# Report on Application of Machine Learning on Sonar Rock vs Mine Dataset

**Date- 25 June 2023**

**Student's Name- Vedant Roy**

**Course Title- Introduction to Machine Learning**

**Course Teacher- Mr. Amit Choudhary**

**Enrollment Number- 12019011621**

**Email ID- vedantroy3@gmail.com**

**Contact Number- 7053969950**

**Link of Google Colab (which has codes of the dataset)-**
**https://colab.research.google.com/drive/15TX3rlVsHEX5zr1Zzo3jQzwdew0scpZH?usp=sharing**

**Link of the dataset-**
**https://drive.google.com/file/d/1g_mQpDVpXcgoHIy25nGK6TQqOsB61Cdb/view?usp=sharing**

**Link of the Video-**
**https://drive.google.com/file/d/17FLFnXWflQVVFbQyFxWkpXMqMJGEwzSu/view?usp=sharing**

# Introduction

The Sonar rock and mine dataset is a valuable resource for exploring the application of machine learning algorithms in the field of underwater object detection and classification. Sonar technology plays a crucial role in various domains such as marine exploration, defense, and underwater robotics. The dataset provides a collection of sonar signals captured from different objects, specifically rocks and mines, along with corresponding class labels.

The goal of this project is to deploy and evaluate multiple machine learning models on the Sonar rock and mine dataset. By leveraging these models, we aim to classify sonar signals accurately and predict the presence of either rocks or mines. The models employed in this project encompass a diverse range of algorithms, including K-Nearest Neighbors (KNN), Decision Tree, Random Forest, Support Vector Machine (SVM), Gradient Boosting Classifier (XGBoost), Linear Regression, and Random Forest Regression.

The primary objective of this project is to assess the performance of these models and gain insights into their strengths and weaknesses for sonar signal classification and prediction tasks. By applying these models on the dataset, we can analyze their accuracy and determine their suitability for practical deployment in real-world scenarios.

In this report, we will present a detailed overview of the model deployment process. We will discuss the steps involved, including data preprocessing, model selection, training, evaluation, and result analysis. Additionally, we will provide a comprehensive assessment of each model's performance and draw conclusions based on the obtained results.

The insights and findings from this project can contribute to advancements in underwater object detection systems and enhance the understanding of sonar signal analysis. Such advancements have significant implications in areas like underwater exploration, surveillance, and safety.

Let's proceed further with the project report

# Methodology

The methodology employed in this project involves a systematic approach to applying various machine learning models on the Sonar rock and mine dataset. The key steps encompass data loading, preprocessing, model selection, training, evaluation, and result analysis. The following sections provide a detailed description of each step.

**1. Data Loading:**
The first step is to load the Sonar rock and mine dataset. The dataset is obtained from the provided link and imported using the pandas library. It is stored in a pandas DataFrame for further processing.

**2. Data Preprocessing:**
Before applying the machine learning models, it is crucial to preprocess the dataset. In this step, the features and labels are separated from the DataFrame. The feature matrix, denoted as X, consists of all the columns except the last one, which represents the class labels. The label vector, denoted as y, contains the class labels.

To ensure compatibility with the models, the class labels are encoded using a LabelEncoder. This conversion transforms the categorical labels ("M" for mine and "R" for rock) into numerical values (0 and 1, respectively). The label encoder is applied to the y vector, replacing the categorical labels with their corresponding numerical representations.

**3. Data Splitting:**
To assess the performance of the models, it is essential to split the dataset into training and testing sets. This step is carried out using the train_test_split function from the sklearn.model_selection module. The dataset is divided into 80% training data and 20% testing data, ensuring a random and representative distribution of samples. The split is performed on both the feature matrix (X) and the label vector (y).

**4. Model Selection and Training:**
In this step, a variety of classification and regression models are selected for deployment on the dataset. For the classification task, the chosen models include K-Nearest Neighbors (KNN), Decision Tree, Random Forest, Support Vector Machine (SVM), and Gradient Boosting Classifier (XGBoost). For the

regression task, Linear Regression and Random Forest Regression models are employed.

Each model is instantiated using its respective class from the corresponding sklearn module. The training data (X_train and y_train) are then used to train each model using the fit method. This step involves learning the underlying patterns and relationships between the features and labels in the training set.

## 5. Model Evaluation:
After training the models, their performance is evaluated using the testing data (X_test and y_test). For the classification models, accuracy is chosen as the evaluation metric, which measures the percentage of correctly classified instances. For the regression models, the coefficient of determination (R-squared) is used as the evaluation metric, indicating the proportion of the variance in the target variable that is predictable from the input variables.

The accuracy and R-squared values are computed using the score method of each trained model, passing the testing data as input.

## 6. Result Analysis:
The final step involves analyzing the results obtained from the model evaluation. The classification model accuracies and regression model accuracies are displayed to provide a comprehensive overview of the performance of each model.

The analysis includes comparing the accuracies of different models to identify the most effective ones for the given dataset. By considering the strengths and weaknesses of each model, we can gain insights into their suitability for the task of sonar signal classification and prediction.

The methodology outlined above provides a structured approach to the deployment of machine learning models on the Sonar rock and mine dataset. It ensures the appropriate handling of data, model selection, training, evaluation, and result analysis, enabling a comprehensive assessment of the models' performance.

# Model Deployment

In this section, we will discuss the deployment of various machine learning models on the Sonar rock and mine dataset. We have utilized the scikit-learn library to implement these models. The code provided below outlines the steps involved in deploying the models.

## Step 1: Importing the necessary libraries

We begin by importing the required libraries for our project. The pandas library is used for data manipulation and analysis, and we import specific modules from scikit-learn for different model implementations. The models we will be deploying include K-Nearest Neighbors (KNN), Decision Tree, Random Forest, Support Vector Machine (SVM), Gradient Boosting Classifier (XGBoost), Linear Regression, and Random Forest Regression.

## Step 2: Load the dataset

We load the Sonar rock and mine dataset using the provided link. The dataset contains various attributes describing sonar signals, and the target variable represents the class labels (rocks or mines). We use the pandas library to read the CSV file and store it in a DataFrame named `df`.

## Step 3: Preprocess the dataset

In this step, we preprocess the dataset by separating the features and labels. We assign the feature matrix to the variable `X`, which includes all columns except the last one. The labels are assigned to the variable `y`, representing the last column of the dataset.

To prepare the labels for model training, we use the LabelEncoder class from scikit-learn. It converts the categorical labels "M" and "R" into numerical values.

## Step 4: Split the dataset

To evaluate the performance of our models, we need to split the dataset into training and testing sets. We use the `train_test_split` function from scikit-learn, which randomly shuffles and divides the data into two portions. The training set comprises 80% of the data, and the remaining 20% is assigned to

the testing set. The feature matrices and label vectors for both sets are stored in `X_train`, `X_test`, `y_train`, and `y_test`, respectively.

## Step 5: Apply the classification models

We proceed to apply five classification models on the dataset.

### Model 1: K-Nearest Neighbors (KNN)
  - We instantiate an object of the `KNeighborsClassifier` class.
  - The model is trained on the training set using the `fit` method.

### Model 2: Decision Tree
  - We create an instance of the `DecisionTreeClassifier` class.
  - The model is trained on the training set using the `fit` method.

### Model 3: Random Forest
  - We initialize an object of the `RandomForestClassifier` class.
  - The model is trained on the training set using the `fit` method.

### Model 4: Support Vector Machine (SVM)
  - We create an instance of the `SVC` class.
  - The model is trained on the training set using the `fit` method.

### Model 5: Gradient Boosting Classifier (XGBoost)
  - We instantiate an object of the `XGBClassifier` class.
  - The model is trained on the training set using the `fit` method.

## Step 6: Apply the regression models

We also apply two regression models on the dataset.

### Model 1: Linear Regression
  - We create an instance of the `LinearRegression` class.
  - The model is trained on the training set using the `fit` method.

### Model 2: Random Forest Regression
  - We initialize an object of the `RandomForestRegressor` class.
  - The model is trained on the training set using

 the `fit` method.

## Step 7: Evaluate the models

To assess the performance of the models, we evaluate their accuracies using appropriate evaluation metrics.

For classification models, we calculate the accuracy scores using the `score` method on the respective testing sets. The accuracy values for each classification model are stored in variables `knn_accuracy`, `dt_accuracy`, `rf_accuracy`, `svm_accuracy`, and `xgb_accuracy`.

For regression models, we calculate the coefficient of determination (R-squared) using the `score` method on the testing sets. The accuracy values for each regression model are stored in variables `linear_reg_accuracy` and `rf_reg_accuracy`.

## Step 8: Display the results

Finally, we display the accuracies of the classification and regression models using print statements.

The classification model accuracies are printed, including K-Nearest Neighbors (KNN), Decision Tree, Random Forest, Support Vector Machine (SVM), and Gradient Boosting Classifier (XGBoost).

The regression model accuracies are also printed, including Linear Regression and Random Forest Regression.

These accuracies provide insights into the performance of each model on the Sonar rock and mine dataset.

By deploying these models and evaluating their accuracies, we gain valuable information about their effectiveness in classifying and predicting the target variable.

# Results and Discussion

The results obtained from applying various machine learning models on the Sonar rock and mine dataset are presented in terms of classification and regression accuracies. The classification model accuracies are as follows:

- K-Nearest Neighbors (KNN) Accuracy: 0.8571
- Decision Tree Accuracy: 0.6429
- Random Forest Accuracy: 0.8571
- Support Vector Machine (SVM) Accuracy: 0.8333
- Gradient Boosting Classifier (XGBoost) Accuracy: 0.8095

On the other hand, the regression model accuracies are as follows:

- Linear Regression Accuracy: 0.2822
- Random Forest Regression Accuracy: 0.3596

In the classification task, K-Nearest Neighbors (KNN), Random Forest, and Support Vector Machine (SVM) achieved relatively high accuracies, with KNN and Random Forest both achieving an accuracy of 0.8571. These models demonstrated good performance in classifying the sonar signals as either rocks or mines. However, the Decision Tree and Gradient Boosting Classifier (XGBoost) models achieved slightly lower accuracies of 0.6429 and 0.8095, respectively.

In the regression task, both Linear Regression and Random Forest Regression achieved lower accuracies compared to the classification models. Linear Regression attained an accuracy of 0.2822, while Random Forest Regression obtained an accuracy of 0.3596. These results suggest that the regression models had limited predictive power in estimating the continuous values related to the sonar signals.

The varying accuracies among the models can be attributed to the differences in their underlying algorithms and modeling techniques. Models like KNN, Random Forest, and SVM are known for their ability to capture complex patterns and relationships in the data, which contributed to their higher classification accuracies. In contrast, the linear regression-based models may have struggled to capture the nonlinear nature of the dataset, leading to lower regression accuracies.

# Conclusion

In conclusion, this project aimed to deploy various machine learning models on the Sonar rock and mine dataset for classification and regression tasks. The results obtained demonstrated the varying performance of the models in accurately classifying the sonar signals and predicting continuous values.

Among the classification models, K-Nearest Neighbors (KNN) and Random Forest exhibited the highest accuracy of 0.8571, indicating their effectiveness in distinguishing between rocks and mines based on the sonar signals. Support Vector Machine (SVM) also performed well with an accuracy of 0.8333. However, the Decision Tree and Gradient Boosting Classifier (XGBoost) models achieved slightly lower accuracies of 0.6429 and 0.8095, respectively.

On the other hand, the regression models, including Linear Regression and Random Forest Regression, yielded lower accuracies compared to the classification models. This suggests that the regression models had limited predictive power in estimating the continuous values associated with the sonar signals.

Overall, the results highlight the importance of selecting appropriate machine learning models based on the nature of the dataset and the specific task at hand. Further exploration and optimization of the models may be necessary to improve their performance and address the challenges associated with the Sonar rock and mine dataset.