# WATER POLLUTION ANALYSIS

Group-12
NaingSoe (13960333)
Tuan Dat Tran (13977666)
Aung Kaung Myat (13727540)
Khun Htun Myint Myat (13921403)

# Abstract

Data mining is an area of computer science that is multidisciplinary. It is a computational technique that finds patterns in relatively big data sets by combining artificial intelligence, machine learning, statistics, and databases. And it's been used to sift through massive volumes of data in order to uncover business insight that can help people solve problems, mitigate risks, and capture new possibilities(Investopedia, n. d).

The purpose of this report is to provide our findings, recommendations and solutions to relevant water quality research institutions through data mining. Provide effective solutions and accurate statistics for severely polluted areas.

Our team will use Wake data mining software for data description, data cleaning, data transformation, data reduction, data modelling and other methods to conduct in-depth research on data. Through k-MEAN and DBSCAN studies, we find that water pollution is closely related to national population density. The higher the population density, the more serious the water pollution. To understand the impact and harm of water pollution on our daily life, and put forward our team's suggestions and methods for the areas with serious water pollution. Because we cherish the principle of protecting the ecology and maintaining the healthy life of the river.

The dataset: https://www.kaggle.com/ozgurdogan646/water-quality-dataset

# **Table of Contents**

# 1   Introduction

Water pollution is when a water source is contaminated with substances that make it impossible to use for drinking, cooking, cleaning, swimming and other activities. Pollutants include chemicals, garbage, bacteria and parasites. All forms of pollution end up in the water (HSPH, n. d).

First, we choose an area with water pollution problem, and we have to solve and dig into this problem. Before deciding to study water quality analysis, we must find suitable data sets to select water quality data sets, because some data do not have relevant information, such as the area of water source and the time period of analysis. In the whole report, we will use data description, data cleansing, data transformation, data reduction, data modelling and other methods to conduct in-depth research on data step by step.

According to these methods, K-MEAN and DBSCAN can help us to know the areas and regions with serious water pollution more quickly and better. Results and charts for critical relational water quality analysis will also be obtained. During the modelling phase, we used the country and population density of the dataset to map areas with high water pollution. So we can categorize this data to change the local environment. The report will also be broken down by water pollution type and severity, so that local governments can provide information when improving the environment.

## 1.1   *Business Scenario*

Our business vision is to use the data mining services we provide to solve local water pollution problems. Because we know the harm and severity of water pollution, and water pollution is not easy to control, which requires a lot of manpower, material resources and time to deal with, our team hopes to provide useful data to relevant organizations through data mining and improve efficiency. In this case, we worked with the relevant surface water research team, which provided us with 20,000 water records and 29 data sets of water quality attributes. With these data, our team will use Weka data mining software to further understand the relationship between these water resources and give some suggestions and solutions based on our results.

# 2   Method

## 2.1   *Min-Max Normalization*

Min-max normalization is one of the most common methods of data normalization. For each function, the minimum value of the function is converted to 0, the maximum value to 1, and each other value is converted to a decimal between 0 and 1. Because it is an essential step in the feature processing process. Data standardization is to eliminate the dimensional influence of different indicators and facilitate the comparability between indicators. Dimensional differences will affect the results of distance calculation in some models （Codecademy, n. d).

## 2.2   *k - means*

K-means clustering is a basic unsupervised learning method. It works based on a simple algorithm that divides a given data set into clusters represented by the letter "K," which is predetermined. In this report, we use K-means to position clustering as points, all observation points or data points are associated with the most recent clustering, calculate and adjust, and then repeat the process with new adjustments until the desired results are obtained (Techopedia, n. d).

## 2.3 *DBSCAN*

DBSCAN is a density-based clustering technique that works on the premise that low-density regions in space separate dense regions in space. In "dense grouping," it groups data points into a single cluster. By examining the local density of data points, it can detect clusters in big spatial data sets. The DBSCAN cluster's most fascinating characteristic is that it excels at detecting outliers. Therefore, it helps us to find noise data and analyze water pollution (Analytics Vidhya, 2020).

## 3 Data Pre-processing and Description

### 3.1 *Dataset Used for the Report*

The data set is based on water-based water quality data shared on EAA and edited, deleted and optimized by Kaggle. Source：https://www.kaggle.com/ozgurdogan646/water-quality-dataset

### 3.2 Dataset Description

The new dataset contains a large amount of data on water-based water quality, including 20,000 records and 29 attributes. The data set mainly analyzes three categories of water quality, namely River water(RW),Ground water(GW) and Lake water(LW). The other detailed properties are shown in Figure 1 (Kaggle, n. d).

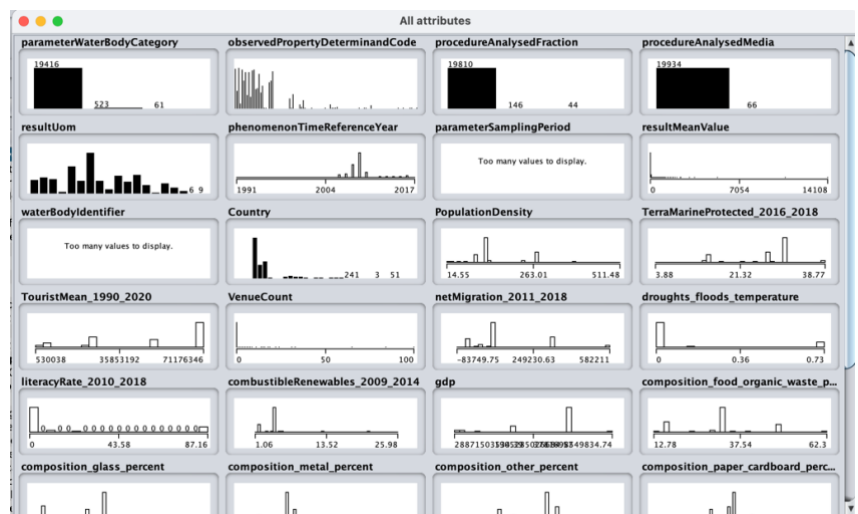| Attribute | Various types |
|---|---|
| Parameter Water Body Category | Nominal |
| Observed Property Determinand Code | Nominal |
| Procedure Analysed Fraction | Nominal |
| Procedure Analysed Media | Nominal |
| Result Uom | Nominal |
| Phenomenon Time Reference Year | Numeric |
| Parameter Sampling Period | Nominal |
| Result Mean Value | Numeric |
| Water Body Identifier | Nominal |
| Country | Nominal |
| Population Density | Numeric |
| Terra Marine Protected_2016_2018 | Numeric |
| TouristMean_1990_2020 | Numeric |
| Venue Count | Numeric |
| netMigration_2011_2018 | Numeric |
| Droughts floods temperature | Numeric |
| literacyRate_2010_2018 | Numeric |
| combustibleRenewables_2009_2014 | Numeric |
| gdp | Numeric |
| Composition food organic waste percent | Numeric |
| Composition glass percent | Numeric |
| Composition metal percent | Numeric |
| Composition other percent | Numeric |
| Composition paper cardboard percent | Numeric |
| Composition plastic percent | Numeric |
| Composition rubber leather percent | Numeric |
| Composition wood percent | Numeric |
| Composition yard garden green waste percent | Numeric |
| Waste treatment recycling percent | Numeric |



Figure-1

*3.3 Data Cleaning*

Once we get the initial data set, we need to edit the data set with 20,000 records and 29 attributes. Because there are missing values and duplicate values in the data set that need to be deleted and optimized. We use Weka for this series of processes.

First, we load the 'CSV' file into Weka and store it as' arff 'for our future identification, checking each attribute and eliminating missing values. We can observe that many attributes have 107 (1%) missing values, such as attribute 11, as shown in Figure 2. To maximize the accuracy of the study results, we used unsupervised filters named " Remove with Value " and "Remove Duplicates" in Weka, which delete missing values and duplicate values from the data set. After deletion, 19,814  records remain. See Figure 3 and 4.
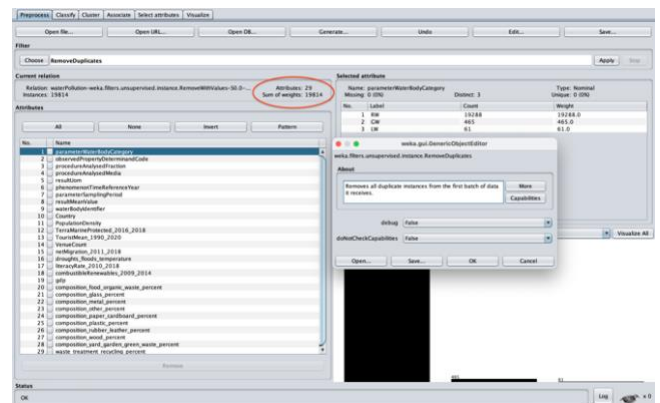


Figure-2



Figure-3



Figure-4

After editing, we deleted some unnecessary attributes and selected 15 attributes for in-depth research because many attributes were not very relevant to what we wanted to study, and there were also many values missing in it. See Figure 5. To confirm that the current data set is complete, unsupervised filters of "Remove Duplicates" for 15 attributes were implemented again, resulting in a reduction of 2,359 incomplete records as shown in Figure 6.
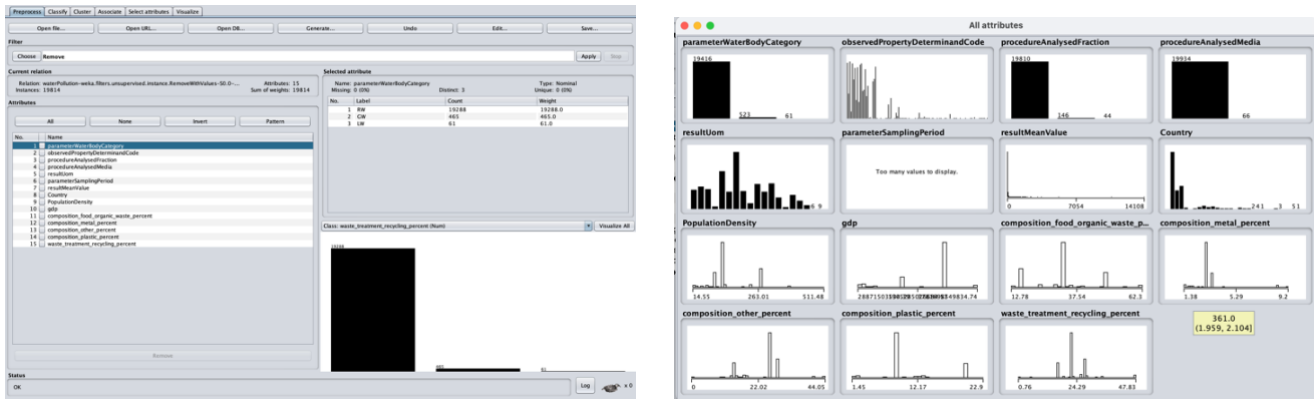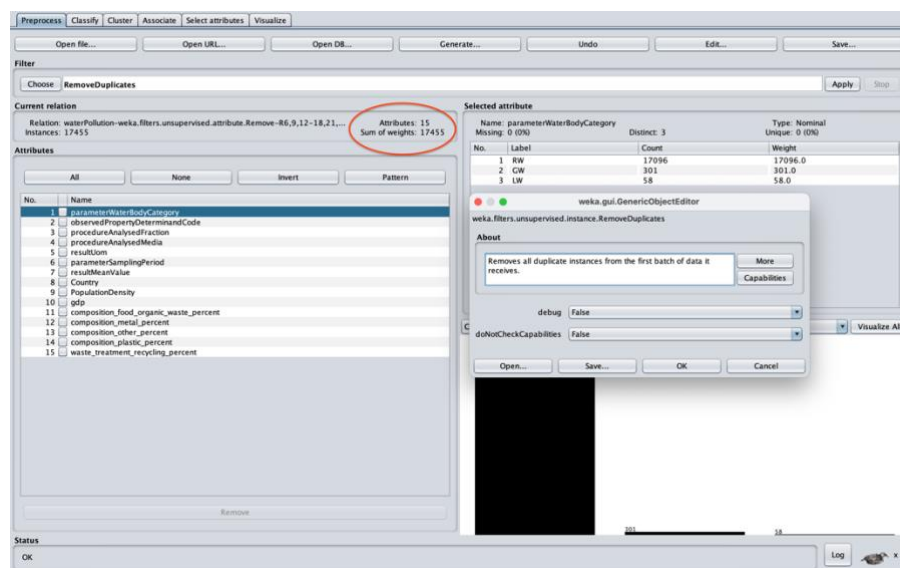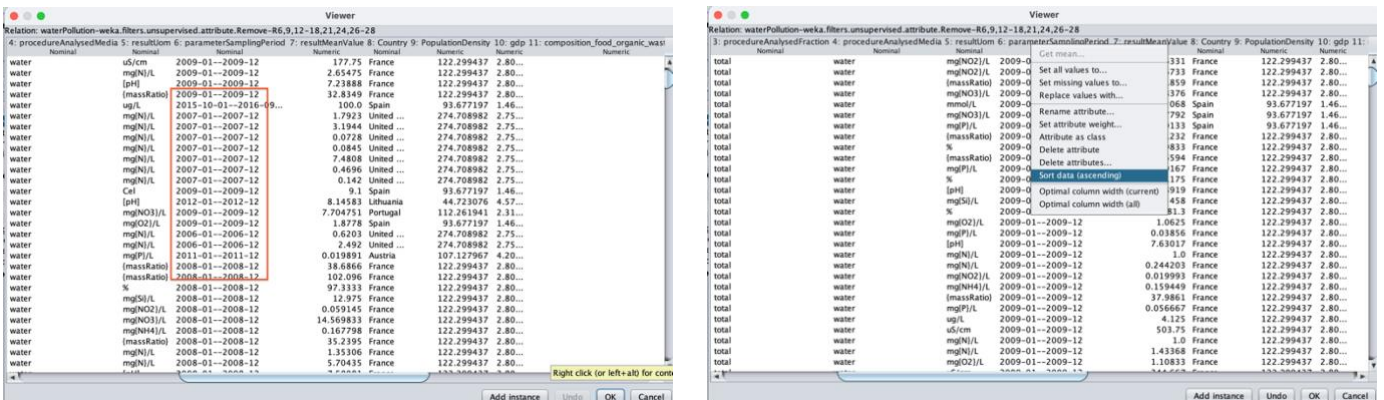
Figure-5



Figure-6

## 3.4 Data Transformation

After the data set is cleaned, we can clearly see that the date of the parameter sampling period is irregular. So we sort the data in 'Edit' by ascending the parameter sampling period. The min-Max Normalization approach is then used to narrow down the data to the specified extent. See Figure 7
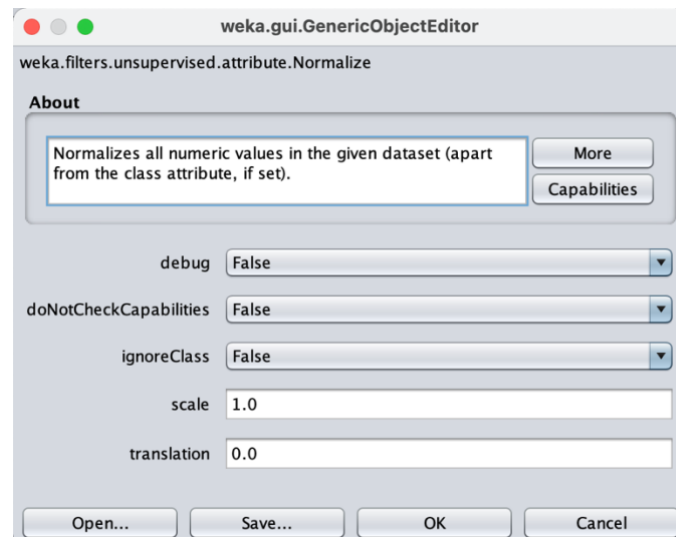
Figure-7

*3.5   Data Reduction*

There is a large amount of data in the dataset, so in order to get an accurate result, we decided to focus on the 2010 data and delete other dates from the new dataset.

First, we open the 'Edit' screen, select the data for 2010 from the 'Parameter Sampling period' property, and delete the data for other dates. And save the file again as another 'arff' file. Then delete the 'parameter sampling period' to make the data clean. See Figure 8.
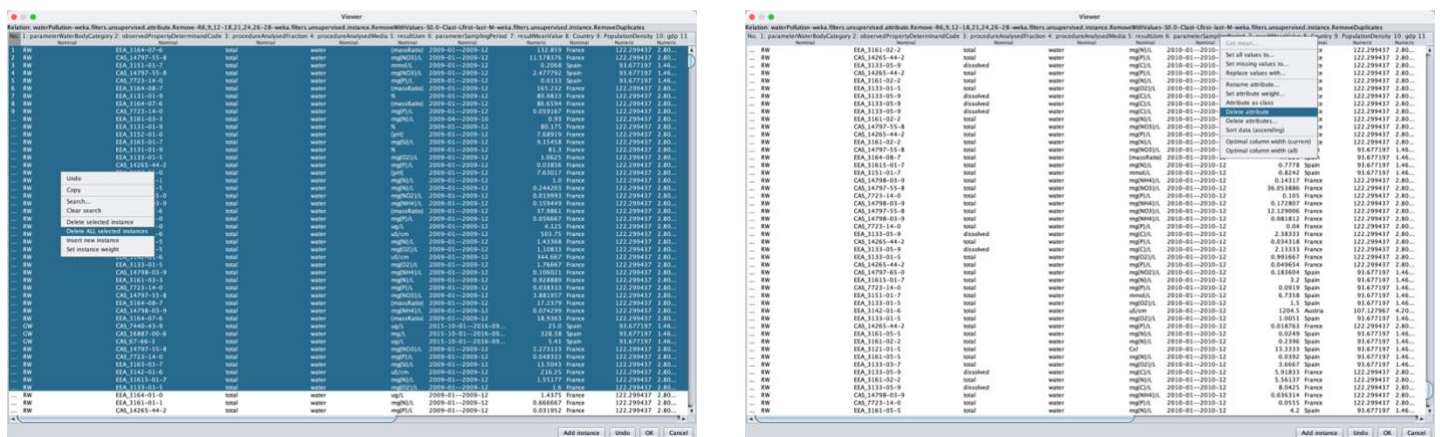


Figure-8

After selecting the 2010 water quality sampling survey data, there are more than 1500 records. The data set is still difficult and complex to study because the survey subjects are so large, so we randomly selected 900 data from it. Name it 'sample.arff' and store it. See figure-9.
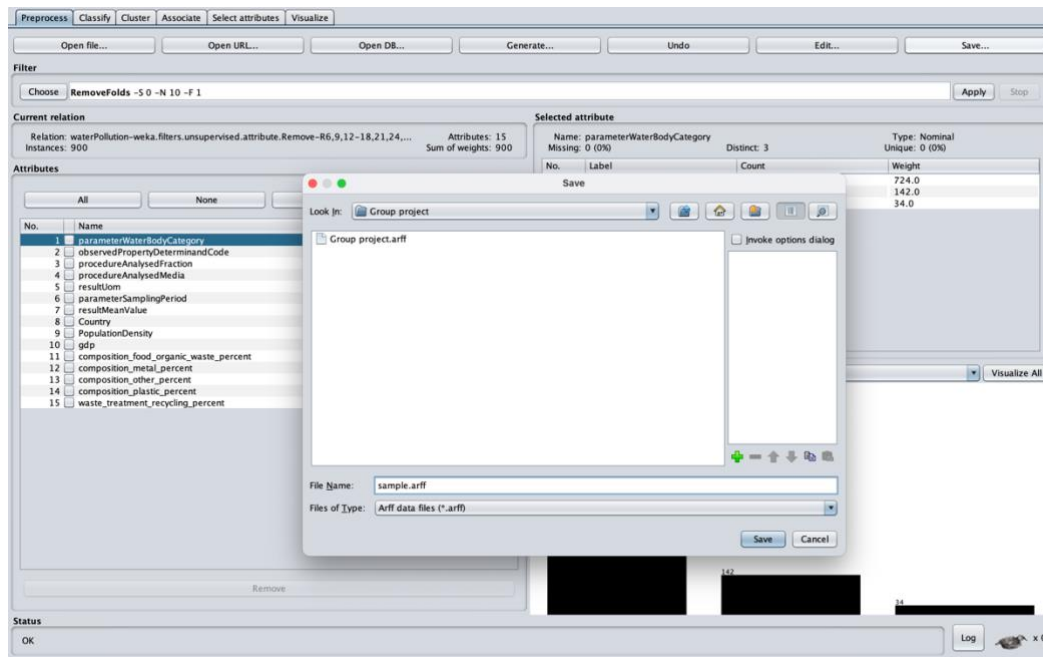
Figure-9

As the data for these columns named 'Procedure Analysed Fraction', 'Procedure Analysed Media' and 'Result Uom' are not relevant to our target, they are deleted. See Figure-10.
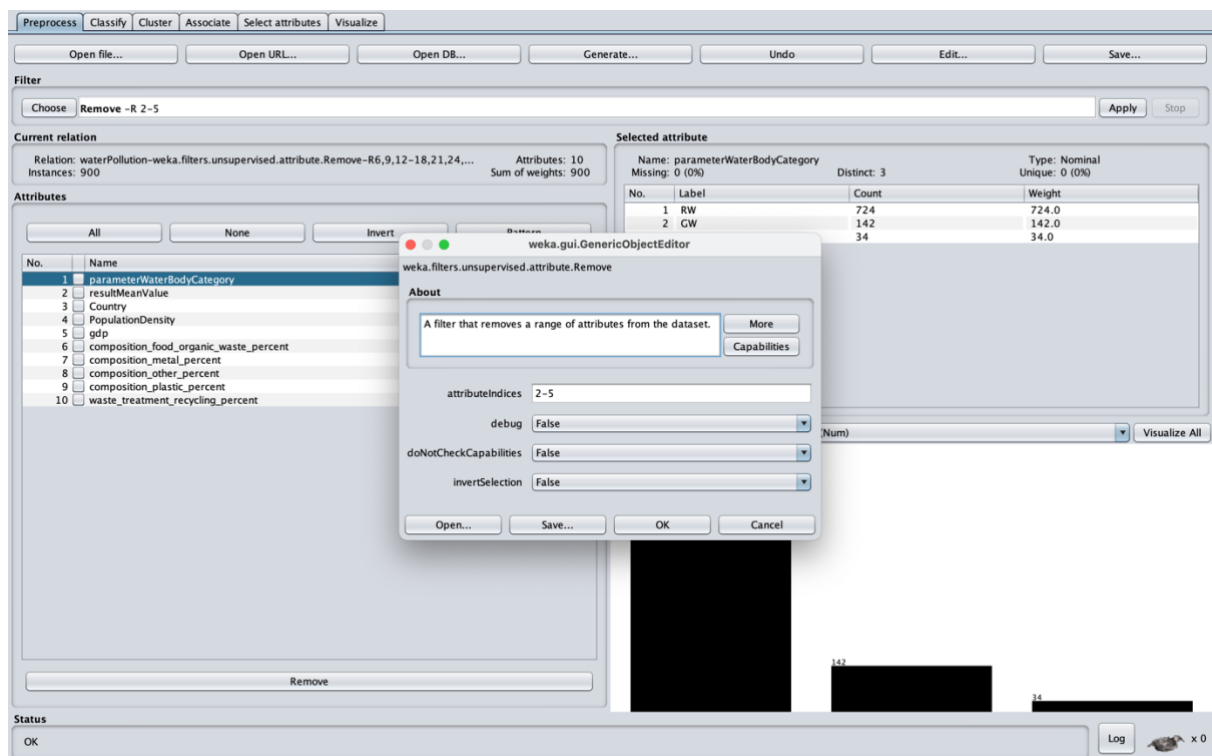


Figure-10

## 4 Data Modelling

### 4.1 Partitioning Clustering Approaches

In order to find out the serious water pollution areas in Europe in a better and faster way, we decided to use k-mean and DBSCAN to do the research. However, in order to make it more specific, we need to use k-mean algorithm to divide the clustering into 6 groups with 123 seeds. As show in figure 11.

### 4.1.1 K-means clustering

As can be seen from figure 12, among the six clusters, cluster 2 is the highest, accounting for 36%, followed by cluster 4, accounting for 31%, followed by cluster 0,3,1 and 5, accounting for 33%. This suggests that water pollution in the United Kingdom and Spain is more serious than in other countries and recommends special action to address water pollution in these waters.
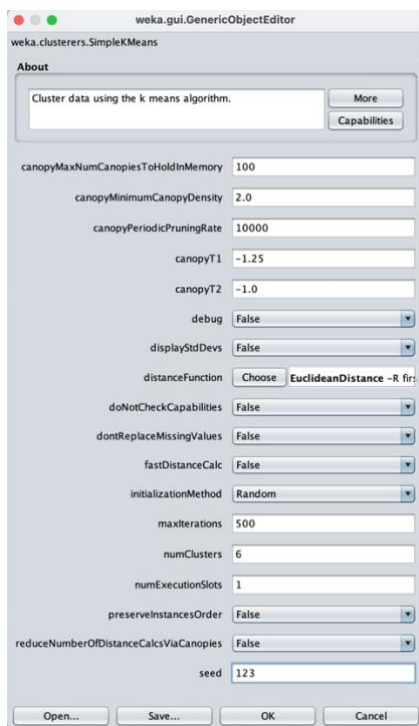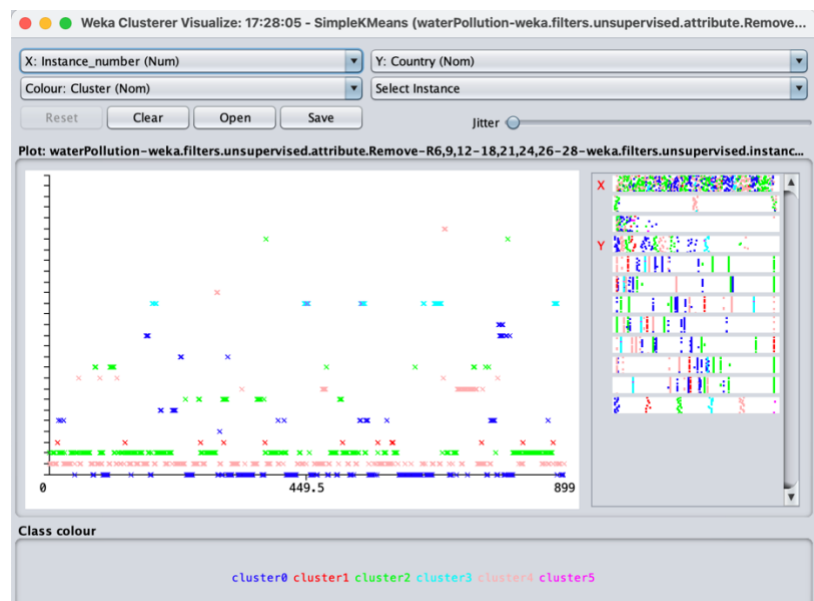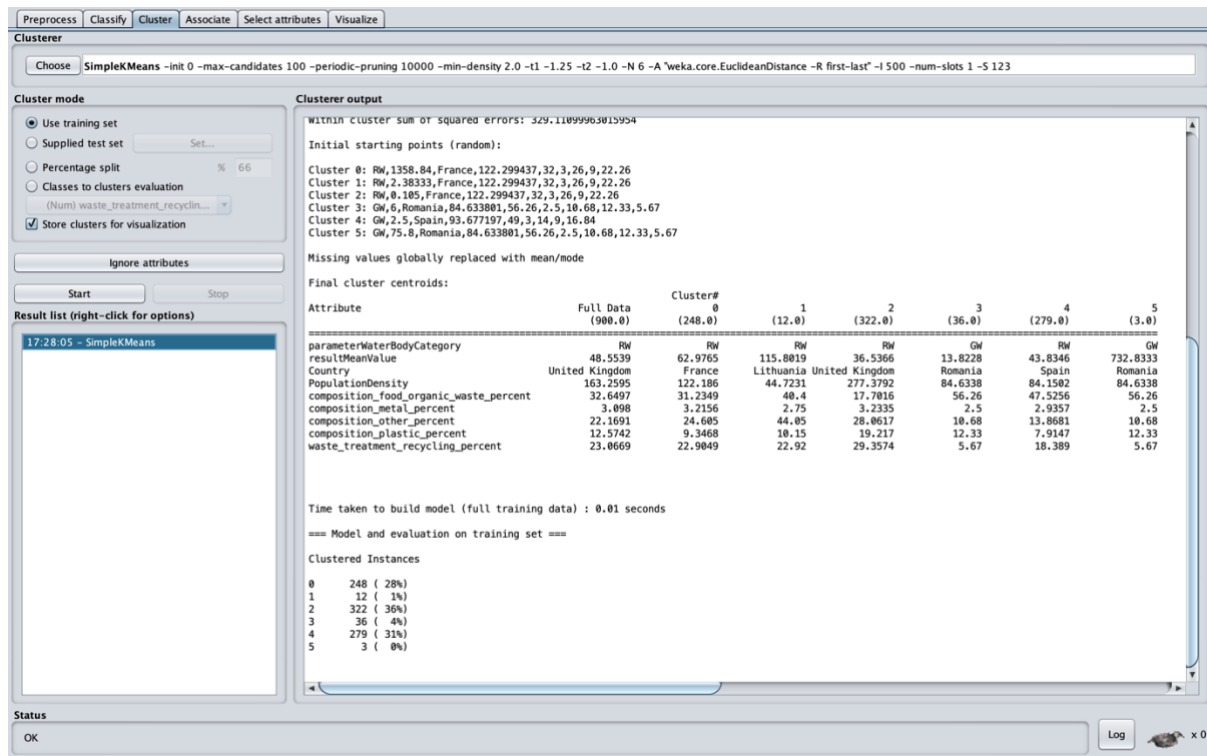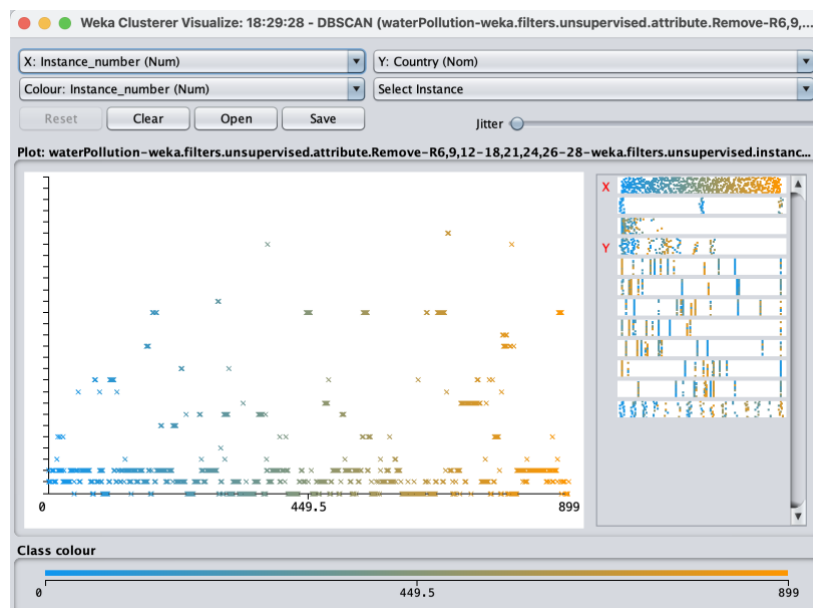
Figure-11

Figure-12

Figure-13

### 4.1.2 DBSCAN

DBSCAN is a clustering method used for machine learning to separate high-density clusters from low-density clusters. Two important parameters of DBSCAN are epsilon and midpoints. (Elutins, 2017).

Epsilon: the radius used to measure the distance between data points.

MinPts: the minPts in the EPS neighbourhood at that point.

First, we set Epsilon to 2 and Minpts to 6, as shown in Figure 14. Then we will set Epsilon to 1 and MinPts to 6, and the results are shown in Figure 15. By comparison, we found that the smaller the Epsilon, the more 'M' values appeared on the image, and conversely, the smaller the Minpts, the less 'M' values appeared on the image. As can be seen from the figure, the areas with heavy colour represent serious water pollution, and the areas with 'M' can be ignored. Countries concerned should treat specific water pollution according to the high frequency areas of K-mean and DBSCAN to protect the environment of the local area.
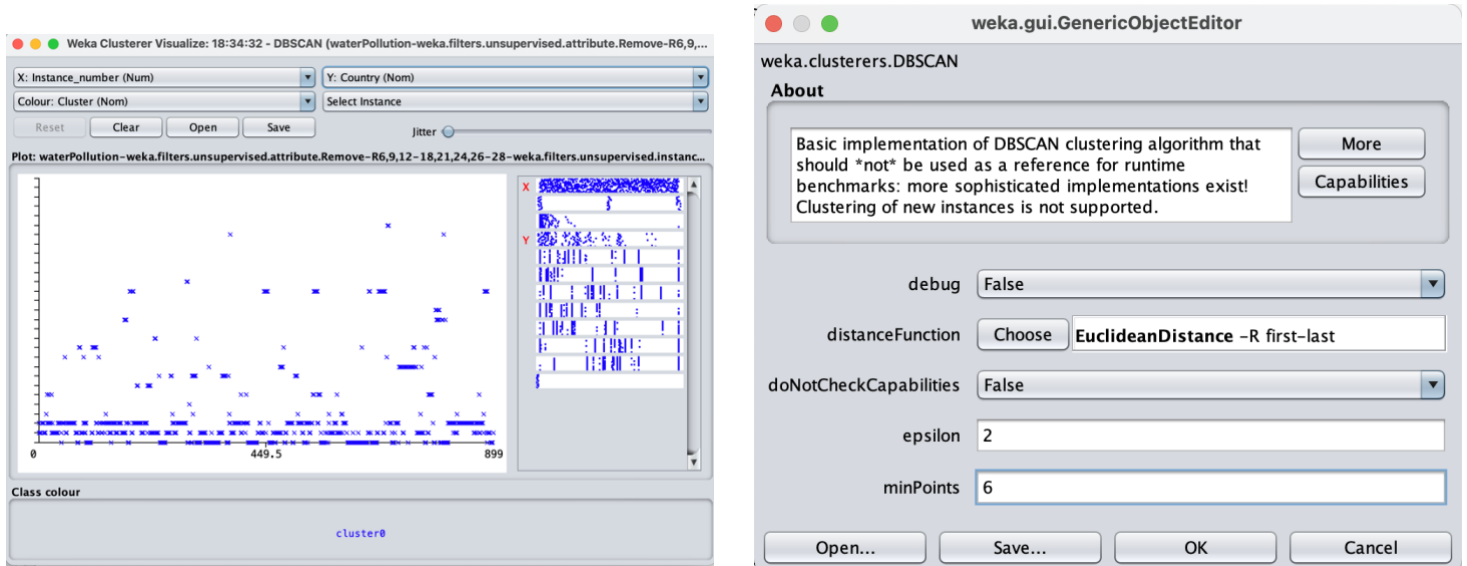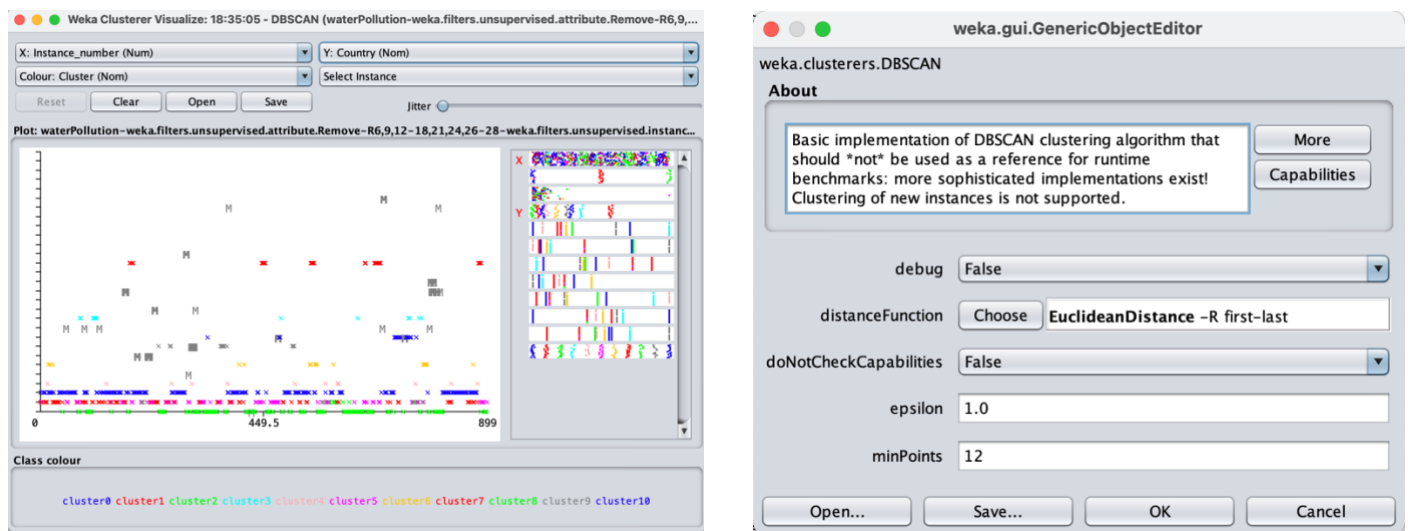


Figure-14



Figure-15

## 4.2    Data Visualization

After data mining, data visualization is needed. In order to convey information clearly and effectively, we use statistical graphs and other tools to express it. See Figure 0. Because

effective visualization helps us analyze and reason about data and evidence. It makes complex data easier to understand and use (Tableau, N. D.).
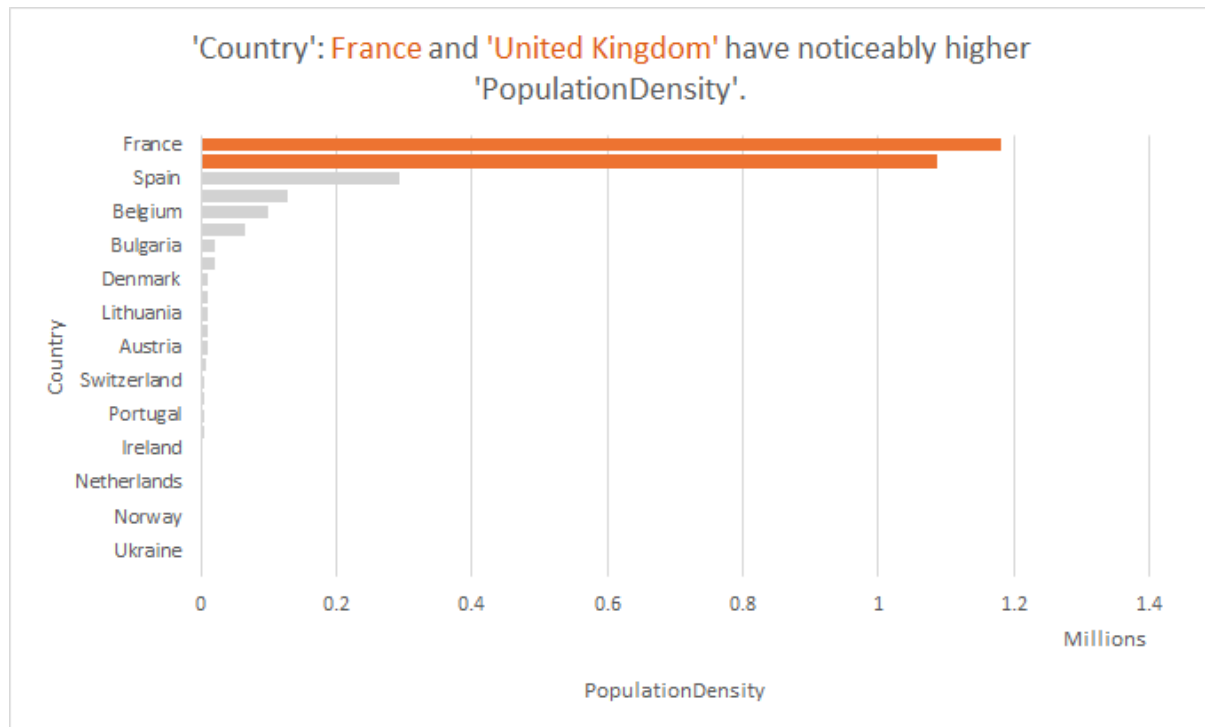


Figure-16

As shown in figure 16, we can see that France has the highest population density compared to other countries (see chart 0), ahead of The UK. Furthermore, the gap in population ratios between France, the United Kingdom and the third country, Spain, was almost double. Therefore, combined with the above data mining (as shown in Figure 0), we can conclude that there is a certain relationship between population density and water pollution. A densely populated country will have more water pollution than a less densely populated country. The reasons are also varied, such as factory sewage, agricultural sewage, and litter and other phenomena caused. Therefore, it is suggested that countries with high population density should reduce water consumption, establish urban sewage treatment system, adjust industrial structure, control agricultural non-point source pollution, rationally develop water resources and develop new water sources.

| Average of resultMeanValue | Column Labels | |
|---|---|---|
| Row Labels | France | 'United Kingdom' |
| uS/cm | 454.9447733 | 282.9649536 |
| '{massRatio}' | 120.1579372 | |
| '\%' | 93.94313604 | 97.87125342 |
| 'mg{NO3}/L' | 14.29107524 | 11.04964289 |
| Cel | 12.08360607 | 9.789306289 |
| 'mg{Si}/L' | 12.44549171 | 5.226174172 |
| [pH] | 7.751358543 | 7.50097068 |
| 'mg{C}/L' | 3.481697733 | 6.782372881 |
| ug/L | 4.345972663 | |
| 'mg{O2}/L' | 3.769420061 | 2.432101923 |
| 'mg{N}/L' | 2.527606751 | 2.515728351 |
| mmol/L | | 1.36333 |
| 'mg{NH4}/L' | 0.105876926 | 0.052776727 |
| 'mg{P}/L' | 0.080565803 | 0.075501874 |
| 'mg{NO2}/L' | 0.086396772 | 0.055581215 |
| 'mg{NH3}/L' | 0.005836 | 0.001036627 |
| Grand Total | 43.14419169 | 24.66543884 |

Figure-17

In addition, water companies can strengthen operations in France and The UK, which are huge water needs. By identifying the two largest countries where companies can focus, we can delve deeper into the water resources of these regions. According to Figure 17, the records for each indicator unit in France and the United Kingdom are very similar. Moreover, with regard to all the numerical indicators of water quality, the water in these two countries can be considered to be of high quality (Naira Hassan Omar, 2019) and can definitely meet the needs of consumers.

## 5    Conclusion

Using public data sets obtained from Kaggle, this report examines three main water quality types in Europe: groundwater (GW), lake water (LW) and river water (RW). We know from data sets that water pollution is closely related to a country's population density, and the more people there are, the more emissions there are. As mentioned above, water pollution in Britain and France is the worst.  we suggest that the government should establish a good sense of water pollution prevention and control in people's minds, if there is no good sense of prevention and control, it will lead to more and more serious water pollution. Therefore, everyone needs to improve the awareness of water pollution prevention, in daily life to minimize the consumption and reuse or recycling of plastic, the correct treatment of chemical cleaners, oil and non-biodegradable items. Everyone should give priority to economical and pollution control. Protect ecology, maintain river health life.

# <u>Reference</u>

Analytics Vidhya.(2020, September 8). *How to Master the Popular DBSCAN Clustering Algorithm for Machine Learning.*

https://www.analyticsvidhya.com/blog/2020/09/how-dbscan-clustering-works/

Codecademy. (n. d). *Normalization.*

https://www.codecademy.com/articles/normalization

HSPH. (n. d). *Water Pollution.*

https://www.hsph.harvard.edu/ehep/82-2/

Investopedia, (n. d). Data mining.

https://www.investopedia.com/terms/d/datamining.asp

Kaggle. (n. d). Water Quality Dataset.

https://www.kaggle.com/ozgurdogan646/water-quality-dataset

Studio. (n. d). *Take control of your R code.*

https://www.rstudio.com/products/rstudio/

Techopedia. (n. d). *K-Means Clustering*.

https://www.techopedia.com/definition/32057/k-means-clustering