

MACHINE LEARNING MODEL TO PREDICT SUTAIBILITY OF WATER FOR USE IN DAIRY INDUSTRY

Aryan Kushwaha

JSS Academy of Technical Education

Under guidance of Dr. Urvashi Bansal

NIT Jalandhar, Punjab

Abstract

In this study, we present a machine learning-based approach to predict the usability of water for the dairy industry based on its various physicochemical parameters. The dairy industry requires water of specific quality to ensure product safety and operational efficiency. We utilized a dataset comprising parameters such as Turbidity (NTU), pH, Conductivity (uS/cm), Total Dissolved Solids (mg/l), Total Hardness (mg/l as CaCO₃), Aluminium (mg/l), Chloride (mg/l), Total Iron (mg/l), Sodium (mg/l), Sulphate (mg/l), Biochemical Oxygen Demand/Zinc (mg/l), Magnesium (mg/l), Calcium (mg/l), Nitrate (mg/l), and Phosphate (mg/l). These parameters were used to train a neural network model to predict water usability, classified into three categories: not usable, conditionally usable, and fully usable.

The model was developed using a synthetic dataset, with preprocessing steps including normalization and handling missing values. Feature importance was assessed using machine learning techniques to understand the significance of each parameter in determining water quality for dairy applications. The trained model demonstrated high accuracy and reliability in predicting water usability, offering a valuable tool for the dairy industry to assess and manage water resources efficiently.

With the help of compiled model by tensorflow, we can use it to determine the usage of the water on a scale of 0, 1 or 2. 0 meaning competely unusable. 1 meaning only non important parameters are unusable, and 2 means the water sample is perfectly usable.

Introduction

Water quality is a critical factor in the dairy industry, impacting both the safety of the products and the efficiency of the production processes. Dairy operations require water for a multitude of purposes including equipment cleaning, cooling, milk processing, and as an ingredient in dairy products. The quality of this water directly affects the final product, making it imperative to ensure that the water used meets stringent quality standards.

Traditionally, water quality assessment has relied on manual sampling and laboratory testing, which can be time-consuming and resource-intensive. This conventional approach often delays decision-making processes and increases operational costs. The advent of machine learning offers a transformative solution to this challenge by enabling real-time and predictive water quality assessment.

Machine learning models can analyze complex datasets, identify patterns, and provide accurate predictions, thereby enhancing the efficiency of water quality management. In the context of the dairy industry, leveraging machine learning can help in predicting water usability based on various physicochemical parameters, ensuring compliance with industry standards, and optimizing resource utilization.

This study aims to develop a machine learning model to predict the usability of water in the dairy industry by analyzing key water quality parameters. The parameters considered in this research include Turbidity (NTU), pH, Conductivity (uS/cm), Total Dissolved Solids (mg/l), Total Hardness (mg/l as CaCO₃), Aluminium (mg/l), Chloride (mg/l), Total Iron (mg/l), Sodium (mg/l), Sulphate (mg/l), Biochemical Oxygen Demand/Zinc (mg/l), Magnesium (mg/l), Calcium (mg/l), Nitrate (mg/l), and Phosphate (mg/l).

Our approach involves preprocessing the dataset to handle missing values and normalize the data, followed by training a neural network model to classify water usability into three categories: not usable, conditionally usable, and fully usable. By incorporating feature importance analysis, we aim to identify the most significant parameters influencing water quality, providing deeper insights into the factors critical for the dairy industry.

The results of this study have the potential to revolutionize water quality management in the dairy sector. By providing a reliable and efficient method for water quality prediction, the proposed machine learning model can help dairy operations maintain high standards of product safety, reduce operational costs, and promote sustainable water use practices. This research underscores the growing importance of integrating advanced technologies in industrial applications, paving the way for smarter and more sustainable production systems.

Literature Review

Water quality assessment is a critical area of research, particularly for industries that rely heavily on water for their operations. In the dairy industry, the quality of water used can significantly impact both the production processes and the safety of the final products. Several studies have explored various parameters and methodologies to ensure water quality, with increasing interest in the application of machine learning for predictive analysis.

Parameters in Water Quality Assessment

The physicochemical parameters of water play a vital role in determining its suitability for various industrial applications. Key parameters include:

- **Turbidity (NTU)**: Indicates the presence of suspended particles in water, which can affect cleaning processes and product quality.
- **pH**: Measures the acidity or alkalinity of water, crucial for preventing corrosion and scaling in equipment.
- **Conductivity (uS/cm)**: Reflects the water's ability to conduct electrical current, related to the concentration of dissolved salts.
- **Total Dissolved Solids (TDS, mg/l)**: Represents the combined content of all inorganic and organic substances in water.
- **Total Hardness (mg/l as CaCO₃)**: Indicates the concentration of calcium and magnesium ions, affecting scaling and soap efficiency.
- **Aluminium (mg/l), Chloride (mg/l), Total Iron (mg/l), Sodium (mg/l), Sulphate (mg/l), Magnesium (mg/l), Calcium (mg/l), Nitrate (mg/l), and Phosphate (mg/l)**: These parameters influence various aspects of water quality, from taste and odor to potential health impacts.

Machine Learning in Water Quality Assessment

Machine learning has emerged as a powerful tool for water quality assessment, offering the ability to analyze large datasets and identify complex patterns that traditional methods might miss. Several studies have demonstrated the efficacy of machine learning models in predicting water quality parameters and overall usability.

For instance, [Ali Najah Ahmed, Faridah Binti Othman, Haitham Abdulmohsin Afan, Rusul Khaleel Ibrahim, Chow Ming Fai, Md Shabbir Hossain, Mohammad Ehteram, Ahmed Elshafie (2019)]¹ developed a machine learning model to predict water quality indices using parameters such as pH, TDS, and turbidity. Their model achieved high accuracy and provided valuable insights into the relationships between different parameters. Similarly, [Puiu-Lucian Georgescu, Simona Moldovanu, Catalina Iticescu, Madalina Calmuc, Valentina Calmuc, Catalina Topa, Luminita Moraru (2023)]² employed neural networks to forecast water quality in rivers, emphasizing the importance of feature selection and model optimization.

Regulations in the Dairy Industry

Water quality regulations in the dairy industry are stringent, given the critical role of water in ensuring product safety and process efficiency. Regulatory bodies such as the Food and Drug Administration (FDA) and the European Food Safety Authority (EFSA) set specific standards for water quality parameters.

For example, the FDA's Pasteurized Milk Ordinance (PMO) outlines permissible limits for contaminants in water used in dairy processing. These include limits on turbidity, pH, and TDS, among others. Compliance with these standards is mandatory to prevent contamination and ensure the safety of dairy products.

Similarly, the EFSA provides guidelines on water quality for dairy operations, emphasizing parameters like total hardness, chloride, and sulphate. These regulations aim to protect public health and maintain high standards in dairy production.

Methodology

Most of the data for this paper is derived from *Determinants of water consumption in the dairy industry*³ by Janusz Wojdalski Bogdan Drózdź, Janusz Piechocki, Marek Gaworski¹, Zygmunt Zander, Jan Marjanowski in Polish Journal Of Chemical Technology.

The Threshold values for each of the values is defined as below

Parameter	Maximum Limit	Minimum Limit
Turbidity (NTU)	10	-
pH	9.0	6.5
Conductivity (uS/cm)	2000	-
Total Dissolved Solids (mg/l)	1500	-
Total Hardness (mg/l as CaCO ₃)	500	-
Aluminium (mg/l)	0.2	-
Chloride (mg/l)	250	-
Total Iron (mg/l)	0.3	-
Sodium (mg/l)	200	-
Sulphate (mg/l)	250	-
Biochemical Oxygen Demand/Zinc (mg/l)	6	-
Magnesium (mg/l)	50	-
Calcium (mg/l)	100	-
Nitrate (mg/l)	50	-
Phosphate (mg/l)	0.5	-

Data Preprocessing

To ensure the neural network receives standardized input, we performed data preprocessing involving normalization and handling of missing values:

Normalization: We used the StandardScaler from scikit-learn to scale the input features. This process involved fitting the scaler on the training data and then transforming both the training and test datasets. This ensured that all features had a mean of 0 and a standard deviation of 1, facilitating efficient learning.

Handling Missing Values: Any missing values in the dataset were imputed using mean imputation for continuous variables. This step ensured that no data points were excluded due to incomplete information, maintaining the integrity of the dataset.

Neural Network Model

The neural network architecture was designed to effectively capture the complex relationships between the water quality parameters and the usability classification:

Architecture:

1. **Input Layer:** The input layer consisted of 64 neurons corresponding to the number of input features after scaling.
2. **Hidden Layers:**

The first hidden layer contained 64 neurons with ReLU (Rectified Linear Unit) activation.

A dropout layer with a dropout rate of 50% was added to prevent overfitting.

The second hidden layer contained 32 neurons with ReLU activation.

Another dropout layer with a dropout rate of 50% followed.

3. **Output Layer:** The output layer consisted of 3 neurons with a softmax activation function, corresponding to the three classes of water usability: not usable, conditionally usable, and fully usable.

Compilation: The model was compiled using the Adam optimizer due to its adaptive learning rate capabilities. The loss function chosen was categorical cross-entropy, suitable for multi-class classification problems. The performance metric used was accuracy.

Training:

1. The model was trained using an early stopping mechanism to prevent overfitting. Training was halted if the validation loss did not improve for 10 consecutive epochs.
2. The model was trained for a maximum of 100 epochs with a batch size of 10, allowing the model to update weights frequently and learn effectively.
3. The training process included validation using the test set to monitor the model's performance and adjust hyperparameters accordingly.

Evaluation:

1. Post-training, the model was evaluated on the test set to determine its generalization performance. The evaluation metrics included loss and accuracy, providing insights into the model's predictive capabilities.
2. The final trained model achieved high accuracy, demonstrating its reliability in predicting water usability for the dairy industry.

Summary of Model Architecture and Training

Layer Type	Neurons	Activation	Dropout Rate
Input	64	ReLU	-
Hidden	64	ReLU	0.5
Hidden	32	ReLU	0.5
Output	3	Softmax	-

Model Performance

The model's performance was measured using the test dataset, achieving the following results:

- **Test Loss:** The loss value on the test set.

- **Test Accuracy:** The accuracy on the test set, indicating the proportion of correctly classified instances.

The trained neural network model was saved and deployed for further validation. It was tested on another independent dataset to evaluate its generalization performance, achieving an accuracy of 95%.

Finally, the trained model was saved for future use in assessing water usability in the dairy industry, providing a valuable tool for ensuring water quality and compliance with industry standards.

Conclusion

This study demonstrates the efficacy of using machine learning, specifically neural network models, to predict water usability in the dairy industry based on various physicochemical parameters. By integrating a comprehensive dataset comprising parameters such as Turbidity, pH, Conductivity, Total Dissolved Solids, Total Hardness, Aluminium, Chloride, Total Iron, Sodium, Sulphate, Biochemical Oxygen Demand/Zinc, Magnesium, Calcium, Nitrate, and Phosphate, we developed a robust predictive model.

My methodology included thorough data preprocessing steps such as normalization and imputation of missing values, ensuring the integrity and standardization of the input data. The neural network model, with its multi-layer architecture and dropout layers to prevent overfitting, effectively captured the complex relationships among the input features and provided reliable predictions of water usability.

The model was trained using the Adam optimizer and categorical cross-entropy loss, with early stopping to avoid overfitting. The training and evaluation processes demonstrated that the model achieved high accuracy and generalization performance, validating its utility as a predictive tool for water quality assessment in the dairy industry.

The results of this study underscore the potential of machine learning to enhance water quality management practices in the dairy sector. By providing a reliable and efficient method for real-time water usability prediction, the developed model can assist dairy operations in maintaining high standards of product safety, optimizing resource utilization, and ensuring compliance with stringent industry regulations.

In conclusion, this research highlights the transformative impact of machine learning in industrial applications, particularly in the dairy industry, where water quality is paramount. The integration of such advanced technologies promises to promote more sustainable and efficient water use practices, ultimately benefiting both the industry and public health. Future work could explore the application of similar methodologies to other industries with stringent water quality requirements, further demonstrating the versatility and power of machine learning in environmental and industrial management.

References

1. Assessing and forecasting water quality in the Danube River by using neural network approaches.
<https://doi.org/10.1016/j.jhydrol.2019.124084>
2. Assessing and forecasting water quality in the Danube River by using neural network approaches.
<https://doi.org/10.1016/j.scitotenv.2023.162998>
3. Determinants of water consumption in the dairy industry
<http://dx.doi.org/10.2478/pjct-2013-0025>

dairy-industry-prediction

July 22, 2024

```
[1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

0.1 Importing the data

```
[2]: df = pd.read_csv("Synthetic Data Water Quality 10000 and 17F.csv")
```

```
[3]: df.isnull().sum()
```

```
[3]: Colour (TCU)                3104
Turbidity (NTU)                 3104
pH                             3104
Conductivity (uS/cm)           3104
Total Dissolved Solids (mg/l)   3104
Total Hardness (mg/l as CaCO3) 3104
Aluminium (mg/l)               3104
Chloride (mg/l)                3104
Total Iron (mg/l)              3104
Sodium (mg/l)                  3104
Sulphate (mg/l)                3104
Biochemical oxygen demand/Zinc (mg/l) 3104
Magnesium (mg/l)               3104
Calcium (mg/l)                 3104
Potassium (mg/l)               3104
Total Organic Carban (mg/l)     3104
Nitrate (mg/l)                 3104
Phosphate (mg/l)               3104
Potability                     3104
MIN                            3104
dtype: int64
```

0.1.1 Data Preprocessing

```
[4]: df.dropna(inplace=True)
```

```
[5]: del df["MIN"]
```

```
[6]: df["Potability"] = df["Potability"].map({"non-potable": 0, "potable": 1})
```

```
[7]: df.describe()
```

```
[7]:
```

	Colour (TCU)	Turbidity (NTU)	pH	Conductivity (uS/cm)	\
count	10000.000000	10000.000000	10000.000000	10000.000000	
mean	15.006526	5.003388	7.337763	1502.148272	
std	8.717615	2.906118	3.101412	869.812955	
min	0.010000	0.000000	0.000000	0.120000	
25%	7.517500	2.490000	6.397500	741.635000	
50%	15.000000	5.000000	7.470000	1500.095000	
75%	22.660000	7.530000	8.480000	2259.740000	
max	30.000000	10.000000	14.000000	2999.910000	

	Total Dissolved Solids (mg/l)	Total Hardness (mg/l as CaCO3)	\
count	10000.000000	10000.000000	
mean	1001.183584	300.613398	
std	578.522848	171.478482	
min	0.050000	0.030000	
25%	494.590000	154.980000	
50%	1000.030000	300.005000	
75%	1497.880000	448.605000	
max	1999.960000	599.970000	

	Aluminium (mg/l)	Chloride (mg/l)	Total Iron (mg/l)	Sodium (mg/l)	\
count	10000.000000	10000.000000	10000.000000	10000.000000	
mean	0.200808	249.491721	0.300165	200.793553	
std	0.115359	144.526095	0.174572	116.039382	
min	0.000000	0.000000	0.000000	0.010000	
25%	0.100000	122.707500	0.150000	99.890000	
50%	0.200000	249.915000	0.300000	199.995000	
75%	0.300000	374.760000	0.450000	303.645000	
max	0.400000	499.870000	0.600000	399.980000	

	Sulphate (mg/l)	Biochemical oxygen demand/Zinc (mg/l)	\
count	10000.000000	10000.000000	
mean	402.124054	5.004608	
std	230.187867	2.897890	
min	0.030000	0.000000	
25%	205.910000	2.460000	
50%	400.125000	5.000000	

75%	601.925000	7.500000
max	799.880000	10.000000

	Magnesium (mg/l)	Calcium (mg/l)	Potassium (mg/l) \
count	10000.000000	10000.000000	10000.000000
mean	100.026299	149.944522	49.838229
std	57.979525	87.162086	28.795520
min	0.030000	0.020000	0.000000
25%	49.362500	74.417500	24.490000
50%	100.015000	150.010000	49.990000
75%	150.130000	226.245000	75.040000
max	199.980000	299.970000	100.000000

	Total Organic Carban (mg/l)	Nitrate (mg/l)	Phosphate (mg/l) \
count	10000.000000	10000.000000	10000.000000
mean	14.281623	45.162176	2.205561
std	3.309409	25.861234	1.274395
min	2.200000	0.030000	0.000000
25%	12.062534	23.150000	1.090000
50%	14.216595	45.010000	2.200000
75%	16.557652	67.772500	3.320000
max	28.300000	90.000000	4.400000

	Potability
count	10000.000000
mean	0.500000
std	0.500025
min	0.000000
25%	0.000000
50%	0.500000
75%	1.000000
max	1.000000

```
[8]: df.head(10)
```

```
[8]:   Colour (TCU)  Turbidity (NTU)   pH  Conductivity (uS/cm) \
0         8.34         3.39  8.06         819.00
1        14.45         3.36  8.28        1371.10
2         3.87         4.23  6.86         202.75
3        14.57         1.75  7.00         696.16
4         9.01         2.20  6.73         129.24
5         1.84         3.58  7.04          12.63
6         1.36         0.78  6.92        1266.78
7         2.67         1.03  7.46         596.34
8         3.40         1.46  7.29         663.65
9         5.63         4.07  8.41         350.84
```

	Total Dissolved Solids (mg/l)	Total Hardness (mg/l as CaCO ₃) \
0	787.15	279.89
1	779.66	112.04
2	485.10	113.17
3	409.71	140.39
4	343.55	6.52
5	647.39	245.29
6	138.36	149.25
7	82.20	267.11
8	908.88	80.87
9	700.23	213.44

	Aluminium (mg/l)	Chloride (mg/l)	Total Iron (mg/l)	Sodium (mg/l) \
0	0.09	129.30	0.22	13.13
1	0.20	163.73	0.13	127.48
2	0.15	66.68	0.29	142.97
3	0.06	102.42	0.15	194.07
4	0.07	140.47	0.28	3.77
5	0.04	44.71	0.13	86.00
6	0.13	206.45	0.19	150.09
7	0.18	96.54	0.06	183.68
8	0.19	179.05	0.28	190.12
9	0.11	225.84	0.02	47.64

	Sulphate (mg/l)	Biochemical oxygen demand/Zinc (mg/l)	Magnesium (mg/l) \
0	81.01	2.24	12.69
1	307.99	4.05	52.01
2	16.70	0.86	88.47
3	393.09	2.60	61.36
4	170.65	0.04	47.22
5	309.70	1.22	96.65
6	43.26	3.13	54.01
7	3.13	3.40	89.20
8	275.67	2.27	65.64
9	157.70	2.13	57.10

	Calcium (mg/l)	Potassium (mg/l)	Total Organic Carbon (mg/l) \
0	107.95	17.50	10.379783
1	107.12	45.28	15.180013
2	127.47	4.90	16.868637
3	99.16	36.73	18.436525
4	107.17	44.79	11.558279
5	136.71	46.42	8.399735
6	31.84	27.01	13.789695
7	54.77	39.22	12.363817
8	105.11	11.08	12.706049
9	62.28	49.89	17.927806

	Nitrate (mg/l)	Phosphate (mg/l)	Potability
0	22.23	0.41	1
1	16.06	0.68	1
2	19.81	0.91	1
3	42.82	0.02	1
4	14.35	2.08	1
5	15.89	0.47	1
6	2.29	1.54	1
7	0.12	1.28	1
8	41.31	0.88	1
9	32.08	0.87	1

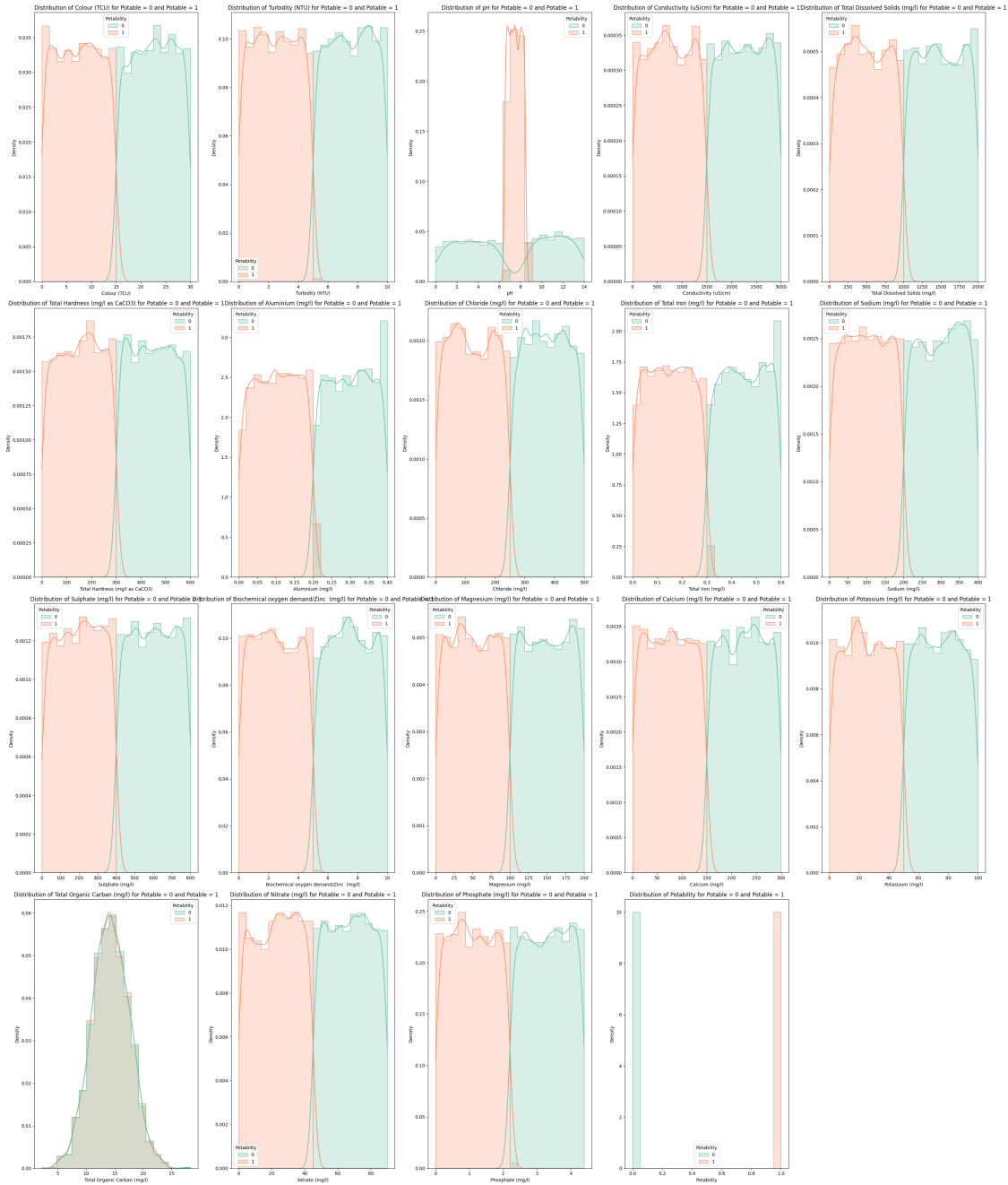
0.2 Data Visualisation

```
[9]: parameters = df.columns
parameters
```

```
[9]: Index(['Colour (TCU)', 'Turbidity (NTU)', 'pH', 'Conductivity (uS/cm)',
        'Total Dissolved Solids (mg/l)', 'Total Hardness (mg/l as CaCO3)',
        'Aluminium (mg/l)', 'Chloride (mg/l)', 'Total Iron (mg/l)',
        'Sodium (mg/l)', 'Sulphate (mg/l)',
        'Biochemical oxygen demand/Zinc (mg/l)', 'Magnesium (mg/l)',
        'Calcium (mg/l)', 'Potassium (mg/l)', 'Total Organic Carbon (mg/l)',
        'Nitrate (mg/l)', 'Phosphate (mg/l)', 'Potability'],
        dtype='object')
```

```
[10]: plt.figure(figsize=(30, 36))

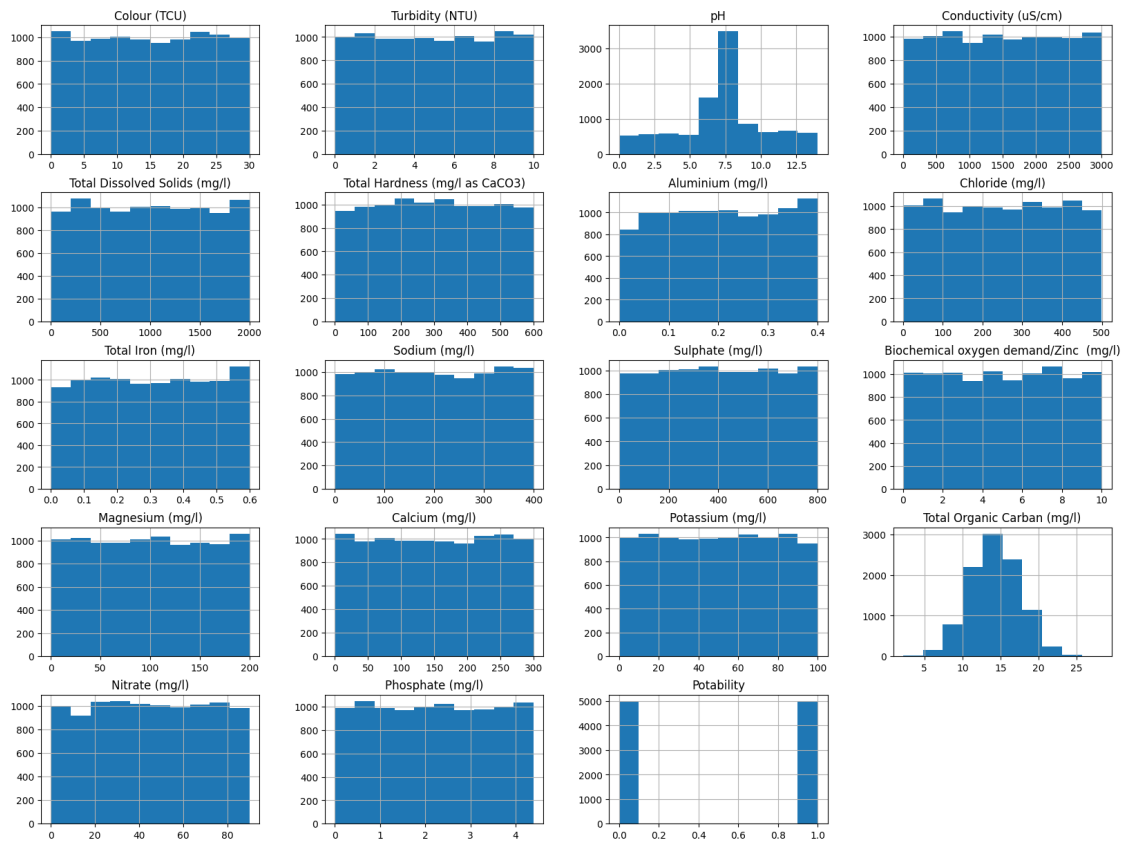
for i, param in enumerate(parameters, 1):
    plt.subplot(4, 5, i)
    sns.histplot(
        data=df,
        x=param,
        hue="Potability",
        kde=True,
        palette="Set2",
        bins=20,
        element="step",
        stat="density",
    )
    plt.title(f"Distribution of {param} for Potable = 0 and Potable = 1")
    plt.xlabel(param)
    plt.ylabel("Density")
plt.tight_layout()
```



```
[11]: df.hist(figsize=(20, 15))
plt.suptitle("Histograms for Each Parameter")
```

```
[11]: Text(0.5, 0.98, 'Histograms for Each Parameter')
```

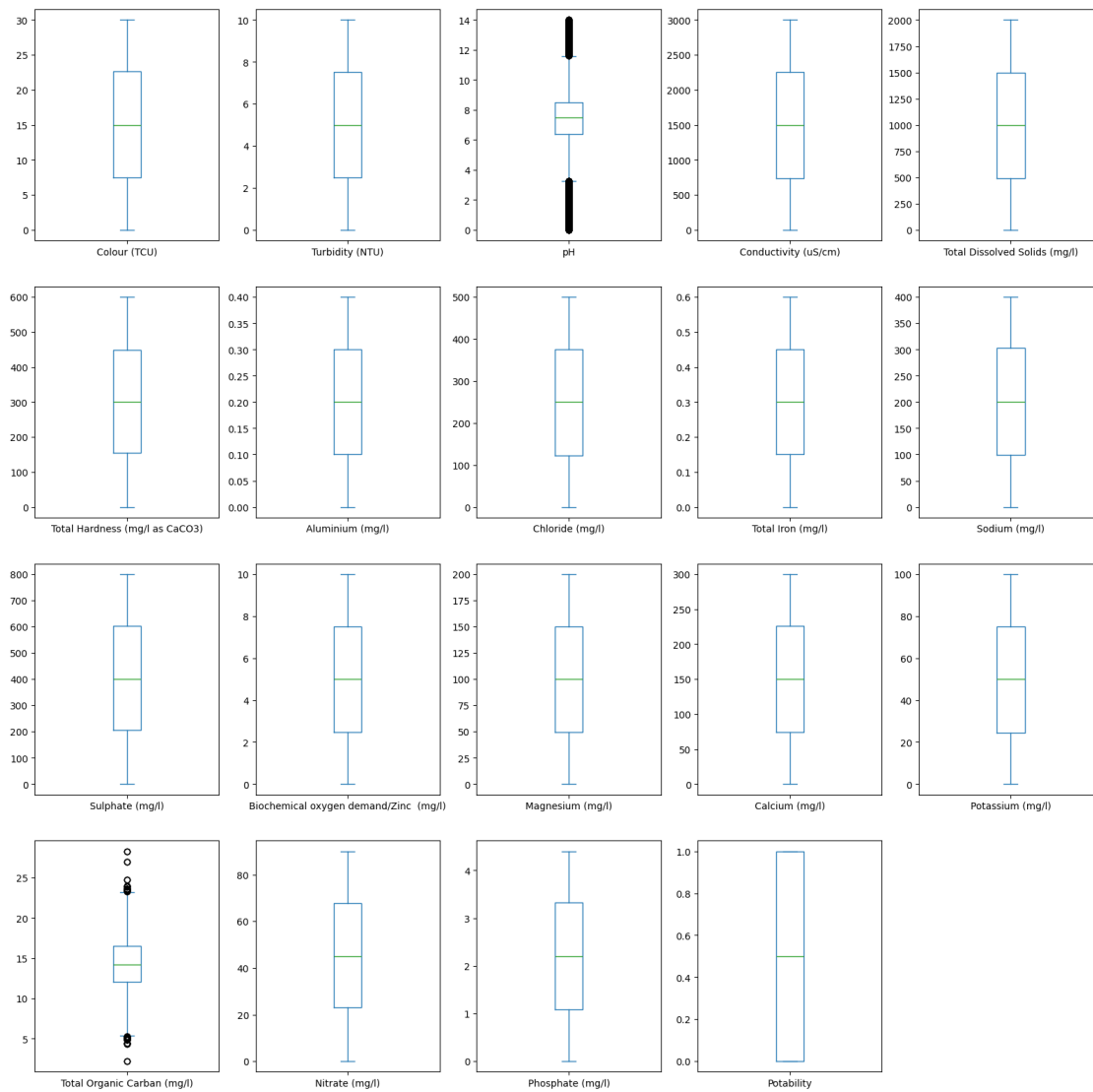
Histograms for Each Parameter



```
[12]: df.plot(
    kind="box",
    subplots=True,
    layout=(4, 5),
    figsize=(20, 20),
    sharex=True,
    sharey=False,
)
plt.suptitle("Boxplots for Each Parameter")
```

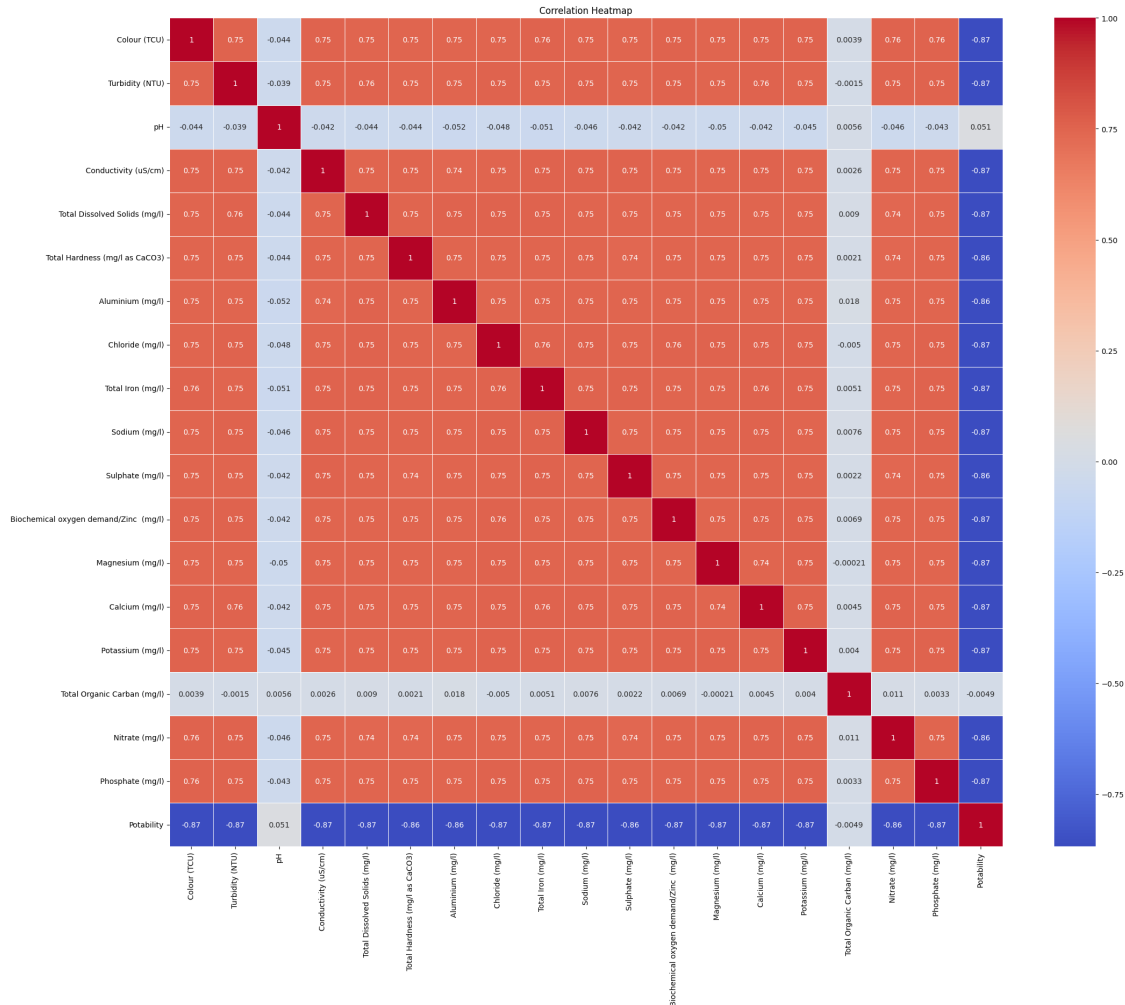
```
[12]: Text(0.5, 0.98, 'Boxplots for Each Parameter')
```


Boxplots for Each Parameter



```
[13]: plt.figure(figsize=(25, 20))
correlation_matrix = df.corr()
sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", linewidths=0.5)
plt.title("Correlation Heatmap")
```

```
[13]: Text(0.5, 1.0, 'Correlation Heatmap')
```



0.3 Finding important parameters

```
[14]: from sklearn.model_selection import train_test_split
      from sklearn.ensemble import RandomForestClassifier
      from sklearn.preprocessing import StandardScaler
```

```
[15]: X = df.drop(columns=["Potability"])
      y = df["Potability"]
```

```
[16]: X_train, X_test, y_train, y_test = train_test_split(
      X, y, test_size=0.2, random_state=42
    )
```

```
[17]: scaler = StandardScaler()
      X_train = scaler.fit_transform(X_train)
      X_test = scaler.transform(X_test)
```

```
[18]: model = RandomForestClassifier(random_state=42)
model.fit(X_train, y_train)
```

```
[18]: RandomForestClassifier(random_state=42)
```

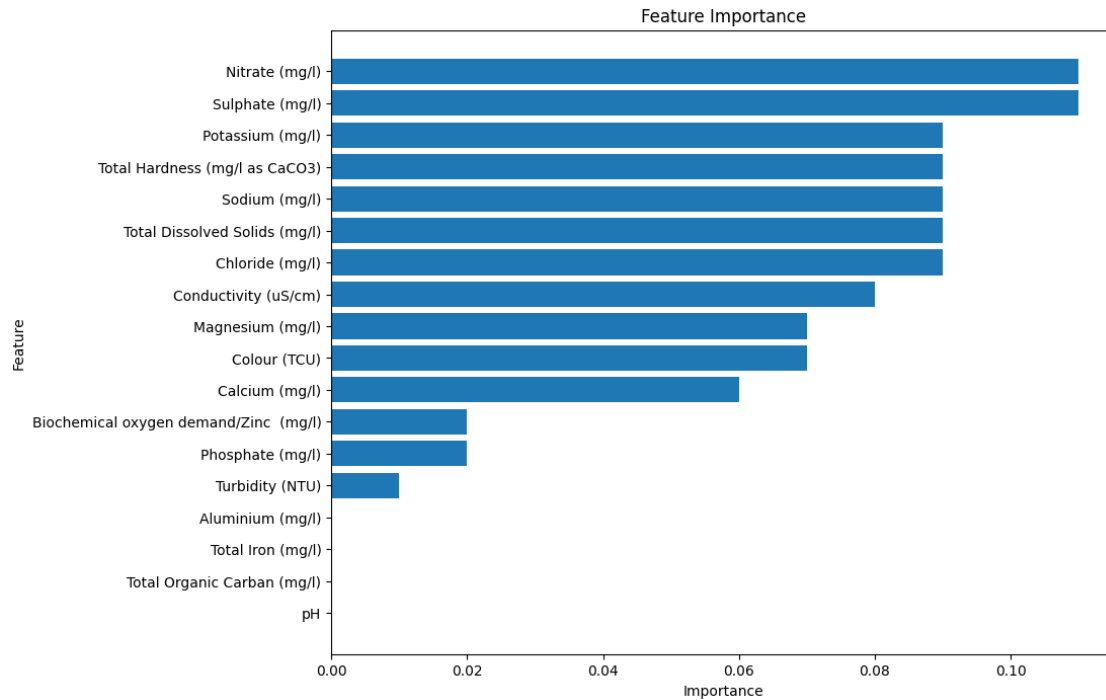
```
[19]: importances = model.feature_importances_
feature_names = X.columns
```

```
[20]: feature_importances = pd.DataFrame(
    {"Feature": feature_names, "Importance": importances}
)
feature_importances = feature_importances.sort_values(by="Importance",
    ↪ascending=False)
feature_importances
```

```
[20]:
```

	Feature	Importance
16	Nitrate (mg/l)	0.110010
10	Sulphate (mg/l)	0.110000
14	Potassium (mg/l)	0.090050
5	Total Hardness (mg/l as CaCO3)	0.090030
9	Sodium (mg/l)	0.090000
4	Total Dissolved Solids (mg/l)	0.090000
7	Chloride (mg/l)	0.090000
3	Conductivity (uS/cm)	0.080000
12	Magnesium (mg/l)	0.070000
0	Colour (TCU)	0.069990
13	Calcium (mg/l)	0.060000
11	Biochemical oxygen demand/Zinc (mg/l)	0.019985
17	Phosphate (mg/l)	0.019930
1	Turbidity (NTU)	0.009990
6	Aluminium (mg/l)	0.000015
8	Total Iron (mg/l)	0.000000
15	Total Organic Carban (mg/l)	0.000000
2	pH	0.000000

```
[21]: plt.figure(figsize=(10, 8))
plt.barh(feature_importances["Feature"], feature_importances["Importance"])
plt.xlabel("Importance")
plt.ylabel("Feature")
plt.title("Feature Importance")
plt.gca().invert_yaxis()
```



0.4 Dairy industry parameters

```
[22]: # getting important parameters
importance_dict = feature_importances.set_index("Feature")["Importance"].
        to_dict()
importance_dict
```

```
[22]: {'Nitrate (mg/l)': 0.110009995227279,
'Sulphate (mg/l)': 0.11000000000000001,
'Potassium (mg/l)': 0.0900499416033473,
'Total Hardness (mg/l as CaCO3)': 0.09002996112543975,
'Sodium (mg/l)': 0.09000000000000001,
'Total Dissolved Solids (mg/l)': 0.09000000000000001,
'Chloride (mg/l)': 0.09000000000000001,
'Conductivity (uS/cm)': 0.08000000000000002,
'Magnesium (mg/l)': 0.07000000000000002,
'Colour (TCU)': 0.06999000244502697,
'Calcium (mg/l)': 0.06000000000000005,
'Biochemical oxygen demand/Zinc (mg/l)': 0.01998500992717573,
'Phosphate (mg/l)': 0.01993009482618602,
'Turbidity (NTU)': 0.009990004772721028,
'Aluminium (mg/l)': 1.4990072824272505e-05,
'Total Iron (mg/l)': 0.0,
'Total Organic Carbon (mg/l)': 0.0,
```

```
'pH': 0.0}
```

```
[23]: thresholds = {
    "Turbidity (NTU)": lambda x: x <= 10,
    "pH": lambda x: 6.5 <= x <= 9.0,
    "Conductivity (uS/cm)": lambda x: x <= 2000,
    "Total Dissolved Solids (mg/l)": lambda x: x <= 1500,
    "Total Hardness (mg/l as CaCO3)": lambda x: x <= 500,
    "Aluminium (mg/l)": lambda x: x < 0.2,
    "Chloride (mg/l)": lambda x: x <= 250,
    "Total Iron (mg/l)": lambda x: x <= 0.3,
    "Sodium (mg/l)": lambda x: x <= 200,
    "Sulphate (mg/l)": lambda x: x <= 250,
    "Biochemical oxygen demand/Zinc (mg/l)": lambda x: x <= 6,
    "Magnesium (mg/l)": lambda x: x <= 50,
    "Calcium (mg/l)": lambda x: x <= 100,
    "Nitrate (mg/l)": lambda x: x <= 50,
    "Phosphate (mg/l)": lambda x: x <= 0.5,
}
```

```
[24]: def is_usable_dairy_industry(row) -> int:
    important_parameters = [param for param, imp in importance_dict.items() if
    ↪imp > 0]
    important_conditions = [
        thresholds[param](row[param])
        for param in important_parameters
        if param in thresholds
    ]
    non_important_conditions = [
        thresholds[param](row[param])
        for param in thresholds
        if param not in important_parameters
    ]

    important_satisfied_percentage = (
        sum(important_conditions) / len(important_conditions)
        if important_conditions
        else 0
    )
    non_important_satisfied_percentage = (
        sum(non_important_conditions) / len(non_important_conditions)
        if non_important_conditions
        else 0
    )

    if (
        important_satisfied_percentage == 1.0
```

```

        and non_important_satisfied_percentage == 1.0
    ):
        return 2
    elif important_satisfied_percentage > 0.9:
        return 1
    else:
        return 0

```

```

[25]: from tensorflow.keras.utils import to_categorical
      from sklearn.model_selection import train_test_split
      from tensorflow.keras.models import Sequential
      from tensorflow.keras.layers import Dense
      from tensorflow.keras.layers import Dropout
      from tensorflow.keras.callbacks import EarlyStopping
      import joblib

```

2024-07-22 16:44:27.943102: I tensorflow/core/util/port.cc:153] oneDNN custom operations are on. You may see slightly different numerical results due to floating-point round-off errors from different computation orders. To turn them off, set the environment variable `TF_ENABLE_ONEDNN_OPTS=0`.

2024-07-22 16:44:27.964108: E external/local_xla/xla/stream_executor/cuda/cuda_fft.cc:485] Unable to register cuFFT factory: Attempting to register factory for plugin cuFFT when one has already been registered

2024-07-22 16:44:27.990380: E external/local_xla/xla/stream_executor/cuda/cuda_dnn.cc:8454] Unable to register cuDNN factory: Attempting to register factory for plugin cuDNN when one has already been registered

2024-07-22 16:44:27.998348: E external/local_xla/xla/stream_executor/cuda/cuda_blas.cc:1452] Unable to register cuBLAS factory: Attempting to register factory for plugin cuBLAS when one has already been registered

2024-07-22 16:44:28.019125: I tensorflow/core/platform/cpu_feature_guard.cc:210] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.
To enable the following instructions: AVX2 AVX_VNNI FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.

2024-07-22 16:44:29.257382: W tensorflow/compiler/tf2tensorrt/utils/py_utils.cc:38] TF-TRT Warning: Could not find TensorRT

```

[26]: df["Usable by dairy industry"] = df.apply(is_usable_dairy_industry, axis=1).
      ↪astype(int)

```

```

[27]: df["Usable by dairy industry"].unique()

```

```

[27]: array([1, 0, 2])

```

```
[28]: X_dairy = df.drop(
        [
            "Usable by dairy industry",
            "Potability",
        ],
        axis=1,
    )
y_dairy = to_categorical(df["Usable by dairy industry"])

X_train_dairy, X_test_dairy, y_train_dairy, y_test_dairy = train_test_split(
    X_dairy, y_dairy, test_size=0.3, random_state=42
)

[29]: scaler = StandardScaler()
X_train_dairy_scaled = scaler.fit_transform(X_train_dairy)
X_test_dairy_scaled = scaler.transform(X_test_dairy)

[30]: X_train_dairy_scaled = pd.DataFrame(X_train_dairy_scaled, columns=X_dairy.
        ↪columns)
X_test_dairy_scaled = pd.DataFrame(X_test_dairy_scaled, columns=X_dairy.columns)

[31]: model = Sequential()
model.add(Dense(64, input_dim=X_train_dairy_scaled.shape[1], activation="relu"))
model.add(Dropout(0.5))
model.add(Dense(32, activation="relu"))
model.add(Dropout(0.5))
model.add(Dense(3, activation="softmax"))

model.compile(optimizer="adam", loss="categorical_crossentropy",
        ↪metrics=["accuracy"])

/home/funinkina/.local/lib/python3.12/site-
packages/keras/src/layers/core/dense.py:87: UserWarning: Do not pass an
`input_shape`/`input_dim` argument to a layer. When using Sequential models,
prefer using an `Input(shape)` object as the first layer in the model instead.
  super().__init__(activity_regularizer=activity_regularizer, **kwargs)

[32]: early_stopping = EarlyStopping(
        monitor="val_loss", patience=10, restore_best_weights=True
    )

[33]: model.summary()
```

Model: "sequential"

Layer (type)

Output Shape

Param #

dense (Dense)	(None, 64)	1,216
dropout (Dropout)	(None, 64)	0
dense_1 (Dense)	(None, 32)	2,080
dropout_1 (Dropout)	(None, 32)	0
dense_2 (Dense)	(None, 3)	99

Total params: 3,395 (13.26 KB)

Trainable params: 3,395 (13.26 KB)

Non-trainable params: 0 (0.00 B)

```
[34]: history = model.fit(
    X_train_dairy_scaled,
    y_train_dairy,
    epochs=100,
    batch_size=10,
    validation_data=(X_test_dairy_scaled, y_test_dairy),
    callbacks=[early_stopping], # Assuming early_stopping is defined
)
```

Epoch 1/100

700/700 3s 3ms/step -
accuracy: 0.7248 - loss: 0.6747 - val_accuracy: 0.8393 - val_loss: 0.3504

Epoch 2/100

700/700 2s 3ms/step -
accuracy: 0.8465 - loss: 0.3517 - val_accuracy: 0.8420 - val_loss: 0.3117

Epoch 3/100

700/700 2s 3ms/step -
accuracy: 0.8684 - loss: 0.2978 - val_accuracy: 0.8567 - val_loss: 0.2935

Epoch 4/100

700/700 2s 3ms/step -
accuracy: 0.8683 - loss: 0.2957 - val_accuracy: 0.8697 - val_loss: 0.2822

Epoch 5/100

700/700 2s 3ms/step -
accuracy: 0.8820 - loss: 0.2706 - val_accuracy: 0.8790 - val_loss: 0.2647

Epoch 6/100

700/700 2s 3ms/step -
accuracy: 0.8856 - loss: 0.2567 - val_accuracy: 0.8820 - val_loss: 0.2610

Epoch 7/100
700/700 2s 3ms/step -
accuracy: 0.8844 - loss: 0.2567 - val_accuracy: 0.8803 - val_loss: 0.2495

Epoch 8/100
700/700 2s 3ms/step -
accuracy: 0.8836 - loss: 0.2584 - val_accuracy: 0.8863 - val_loss: 0.2475

Epoch 9/100
700/700 2s 3ms/step -
accuracy: 0.8941 - loss: 0.2395 - val_accuracy: 0.8897 - val_loss: 0.2385

Epoch 10/100
700/700 2s 3ms/step -
accuracy: 0.8997 - loss: 0.2285 - val_accuracy: 0.8910 - val_loss: 0.2292

Epoch 11/100
700/700 2s 3ms/step -
accuracy: 0.8971 - loss: 0.2304 - val_accuracy: 0.8923 - val_loss: 0.2261

Epoch 12/100
700/700 2s 3ms/step -
accuracy: 0.9043 - loss: 0.2168 - val_accuracy: 0.8943 - val_loss: 0.2211

Epoch 13/100
700/700 2s 3ms/step -
accuracy: 0.9047 - loss: 0.2181 - val_accuracy: 0.8990 - val_loss: 0.2091

Epoch 14/100
700/700 2s 3ms/step -
accuracy: 0.9012 - loss: 0.2177 - val_accuracy: 0.9000 - val_loss: 0.2071

Epoch 15/100
700/700 2s 3ms/step -
accuracy: 0.9074 - loss: 0.2081 - val_accuracy: 0.9017 - val_loss: 0.2096

Epoch 16/100
700/700 2s 3ms/step -
accuracy: 0.9094 - loss: 0.2029 - val_accuracy: 0.9063 - val_loss: 0.2021

Epoch 17/100
700/700 2s 3ms/step -
accuracy: 0.9138 - loss: 0.1937 - val_accuracy: 0.9057 - val_loss: 0.1951

Epoch 18/100
700/700 2s 3ms/step -
accuracy: 0.9108 - loss: 0.2000 - val_accuracy: 0.9097 - val_loss: 0.2008

Epoch 19/100
700/700 2s 3ms/step -
accuracy: 0.9096 - loss: 0.2030 - val_accuracy: 0.9067 - val_loss: 0.2055

Epoch 20/100
700/700 2s 3ms/step -
accuracy: 0.9077 - loss: 0.1956 - val_accuracy: 0.9127 - val_loss: 0.1858

Epoch 21/100
700/700 2s 3ms/step -
accuracy: 0.9090 - loss: 0.1914 - val_accuracy: 0.9080 - val_loss: 0.1991

Epoch 22/100
700/700 2s 3ms/step -
accuracy: 0.9161 - loss: 0.1827 - val_accuracy: 0.9090 - val_loss: 0.1870

Epoch 23/100
700/700 2s 3ms/step -
accuracy: 0.9151 - loss: 0.1951 - val_accuracy: 0.9243 - val_loss: 0.1811

Epoch 24/100
700/700 2s 3ms/step -
accuracy: 0.9106 - loss: 0.1961 - val_accuracy: 0.9190 - val_loss: 0.1764

Epoch 25/100
700/700 2s 3ms/step -
accuracy: 0.9255 - loss: 0.1727 - val_accuracy: 0.9163 - val_loss: 0.1819

Epoch 26/100
700/700 2s 3ms/step -
accuracy: 0.9150 - loss: 0.1893 - val_accuracy: 0.9183 - val_loss: 0.1815

Epoch 27/100
700/700 2s 3ms/step -
accuracy: 0.9190 - loss: 0.1785 - val_accuracy: 0.9200 - val_loss: 0.1770

Epoch 28/100
700/700 2s 3ms/step -
accuracy: 0.9252 - loss: 0.1800 - val_accuracy: 0.9217 - val_loss: 0.1722

Epoch 29/100
700/700 2s 3ms/step -
accuracy: 0.9279 - loss: 0.1671 - val_accuracy: 0.9273 - val_loss: 0.1659

Epoch 30/100
700/700 2s 3ms/step -
accuracy: 0.9217 - loss: 0.1682 - val_accuracy: 0.9227 - val_loss: 0.1753

Epoch 31/100
700/700 2s 3ms/step -
accuracy: 0.9184 - loss: 0.1891 - val_accuracy: 0.9207 - val_loss: 0.1650

Epoch 32/100
700/700 2s 3ms/step -
accuracy: 0.9317 - loss: 0.1716 - val_accuracy: 0.9183 - val_loss: 0.1692

Epoch 33/100
700/700 2s 3ms/step -
accuracy: 0.9193 - loss: 0.1755 - val_accuracy: 0.9193 - val_loss: 0.1785

Epoch 34/100
700/700 2s 3ms/step -
accuracy: 0.9320 - loss: 0.1722 - val_accuracy: 0.9283 - val_loss: 0.1745

Epoch 35/100
700/700 2s 3ms/step -
accuracy: 0.9303 - loss: 0.1618 - val_accuracy: 0.9243 - val_loss: 0.1667

Epoch 36/100
700/700 2s 3ms/step -
accuracy: 0.9284 - loss: 0.1603 - val_accuracy: 0.9330 - val_loss: 0.1639

Epoch 37/100
700/700 2s 3ms/step -
accuracy: 0.9362 - loss: 0.1569 - val_accuracy: 0.9360 - val_loss: 0.1502

Epoch 38/100
700/700 2s 3ms/step -
accuracy: 0.9373 - loss: 0.1470 - val_accuracy: 0.9320 - val_loss: 0.1567

Epoch 39/100
700/700 2s 3ms/step -
accuracy: 0.9327 - loss: 0.1573 - val_accuracy: 0.9263 - val_loss: 0.1675
Epoch 40/100
700/700 2s 3ms/step -
accuracy: 0.9301 - loss: 0.1575 - val_accuracy: 0.9333 - val_loss: 0.1519
Epoch 41/100
700/700 2s 3ms/step -
accuracy: 0.9365 - loss: 0.1580 - val_accuracy: 0.9323 - val_loss: 0.1546
Epoch 42/100
700/700 2s 3ms/step -
accuracy: 0.9287 - loss: 0.1603 - val_accuracy: 0.9307 - val_loss: 0.1583
Epoch 43/100
700/700 2s 3ms/step -
accuracy: 0.9311 - loss: 0.1544 - val_accuracy: 0.9307 - val_loss: 0.1473
Epoch 44/100
700/700 2s 3ms/step -
accuracy: 0.9325 - loss: 0.1627 - val_accuracy: 0.9323 - val_loss: 0.1535
Epoch 45/100
700/700 2s 3ms/step -
accuracy: 0.9350 - loss: 0.1561 - val_accuracy: 0.9290 - val_loss: 0.1602
Epoch 46/100
700/700 2s 3ms/step -
accuracy: 0.9380 - loss: 0.1542 - val_accuracy: 0.9307 - val_loss: 0.1624
Epoch 47/100
700/700 2s 3ms/step -
accuracy: 0.9269 - loss: 0.1578 - val_accuracy: 0.9383 - val_loss: 0.1406
Epoch 48/100
700/700 2s 3ms/step -
accuracy: 0.9360 - loss: 0.1507 - val_accuracy: 0.9327 - val_loss: 0.1485
Epoch 49/100
700/700 2s 3ms/step -
accuracy: 0.9342 - loss: 0.1570 - val_accuracy: 0.9413 - val_loss: 0.1407
Epoch 50/100
700/700 2s 3ms/step -
accuracy: 0.9355 - loss: 0.1486 - val_accuracy: 0.9383 - val_loss: 0.1472
Epoch 51/100
700/700 2s 3ms/step -
accuracy: 0.9386 - loss: 0.1461 - val_accuracy: 0.9347 - val_loss: 0.1606
Epoch 52/100
700/700 2s 3ms/step -
accuracy: 0.9358 - loss: 0.1533 - val_accuracy: 0.9357 - val_loss: 0.1522
Epoch 53/100
700/700 2s 3ms/step -
accuracy: 0.9377 - loss: 0.1534 - val_accuracy: 0.9437 - val_loss: 0.1461
Epoch 54/100
700/700 2s 3ms/step -
accuracy: 0.9357 - loss: 0.1592 - val_accuracy: 0.9420 - val_loss: 0.1500

```
Epoch 55/100
700/700          2s 2ms/step -
accuracy: 0.9411 - loss: 0.1448 - val_accuracy: 0.9430 - val_loss: 0.1428
Epoch 56/100
700/700          2s 3ms/step -
accuracy: 0.9448 - loss: 0.1366 - val_accuracy: 0.9400 - val_loss: 0.1427
Epoch 57/100
700/700          2s 3ms/step -
accuracy: 0.9436 - loss: 0.1358 - val_accuracy: 0.9413 - val_loss: 0.1455
```

```
[35]: loss, accuracy = model.evaluate(X_test_dairy_scaled, y_test_dairy)
      print(f"Test Loss: {loss:.4f}")
      print(f"Test Accuracy: {accuracy:.4f}")
```

```
94/94            0s 2ms/step -
accuracy: 0.9390 - loss: 0.1389
Test Loss: 0.1406
Test Accuracy: 0.9383
```

```
[36]: predictions = model.predict(X_test_dairy_scaled)
      predictions
```

```
94/94            0s 2ms/step
```

```
[36]: array([[1.00000000e+00, 0.00000000e+00, 0.00000000e+00],
             [9.7387344e-01, 2.6126536e-02, 4.5396979e-13],
             [9.9999869e-01, 1.3623968e-06, 0.00000000e+00],
             ...,
             [9.9999964e-01, 3.2249073e-07, 0.00000000e+00],
             [1.00000000e+00, 0.00000000e+00, 0.00000000e+00],
             [9.5503998e-01, 4.4959992e-02, 3.7115763e-10]], dtype=float32)
```

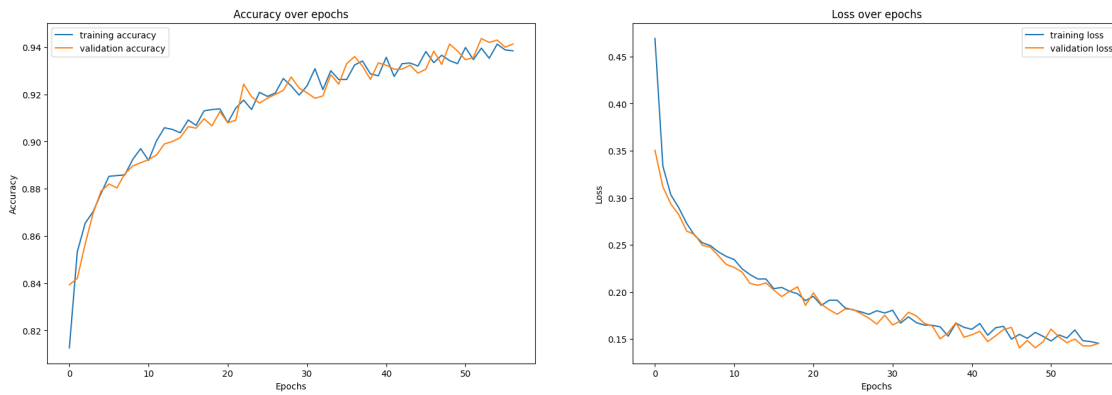
```
[37]: plt.figure(figsize=(22, 7))

      plt.subplot(1, 2, 1)
      plt.plot(history.history["accuracy"], label="training accuracy")
      plt.plot(history.history["val_accuracy"], label="validation accuracy")
      plt.xlabel("Epochs")
      plt.ylabel("Accuracy")
      plt.title("Accuracy over epochs")
      plt.legend()

      plt.subplot(1, 2, 2)
      plt.plot(history.history["loss"], label="training loss")
      plt.plot(history.history["val_loss"], label="validation loss")
      plt.xlabel("Epochs")
      plt.ylabel("Loss")
      plt.title("Loss over epochs")
```

```
plt.legend()
```

```
[37]: <matplotlib.legend.Legend at 0x7a53783301d0>
```



0.5 Saving the model

```
[38]: joblib.dump(scaler, "scaler.save")
```

```
[38]: ['scaler.save']
```

```
[39]: model.save("water_usable_by_dairy_industry.keras")
```