



# Water Quality



# TABLE OF CONTENTS

## Contents

<b>ABSTRACT .....</b>	<b>1</b>
<b>INTRODUCTION .....</b>	<b>2</b>
WHAT IS THE WATER QUALITY? .....	2
CONTEXT .....	2
<b>CONTENT .....</b>	<b>3</b>
1. PH VALUE: .....	3
2. HARDNESS: .....	3
3. SOLIDS (TOTAL DISSOLVED SOLIDS - TDS):.....	3
4. CHLORAMINES:.....	4
5. SULFATE:.....	4
6. CONDUCTIVITY: .....	4
7. ORGANIC CARBON:.....	5
8. TRIHALOMETHANES:.....	5
9. TURBIDITY: .....	5
10. POTABILITY: .....	6
<b>II. METHODOLOGY .....</b>	<b>6</b>
3.1. DESCRIPTION OF THE DATASET .....	6
<b>DATA ANALYSIS: .....</b>	<b>7</b>
DATA TABLE .....	7
TREE .....	8
TREE VIEWER .....	8
DISTRIBUTION .....	8
FREE VIZ .....	9
IMPUTE .....	10
CORRELATION .....	10
SCATTER PLOT .....	11
FEATURE STATCICS .....	11
RANK .....	12
<b>DATA MODELING: .....</b>	<b>13</b>
PREPROCESS .....	13
OUTLIERS.....	14
DATA SAMPLER .....	14
TEST AND SCORE.....	15
CONCLUSION MATRIX .....	15
PREDICTION.....	16
SAVE DATA .....	17
SAVE MODEL .....	17
<b>PREDICTION.....</b>	<b>18</b>
LOAD MODEL .....	18
<b>SUMMARY .....</b>	<b>19</b>
<b>REFERENCES.....</b>	<b>19</b>

# TABLE OF FIGURES

FIGURE 1 DATA TABLE .....	7
FIGURE 2 TREE .....	8
FIGURE 3 TREE VIWER .....	8
FIGURE 4 DISTRIBUTION PH. ....	8
FIGURE 5 DISTRIBUTION SOLIDS.....	8
FIGURE 6 FREE VIZ .....	9
FIGURE 7 IMPUTE .....	10
FIGURE 8 CORRELATION .....	10
FIGURE 9 SCATTER PLOT.....	11
FIGURE 10 FEATURE STATCICS .....	11
FIGURE 11 RANK .....	12
FIGURE 12 PREPROCESS .....	13
FIGURE 13 OUTLIERS.....	14
FIGURE 14 DATA SAMPLER .....	14
FIGURE 15 TEST SCORE.....	15
FIGURE 16 CONCLUSION MATRIX .....	15
FIGURE 17 PREDICTION .....	16
FIGURE 18 SAVE DATA.....	17
FIGURE 19 SAVE MODEL .....	17
FIGURE 20 LOAD MODEL.....	18

# Abstract

The second most valuable natural resource after air is water. Although water makes up the majority of the earth's surface, very little of it is actually usable, making it a very limited resource. Therefore, care must be taken when using this valuable and limited resource. Water must be suitable before use because it is needed for a variety of purposes. Additionally, water sources must be checked regularly to see if they are safe or not. Water bodies in poor condition face a threat to the ecosystem as well as being a sign of environmental degradation. In industries, poor water quality can result in risks and significant financial loss.

Water quality thus is essential for both environmental and economic reasons. Analysis of the water's quality is therefore necessary before using it for any purpose. After many years of study, there are now some established protocols for water quality analysis. There are rules for sample collection, storage, and analysis. Here, the typical chain of events is briefly discussed for the benefit of researchers and analysts. financial loss

# Introduction

## What is the water quality?

Water quality refers to the chemical, physical, and biological characteristics of water based on the standards of its usage. It is most frequently used by reference to a set of standards against which compliance, generally achieved through treatment of the water, can be assessed. The most common standards used to monitor and assess water quality convey the health of ecosystems, safety of human contact, extend of water pollution and condition of drinking water. Water quality has a significant impact on water supply and oftentimes determines supply options.

## Context

Access to safe drinking-water is essential to health, a basic human right and a component of effective policy for health protection. This is important as a health and development issue at a national, regional and local level. In some regions, it has been shown that investments in water supply and sanitation can yield a net economic benefit, since the reductions in adverse health effects and health care costs outweigh the costs of undertaking the interventions.



# Content

The water\_potability.csv file contains water quality metrics for 3276 different water bodies.

## 1. pH value:

PH is an important parameter in evaluating the acid–base balance of water. It is also the indicator of acidic or alkaline condition of water status. WHO has recommended maximum permissible limit of pH from 6.5 to 8.5. The current investigation ranges were 6.52–6.83 which are in the range of WHO standards.

## 2. Hardness:

Hardness is mainly caused by calcium and magnesium salts. These salts are dissolved from geologic deposits through which water travels. The length of time water is in contact with hardness producing material helps determine how much hardness there is in raw water. Hardness was originally defined as the capacity of water to precipitate soap caused by Calcium and Magnesium.

## 3. Solids (Total dissolved solids - TDS):

Water has the ability to dissolve a wide range of inorganic and some organic minerals or salts such as potassium, calcium, sodium, bicarbonates, chlorides, magnesium, sulfates etc. These minerals

produced un-wanted taste and diluted color in appearance of water. This is the important parameter for the use of water. The water with high TDS value indicates that water is highly mineralized. Desirable limit for TDS is 500 mg/l and maximum limit is 1000 mg/l which prescribed for drinking purpose.

#### 4. Chloramines:

Chlorine and chloramine are the major disinfectants used in public water systems. Chloramines are most commonly formed when ammonia is added to chlorine to treat drinking water. Chlorine levels up to 4 milligrams per liter (mg/L or 4 parts per million (ppm)) are considered safe in drinking water.

#### 5. Sulfate:

Sulfates are naturally occurring substances that are found in minerals, soil, and rocks. They are present in ambient air, groundwater, plants, and food. The principal commercial use of sulfate is in the chemical industry. Sulfate concentration in seawater is about 2,700 milligrams per liter (mg/L). It ranges from 3 to 30 mg/L in most freshwater supplies, although much higher concentrations (1000 mg/L) are found in some geographic locations.

#### 6. Conductivity:

Pure water is not a good conductor of electric current rather's a good insulator. Increase in ions concentration enhances the electrical conductivity of water. Generally, the amount of dissolved solids in water determines the electrical conductivity. Electrical conductivity

(EC) actually measures the ionic process of a solution that enables it to transmit current. According to WHO standards, EC value should not exceeded 400  $\mu\text{S}/\text{cm}$ .

## 7. Organic carbon:

Total Organic Carbon (TOC) in source waters comes from decaying natural organic matter (NOM) as well as synthetic sources. TOC is a measure of the total amount of carbon in organic compounds in pure water. According to US EPA < 2 mg/L as TOC in treated / drinking water, and < 4 mg/Lit in source water which is use for treatment.

## 8. Trihalomethanes:

THMs are chemicals which may be found in water treated with chlorine. The concentration of THMs in drinking water varies according to the level of organic material in the water, the amount of chlorine required to treat the water, and the temperature of the water that is being treated. THM levels up to 80 ppm is considered safe in drinking water.

## 9. Turbidity:

The turbidity of water depends on the quantity of solid matter present in the suspended state. It is a measure of light emitting properties of water and the test is used to indicate the quality of waste discharge with respect to colloidal matter. The mean turbidity value obtained for Wondo Genet Campus (0.98 NTU) is lower than the WHO recommended value of 5.00 NTU.



## 10. Potability:

Indicates if water is safe for human consumption where 1 means Potable and 0 means Not potable.

---

## II. METHODOLOGY

### 3.1. Description of the Dataset

**ppm:** parts per million

**µg/L:** microgram per liter

**mg/L:** milligram per liter

**Column description:**

1. **ph:** pH of 1. water (0 to 14).
2. **Hardness:** Capacity of water to precipitate soap in mg/L.
3. **Solids:** Total dissolved solids in ppm.
4. **Chloramines:** Amount of Chloramines in ppm.
5. **Sulfate:** Amount of Sulfates dissolved in mg/L.
6. **Conductivity:** Electrical conductivity of water in µS/cm.
7. **Organic carbon:** Amount of organic carbon in ppm.
8. **Trihalomethanes:** Amount of Trihalomethanes in µg/L.
9. **Turbidity:** Measure of light emitting property of water in NTU.
10. **Potability:** Indicates if water is safe for human consumption.  
Potable -1 and not potable -0\

# Data analysis:

## Data Table

Data Table Displays attribute-value data in a spreadsheet.

Info

1275 instances

9 features (4 of them missing data)

Integer with 5 values

16 more attributes

Visualize

Filter variable table (if present)

Visualize missing values

Color by missing data

Selection

Feature selection

	Pathology	Age	Intermass	Solids	Cholesterol	Satets	Creativity	Deposits	Trichomonas	Turbidity
1	0	254.69	20751.2	7.30021	388.516	364.320	10.2478	86.291	2.06314	
2	0	1.71600	129.423	18630.1	6.63525	590.865	15.18	56.1291	4.50066	
3	0	5.03912	224.236	19609.5	9.27588	418.806	16.8686	86.4201	3.05535	
4	0	5.17677	214.373	27018.4	8.05933	316.086	16.1767	100.147	4.62677	
5	0	9.09222	181.102	17970	6.5466	310.136	398.411	11.5583	4.07938	
6	0	5.58469	188.313	28740.7	7.54407	326.675	280.468	8.39973	2.55971	
7	0	10.2229	248.072	28769.7	7.51181	333.663	280.652	13.7897	2.67299	
8	0	8.63505	203.362	13672.1	4.56301	305.31	474.688	12.3638	4.40142	
								12.706	3.55922	
								7.9078	4.37056	
								5.5888	3.86229	
12	0	7.97452	218.683	18767.7	8.11038	364.098	14.5757	76.4858	4.01177	
13	0	7.11562	158.705	18750.8	3.68604	282.344	347.715	15.0245	3.68576	
14	0	150.175	27331.4	6.83622	295.416	379.762	19.5708	70.51	4.41337	
15	0	7.89623	205.345	20388	5.07256	464.645	13.2283	70.3002	4.77738	
16	0	6.14727	189.733	41063.2	9.6256	304.486	516.143	11.5398	4.37635	
17	0	7.02579	211.049	30560.6	10.2948	215.161	20337	56.5516	4.66813	
18	0	9.18156	271.614	24061.3	6.40488	346.351	13.3073	21.4574	4.50566	
19	0	8.07546	279.257	19460.4	6.20432	431.448	12.8888	63.8212	2.65609	
20	0	7.37105	214.487	25650.3	4.43267	315.754	469.915	17.5037	2.5603	
21	0	227.435	27305.6	10.3339		954.82	16.3317	45.3828	4.13342	
22	0	6.60021	168.284	30544.4	5.85677	310.331	523.671	17.8842	3.7437	
23	0	215.978	17107.2	5.60706	326.344	436.256	14.1091	59.8055	5.45025	
24	0	3.95240	196.903	21167.5	6.99631	444.478	16.609	30.1077	4.50532	
25	0	5.4003	152.739	17286.6	10.2569	328.358	472.874	11.2564	4.62779	
26	0	6.51442	180.767	21218.7	0.57084	323.586	413.29	14.9	5.20039	
27	0	5.44506	207.506	35424.8	8.78215	384.007	441.786	15.8059	4.1814	
28	0	145.768	11274.9	7.90644	304.000	290.991	17.7795	49.5760	4.00487	
29	0	268.431	26363	7.70006	315.389	364.48	10.349	53.0064	3.99156	
30	0	148.153	15193.4	9.04603	307.012	563.805	16.5687	52.6762	6.03038	
31	0	7.38165	209.626	15106.2	6.99468	338.336	342.111	7.9226	5.00036	
32	0	8.04549	190.797	19677.9	6.75754	452.836	16.099	47.082	2.85747	
33	0	10.4233	117.791	22326.2	8.1615	307.708	412.967	12.8007	5.05731	
34	0	7.61415	235.045	12555.9	6.84585	387.175	411.963	10.2548	3.16052	

Figure 1 data table

Info

1276 instances

9 features (4 of them missing data)

Integer with 5 values

16 more attributes

Visualize

Filter variable table (if present)

Visualize missing values

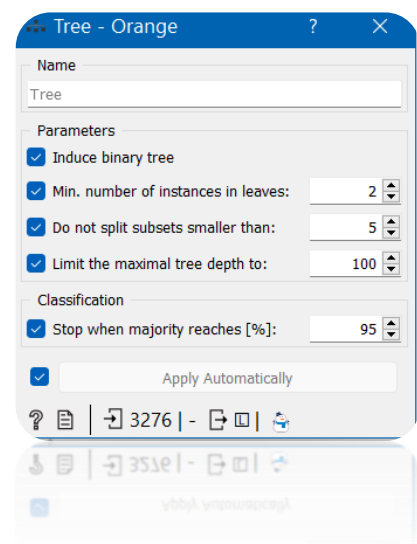
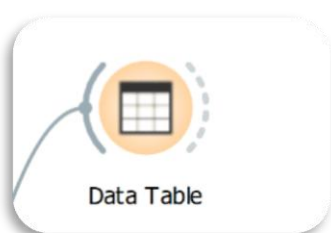
Color by missing data

Selection

Feature selection

	Pathology	Age	Intermass	Solids	Cholesterol	Satets	Creativity	Deposits	Trichomonas	Turbidity
12	0	7.97452	218.683	18767.7	8.11038	364.098	14.5757	76.4858	4.01177	
13	0	7.11562	158.705	18750.8	3.68604	282.344	347.715	15.0245	3.68576	
14	0	150.175	27331.4	6.83622	295.416	379.762	19.5708	70.51	4.41337	
15	0	7.89623	205.345	20388	5.07256	464.645	13.2283	70.3002	4.77738	
16	0	6.14727	189.733	41063.2	9.6256	304.486	516.143	11.5398	4.37635	
17	0	7.02579	211.049	30560.6	10.2948	215.161	20337	56.5516	4.66813	
18	0	9.18156	271.614	24061.3	6.40488	346.351	13.3073	21.4574	4.50566	
19	0	8.07546	279.257	19460.4	6.20432	431.448	12.8888	63.8212	2.65609	
20	0	7.37105	214.487	25650.3	4.43267	315.754	469.915	17.5037	2.5603	
21	0	227.435	27305.6	10.3339		954.82	16.3317	45.3828	4.13342	
22	0	6.60021	168.284	30544.4	5.85677	310.331	523.671	17.8842	3.7437	
23	0	215.978	17107.2	5.60706	326.344	436.256	14.1091	59.8055	5.45025	
24	0	3.95240	196.903	21167.5	6.99631	444.478	16.609	30.1077	4.50532	
25	0	5.4003	152.739	17286.6	10.2569	328.358	472.874	11.2564	4.62779	
26	0	6.51442	180.767	21218.7	0.57084	323.586	413.29	14.9	5.20039	
27	0	5.44506	207.506	35424.8	8.78215	384.007	441.786	15.8059	4.1814	
28	0	145.768	11274.9	7.90644	304.000	290.991	17.7795	49.5760	4.00487	
29	0	268.431	26363	7.70006	315.389	364.48	10.349	53.0064	3.99156	
30	0	148.153	15193.4	9.04603	307.012	563.805	16.5687	52.6762	6.03038	
31	0	7.38165	209.626	15106.2	6.99468	338.336	342.111	7.9226	5.00036	
32	0	8.04549	190.797	19677.9	6.75754	452.836	16.099	47.082	2.85747	
33	0	10.4233	117.791	22326.2	8.1615	307.708	412.967	12.8007	5.05731	
34	0	7.61415	235.045	12555.9	6.84585	387.175	411.963	10.2548	3.16052	

Figure 1 data table



# Tree

We must make tree to display the tree viewer

## Tree viewer

It's show's what is the most feature has impact on target

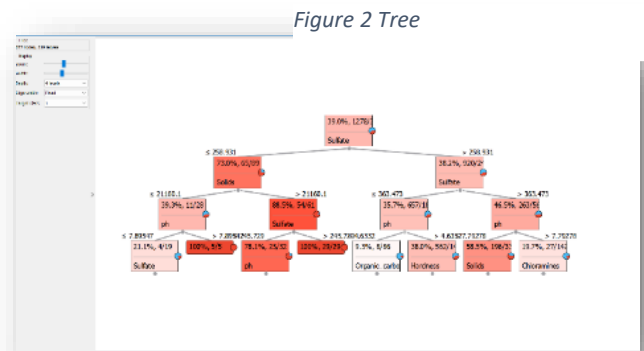


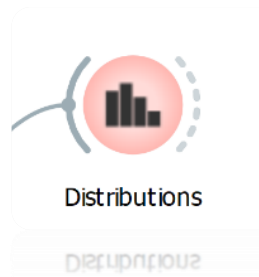
Figure 3 Tree viewer

## Distribution

It is used to present the data selected and data limitation graphically.

Example like PH & solids

PH



solids

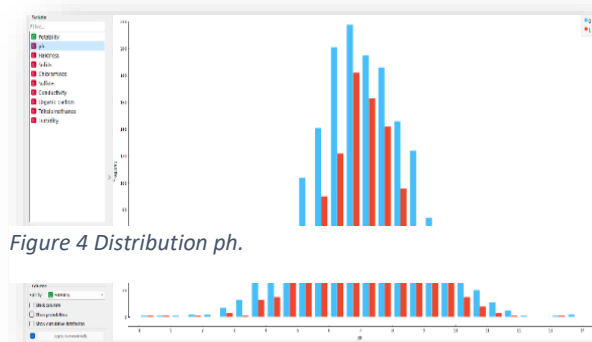


Figure 4 Distribution ph.

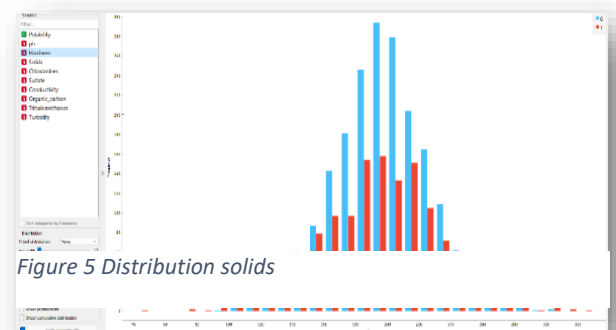


Figure 5 Distribution solids

## Free viz

A widget that displays the dataset's most significant features as circles with the largest radius is an important factor, and a smaller radius has less impact.

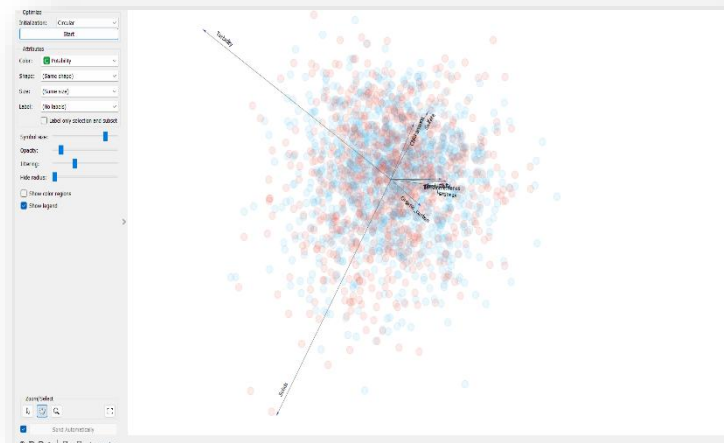
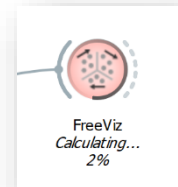
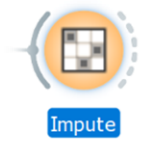


Figure 6 Free Viz



# Impute

Get rid from any missing data & substitutes missing values by values.



Don't Impute does nothing with the missing values.

Average/Most-frequent uses the average value (for continuous attributes) or the most common value (for discrete attributes).

As a distinct value creates new values to substitute the missing ones.

Model-based imputer constructs a model for predicting the missing value, based on values of other attributes

Random values computes the distributions of values for each attribute and then imputes by picking random values from them.

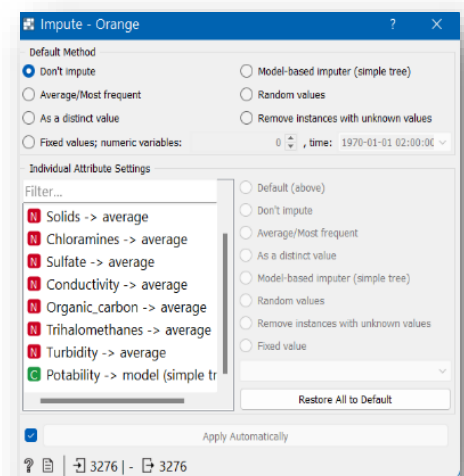


Figure 7 Impute

# Correlation

It is a widget that illustrates the relationship between two features from the selected data set so that we can evaluate their connection.

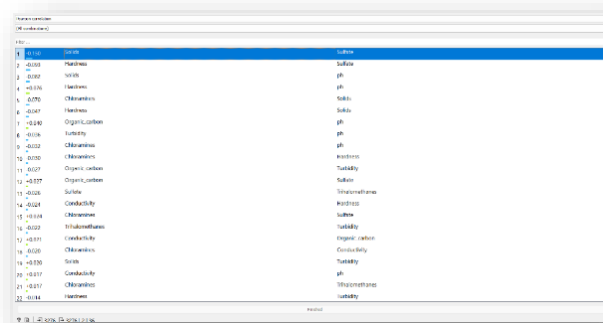


Figure 8 correlation



# Scatter plot

A widget that displays the relationship between two numerical data points from the chosen dataset.



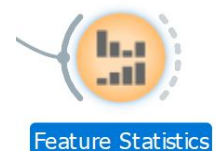
Figure 9 scatter plot

# Feature statistics

The purpose of this widget is to display statistics about each feature in the chosen dataset, allowing the user to see the (Mean, Dispersion, min, max and missing values)



Figure 10 feature statistics





# Rank

A tool for ranking the dataset's most significant features in order of highest to lowest effect.

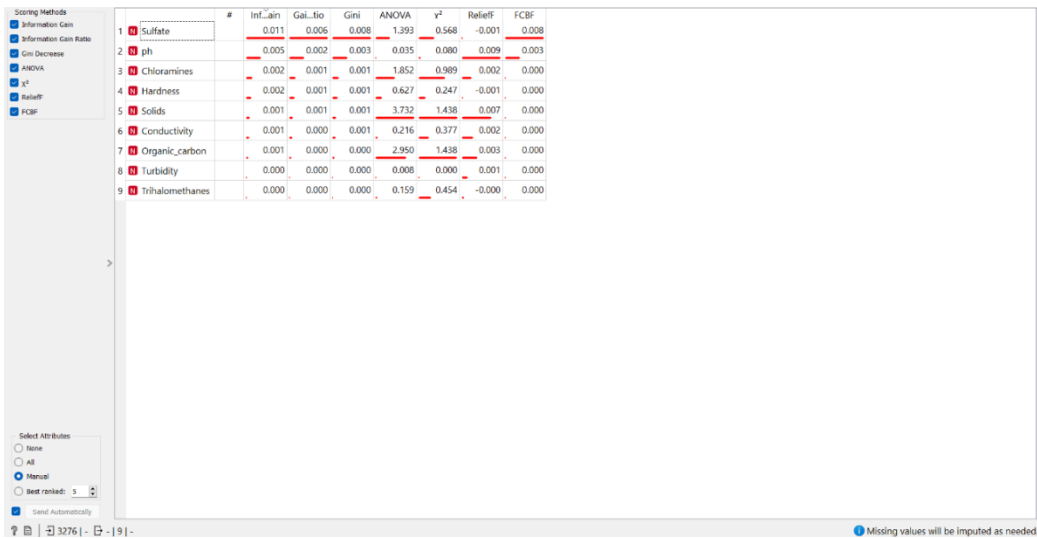


Figure 11 Rank

# Data modeling:

## Preprocess

- Preprocessing is crucial for achieving better-quality analysis results. The Preprocess widget offers several preprocessing methods that can be combined in a single preprocessing pipeline. Some methods are available as separate widgets, which offer advanced techniques and greater parameter tuning.
- Divide by number of values is similar to treat as ordinal, but the final values will be divided by the total number of values and hence the range of the new continuous variable will be  $[0, 1]$ .

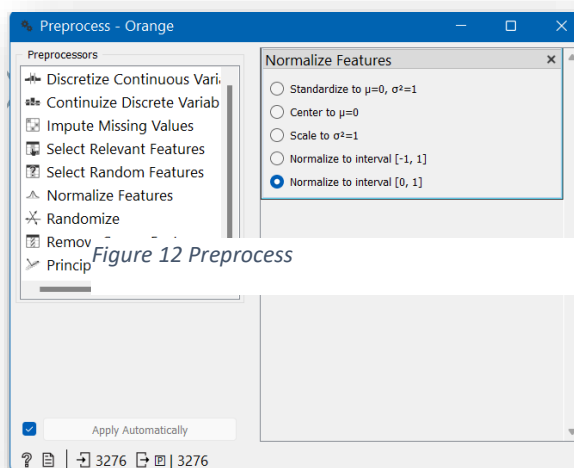


Figure 12 Preprocess



# Outliers

Weird data should be used to remove any abnormal data

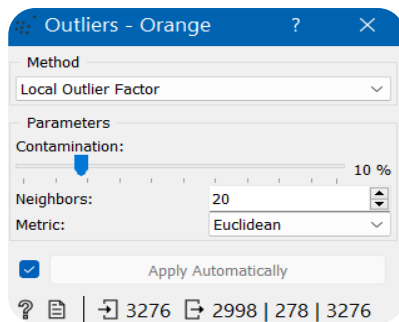
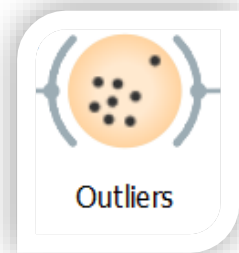


Figure 13 outliers



## Data sampler

The Data Sampler widget implements several data sampling methods. It outputs a sampled and a complementary dataset (with instances from the input set that are not included in the sampled dataset). The output is processed after the input dataset is provided and Sample Data is pressed.

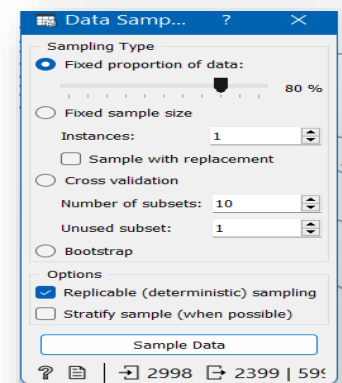


Figure 14 data sampler

# Test and score

Evaluation Results: results of testing classification algorithms and we chose the best one that has high accuracy

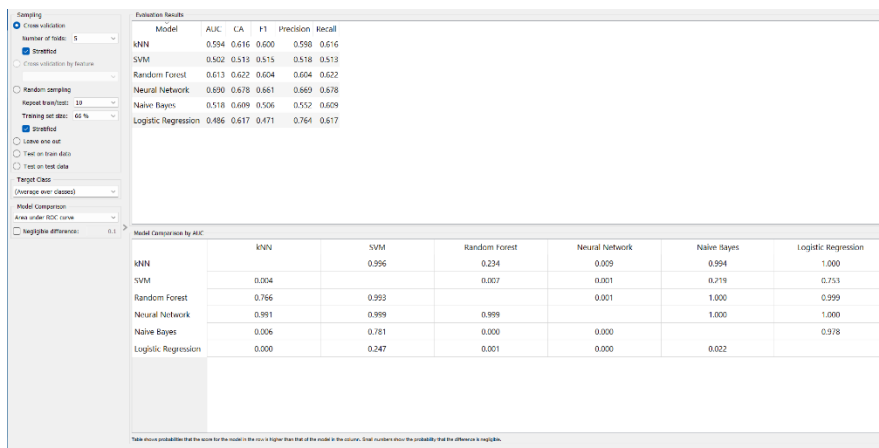


Figure 15 test score

# Conclusion matrix

Evaluation results: results of testing classification algorithms. The widget tests learning algorithms

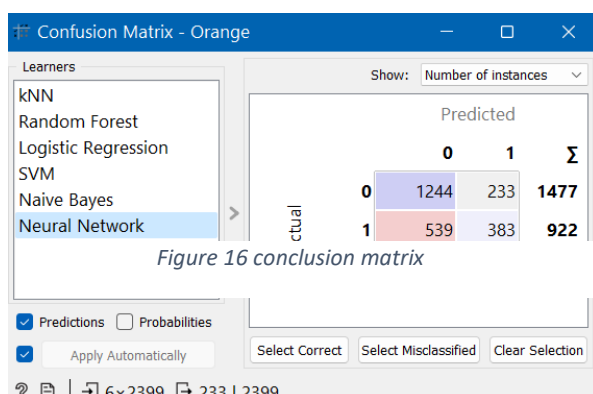
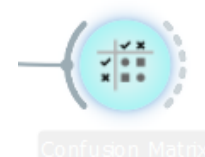


Figure 16 conclusion matrix



# Prediction

The widget receives a dataset and one or more predictors (predictive models, not learning algorithms - see the example below). It outputs the data and the predictions.

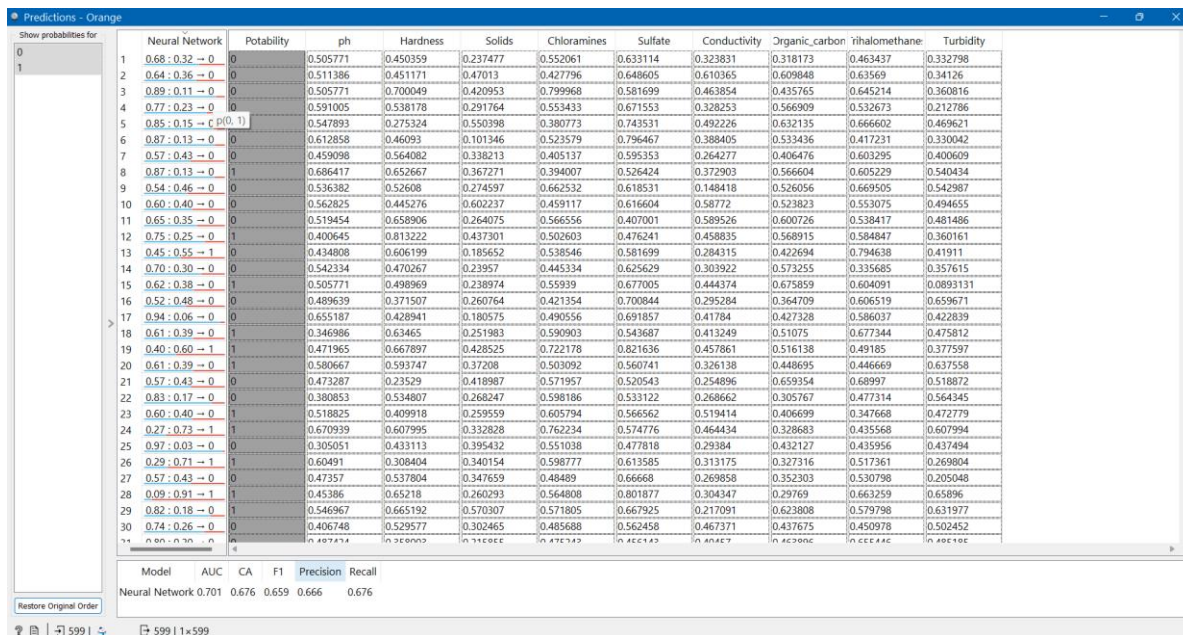


Figure 17 prediction

## Save data

The Save Data widget considers a dataset provided in the input channel and saves it to a data file with a specified name. It can save the data as:

Excel spreadsheets (.xlsx)

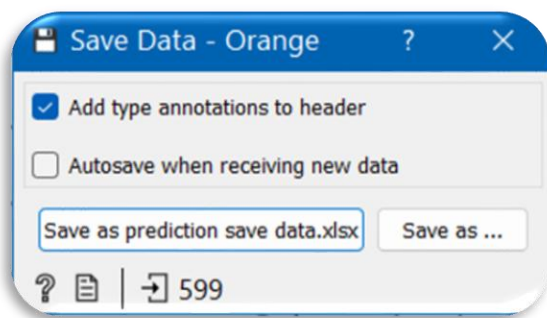
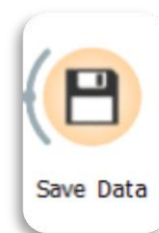


Figure 18 save data



## Save model

If the file is saved to the same directory as the workflow or in the subtree of that directory, the widget remembers the relative path. Otherwise, it will store an absolute path, but disable auto save for security reasons.

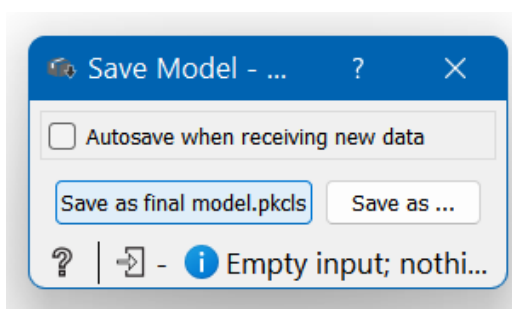
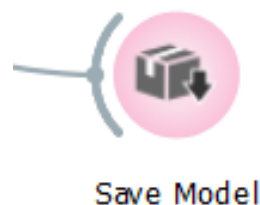


Figure 19 save model

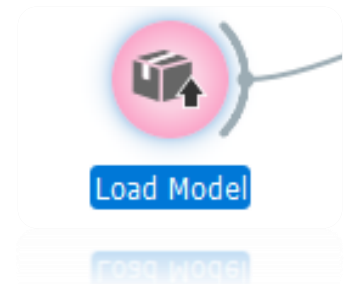
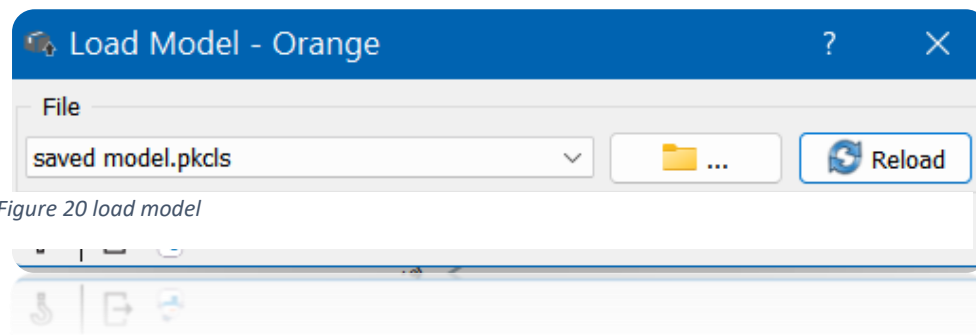




# Prediction

## Load model

When you want to use a custom-set model that you've saved before, open the Load Model widget and select the desired file with the Browse icon. This widget loads the existing model into Predictions widget. Datasets used with Load Model have to contain compatible attributes!



## Summary

Access to safe drinking water is one of the essential needs of all human beings. From a legal point of view, access to drinking water is one of the fundamental human rights. Many factors affect water quality, it is also one of the major research areas in machine learning.

So, this is how you can analyze the quality of water and train a machine learning model to classify safe and unsafe water for drinking. Access to safe drinking water is one of the essential needs of all human beings. From a legal point of view, access to drinking water is one of the fundamental human rights. Many factors affect water quality, it is also one of the major research areas in machine learning.

## References

1-Water Quality Analysis | Aman Kharwal (thecleverprogrammer.com)

2-[https://en.wikipedia.org/wiki/Water\\_quality](https://en.wikipedia.org/wiki/Water_quality)

3-Orange Data Mining - Data Mining

4- Find Open Datasets and Machine Learning Projects | Kaggle