering. *IEEE Transactions on Visualization and Computer Graphics* 24 (2018), 142–151.
- [21] MIMAROGLU, S., AND AKSEHIRLI, E. Diclens: Divisive clustering ensemble with automatic cluster number. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 9, 2 (March 2012), 408–420.
  - [22] MINAEI, B., TOPCHY, A., AND PUNCH, W. Ensembles of partitions via data resampling. pp. 188–192.
  - [23] MONTI, S., TAMAYO, P., AND MESIROV, J. Consensus clustering: A resampling-based method for class discovery and visualization of gene expression microarray data. *Machine Learning* 52 (07 2003), 91–118.
  - [24] RUI XU, AND WUNSCH, D. Survey of clustering algorithms. *IEEE Transactions on Neural Networks* 16, 3 (May 2005), 645–678.
  - [25] SCHUBERT, E., KOOS, A., EMRICH, T., ZÜFLE, A., SCHMID, K. A., AND ZIMEK, A. A framework for clustering uncertain data. *Proc. VLDB Endow.* 8, 12 (Aug. 2015), 1976–1979.
  - [26] SEDLMAIR, M., HEINZL, C., BRUCKNER, S., PIRINGER, H., AND MÖLLER, T. Visual parameter space analysis: A conceptual framework. *IEEE Transactions on Visualization and Computer Graphics* 20, 12 (Dec 2014), 2161–2170.
  - [27] VEGA-PONS, S., AND RUIZ-SHULCLOPER, J. A survey of clustering ensemble algorithms. *International Journal of Pattern Recognition and Artificial Intelligence* 25 (May 2011), 337–372.
  - [28] WILKERSON, M. D., AND HAYES, D. N. ConsensusClusterPlus: a class discovery tool with confidence assessments and item tracking. *Bioinformatics* 26, 12 (04 2010), 1572–1573.
  - [29] WONG, K. A short survey on data clustering algorithms. In *2015 Second International Conference on Soft Computing and Machine Intelligence (ISCMCI)* (Nov 2015), pp. 64–68.
  - [30] XU, D., AND TIAN, Y. A comprehensive survey of clustering algorithms. *Annals of Data Science* 2 (08 2015).
  - [31] YANG, Y. Chapter 4 - ensemble learning. In *Temporal Data Mining Via Unsupervised Ensemble Learning*, Y. Yang, Ed. Elsevier, 2017, pp. 35 – 56.