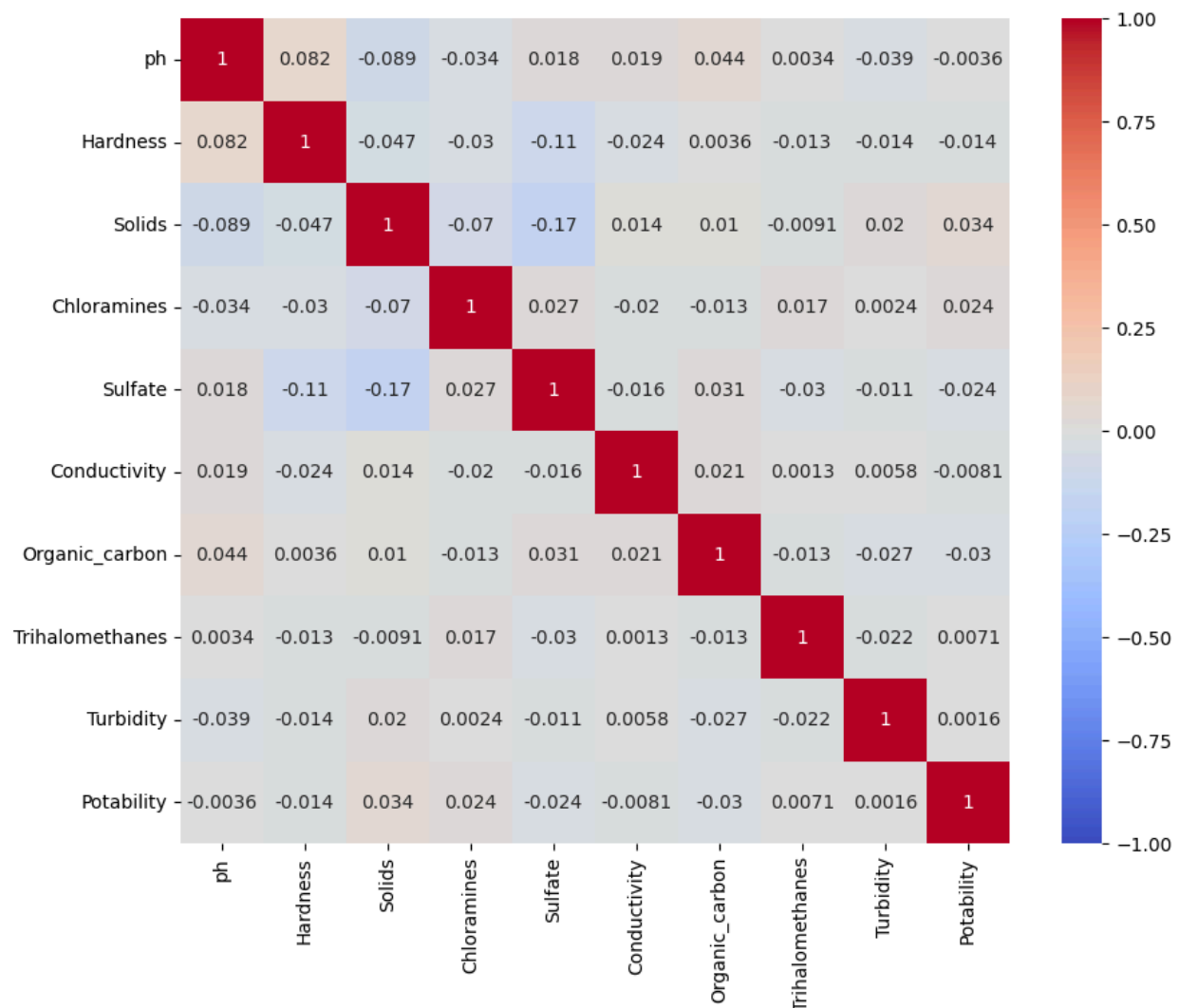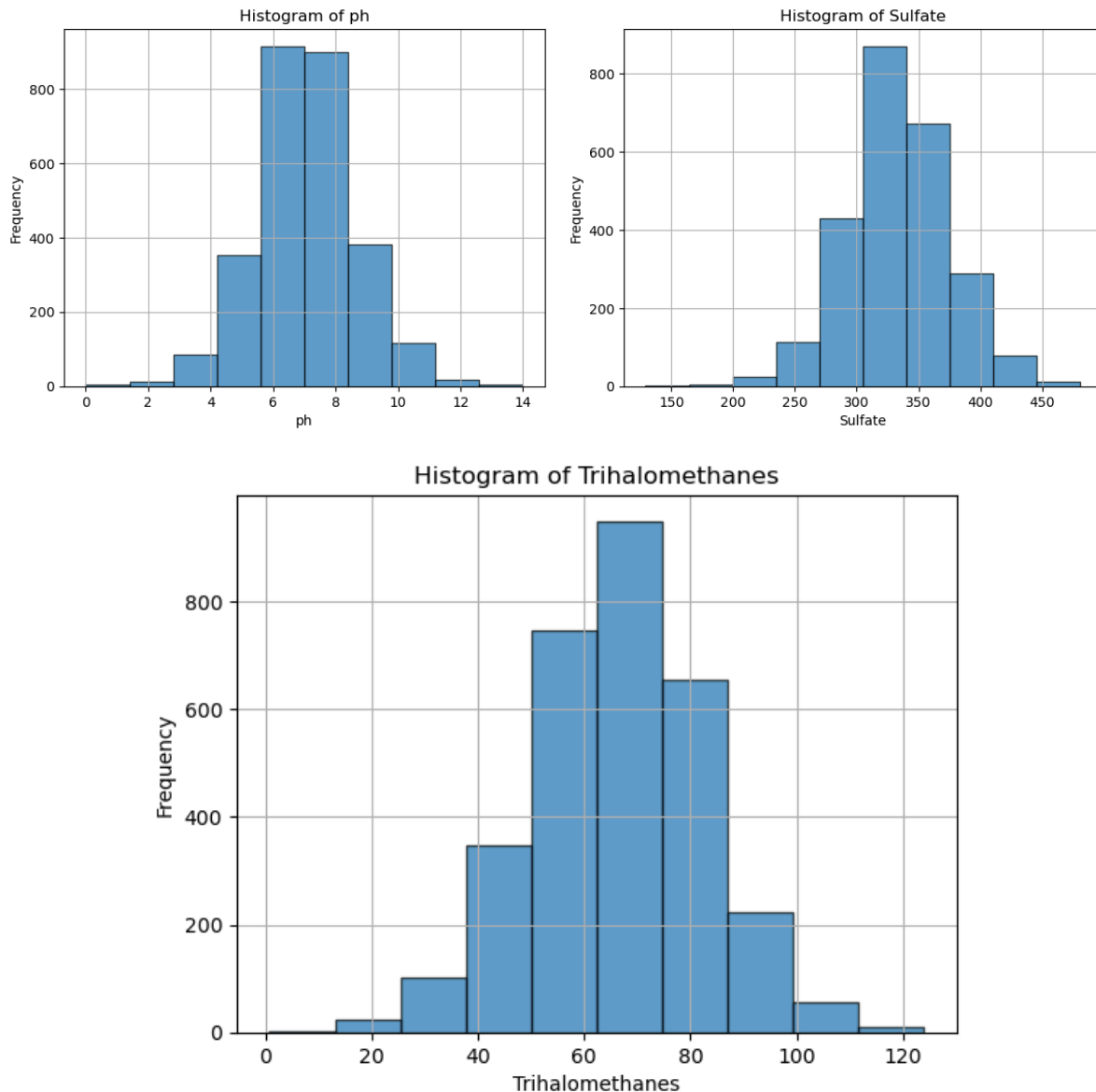# Assignment: 1

**MLFlow implementation:**

1. **Data Collection & Preprocessing:**

● We have obtained water_potability data from kaggle([water_potability kaggle data](#)). In that data there are various features (ph, Hardness, Solids, Chloramines, Sulfate, Conductivity, Organic_carbon, Trihalomethanes and Turbidity) that determined that the water is drinkable or not that shows in potability column.

● So, first determine that which kind of the data in the column. So using the data.info() and data.describe() we can determine the basic things like statistics and data types. So, Here in the data set all independent features are in float data type. And dependent feature is integer data type.

● And then determine the heatmap of correlation matrix but we cannot find anything specific in that. All features are uniquely distinguish.

● Here plot the histogram for the feature who contains the null value.







● After that find the Null, NA and Nan value in the column and replace by mean or median. Here data is normally distributed(As we can see in the above plots of histogram) for feature ph,Sulfate and Trihalomethanes so we can replace the null or na value with the mean.
● And to scale the feature with outliers I've used MinMaxScaler.

## 2. Model Implementation:

**Ridge Regression:**

- Ridge Regression extends linear regression to improve model robustness. It introduces L2 regularization, adding a penalty based on the sum of squared coefficients (scaled by the hyperparameter "alpha").
- By adjusting alpha, we control the trade-off between fitting the data well and keeping coefficients small. Ridge is particularly useful for handling multicollinearity and stabilizing vanilla linear regression against outliers and overfitting.

**Lasso Regression:**

- Lasso Regression is another regularization technique which I've applied. It encourages sparse models by using L1 regularization.
- The L1 term penalizes absolute coefficient values, leading some coefficients to become exactly zero (feature selection). Lasso simplifies models and prevents overfitting, similar to Ridge.

**Result Table:**

| Alpha Value | Ridge Regression | | Lasso Regression | |
|---|---|---|---|---|
| | **RMSE** | **R-squared** | **RMSE** | **R-squared** |
| 1 | 0.4833 | 0.0001 | 0.4839 | -0.0026 |
| **10** | 0.4833 | 0.0001 | 0.4839 | -0.0025 |
| **100** | 0.4833 | 0.0001 | 0.4838 | -0.0020 |
| **1000** | 0.4833 | -0.0001 | 0.4839 | -0.0022 |

## 3. MLflow Integration:

**Dashboard:**

**Model metrics:**



**Model 1 metrics (Alpha = 1000):**

**Model 2 metrics (alpha = 100):**

# Model 3 metrics (alpha = 10)

**Model 4 metrics (Alpha = 1):**