# MILESTONE FOUR

## 1. Introduction:

The goal was to develop a predictive model capable of accurately classifying instances in the dataset. To achieve this objective, we explored the performance of a Dummy Classifier as a baseline, followed by more sophisticated models including Gradient Boosting and Support Vector Machines (SVM). The evaluation was conducted using various metrics such as accuracy, precision, recall, and F1-score. Additionally, confusion matrices were employed to provide insights into the model's behavior.

Instead of using the RandomForestClassifier and Logistic Regression as in the previous milestone, I'll train a DummyClassifier as a baseline model to see how it performs compared to the previous models. This will provide a benchmark for comparison and help evaluate the effectiveness of more advanced models.

## 2. Data Preprocessing:

- **Feature Engineering**: The program demonstrates several feature engineering steps aimed at enhancing the predictive power of the features. Here's a breakdown of the feature engineering techniques used:

    - **Polynomial Feature Transformation**:

        The program utilizes the PolynomialFeatures transformer from scikit-learn to generate polynomial features up to the specified degree. Polynomial features are created by multiplying different combinations of existing features, allowing the model to capture nonlinear relationships between variables.

    - **Variance Thresholding:**

        The VarianceThreshold transformer is applied to remove features with low variance. Features with near-constant values contribute little to the model's predictive power and can be safely discarded. The threshold parameter (0.1 in this case) determines the minimum variance required for a feature to be retained.

- **Feature Scaling:**

    Feature scaling is performed to ensure that all features have a similar scale, preventing certain features from dominating the model training process. The StandardScaler from scikit-learn is used to standardize the features by subtracting the mean and dividing by the standard deviation.

- **Train-Test Split:** The dataset was split into training and testing sets to train and evaluate the models, respectively.

## 3. Model Training:

- **Baseline Model (Dummy Classifier):** A Dummy Classifier was trained using the most frequent strategy.

- **Advanced Models:** Gradient Boosting and SVM were selected as advanced models due to their effectiveness in handling complex datasets.

- **Hyperparameter Tuning:** GridSearchCV was employed to find the optimal hyperparameters for each model.

## 4. Model Evaluation:

- **Performance Metrics:** Accuracy, precision, recall, and F1-score were computed to assess model performance.

- **Confusion Matrices:** Confusion matrices were visualized to gain insights into model behavior, particularly in terms of true positive, true negative, false positive, and false negative predictions.

## 5. Results and Analysis:

- **Dummy Classifier:** The Dummy Classifier served as a baseline with an accuracy of 0.43615384615384617 However, its performance was limited due to its simplistic nature.

- **Gradient Boosting and SVM:** Both Gradient Boosting and SVM outperformed the Dummy Classifier, Gradient Boosting achieving accuracies of 0.6074635374250388 and SVM 0.5451221588805805, respectively. They demonstrated improved precision, recall, and F1-score compared to the baseline.

- **Voting Classifiers:** Soft and hard voting classifiers were constructed using the best performing models (Gradient Boosting and SVM). These ensemble methods further improved classification performance.

## 6. Recommendations:

- **Feature Engineering:** Experiment with additional feature engineering techniques to potentially enhance model performance.

- **Model Selection:** Explore alternative classification algorithms and ensemble methods to identify the most suitable model for the dataset.

- **Further Analysis:** Investigate misclassified instances to identify patterns and potential areas for improvement.

- **Deployment:** Deploy the selected model in a production environment for real-world application.

## 7. Conclusion:

In conclusion, this study demonstrates the effectiveness of advanced classification models in accurately predicting the wine dataset. By leveraging techniques such as hyperparameter tuning and ensemble methods, significant improvements in model performance were achieved compared to the baseline. Moving forward, continuous refinement and optimization of the chosen model are recommended to ensure its robustness in real-world scenarios.